(72) Inventors: SCHIEBINGER, Geoffrey; c/o 415 Main
Street, Cambridge, Massachusetts 02142 (US). SHU, Jian;
c/o 415 Main Street, Cambridge, Massachusetts 02142
(US). TABAKA, Marcin; c/o 415 Main Street, Cambridge,
Massachusetts 02 142 (US). CLEARY, Brian; c/o 77 Mass¬
achusetts Avenue, Cambridge, Massachusetts 02139 (US).
REGEV, Aviv; c/o 415 Main Street, Cambridge, Massa¬
chusetts 02142 (US). LANDER, Eric S.; c/o 415 Main
Street, Cambridge, Massachusetts 02142 (US).

(54) Title: METHODS AND SYSTEMS FOR RECONSTRUCTION OF DEVELOPMENTAL LANDSCAPES BY OPTIMAL
TRANSPORT ANALYSIS

(57) Abstract: Methods and compositions for producing induced pluripo-
tent stem cell by introducing nucleic acids encoding one or more transcrip¬
tion factors including Obox6 into a target cell.

FTG. 1

GM, KE, LR, LS, MW, MZ, NA, RW, SD, SL, ST, SZ, TZ, UG, ZM, ZW), Eurasian (AM, AZ, BY, KG, KZ, RU, TJ, TM), European (AL, AT, BE, BG, CH, CY, CZ, DE, DK, EE, ES, FI, FR, GB, GR, HR, HU, IE, IS, IT, LT, LU, LV, MC, MK, MT, NL, NO, PL, PT, RO, RS, SE, SI, SK, SM, TR), OAPI (BF, BJ, CF, CG, CI, CM, GA, GN, GQ, GW, KM, ML, MR, NE, SN, TD, TG).

**Published:**
— *with international search report (Art. 21(3))*
— *before the expiration of the time limit for amending the claims and to be republished in the event of receipt of amendments (Rule 48.2(h))*

# METHODS AND SYSTEMS FOR RECONSTRUCTION OF DEVELOPMENTAL LANDSCAPES BY OPTIMAL TRANSPORT ANALYSIS

## CROSS-REFERENCE TO RELATED APPLICATIONS

[0001]    This application claims the benefit of U.S. Provisional Application Nos. 62/560,674, filed September 19, 2017 and 62/561,047, filed September 20, 2017. The entire contents of the above-identified applications are hereby fully incorporated herein by reference.

## TECHNICAL FIELD

[0002]    The subject matter disclosed herein is generally directed to methods and systems for analyzing the fates and origins of cells along developmental trajectories using optimal transport analysis of single-cell RNA-seq information over a given time course.

## BACKGROUND

[0003]    In the mid-20th century, Waddington introduced two images to describe cellular differentiation during development: first, trains moving along branching railroad tracks and, later, marbles following probabilistic trajectories as they roll through a developmental landscape of ridges and valleys (1, 2). These metaphors have powerfully shaped biological thinking in the ensuing decades. The recent advent of massively parallel single-cell RNA sequencing (scRNA-Seq) (3-7) now offers the prospect of empirically reconstructing and studying the actual "landscapes", "fates" and "trajectories" associated with complex processes of cellular differentiation and de-differentiation—such as organismal development, long-term physiological responses, and induced reprogramming—based on snapshots of expression profiles from heterogeneous cell populations undergoing dynamic transitions (6-1 1).

[0004]    To understand such processes in detail, general approaches are needed to answer key questions. For any given system, we would like to know: What classes of cells are present at each stage? For the cells in each class, what was their origin at earlier stages, what are their potential fates at later stages, and what is the actual outcome of a given cell? To what extent are events along a path synchronous or asynchronous? What are the genetic regulatory programs that control each path? What are the intercellular interactions between classes of cells? Answering these questions would provide insights into the nature of developmental processes: How

deterministic or stochastic is the process—that is: if, and how early, does it become determined that a particular cell or an entire cell class is destined to a specific fate? For a given origin and target fate, is there only a single path to the target, or are there multiple developmental paths? To what extent is the process cell-intrinsic, driven by intracellular mechanisms that do not require ongoing external inputs, or externally regulated, being affected by other contemporaneous cells? For artificial processes such as induced reprogramming, there are additional questions: What off-target cell classes arise? To what extent do cells activate normal developmental programs vs. unnatural hybrid programs? How can the efficiency of reprogramming be improved?

[0005]     Experimental approaches to such questions have typically involved studying bulk populations or identifying subsets of cells based on activation of one or a few genes at a specific time (e.g., reporter genes or cell-surface markers) and tracing their subsequent fate. These experiments are severely limited, however, by the need to choose subsets of cells a priori and develop distinct reagents to study each subset. For example, studies of cellular reprogramming from fibroblasts to induced pluripotent cells (iPSCs) have largely relied on RNA- and chromatin-profiling studies of bulk cell populations, together with fate-tracing of cells based on a limited set of markers (e.g., Thyl and CD44 as markers of the fibroblast state, and ICAM1, Oct4, and Nanog as markers of partial reprogramming) (12-16).

[0006]     Computational approaches based on single-cell gene expression profiles offer a complementary approach with broader molecular scope, because one can readily define classes of cells based on any expression profile at any stage. The remaining challenge is to reliably infer their trajectories across stages.

[0007]     Several pioneering papers have introduced methods to infer cellular trajectories (9, 10, 17-29). Early studies recognized that cellular profiles from heterogeneous populations can provide information about the temporal order of asynchronous processes—enabling intermediate transitional cells to be ordered in "pseudotime" along "trajectories", based on their state of cell differentiation (18). Some approaches relied on k-nearest neighbor graphs (18) or binary trees (9). More recently, diffusion maps have been used to order cell state transitions. In this case, single-cell profiles are assigned to densely populated paths through diffusion map space (20, 21). Each such path is interpreted as a transition between cellular fates, with trajectories determined by curve fitting, and cells "pseudotemporally ordered" based on the diffusion distance to the

endpoints of each path. Whereas initial efforts focused mostly on single paths, more recent work has grappled with challenges of branching, which is critical for understanding developmental decisions (10, 11, 21).

[0008]    While these pioneering approaches have shed important light on various biological systems, many important challenges remain. First, because many methods were initially designed to extract information about stationary processes (such as the cell cycle or adult stem cell differentiation) in which all stages exist simultaneously, they neither directly model nor explicitly leverage the temporal information in a developmental time course *(29)*. Second, a single cell can undergo multiple temporal processes at once. These processes can dramatically impact the performance of these models, with a notable example being the impact of cell proliferation and death *(29)*. Third, many of the methods impose strong structural constraints on the model, such as one-dimensional trajectories and zero-dimensional branch points. This is of particular concern if development follows the flexible "marble" rather than the regimented "tracks" models, in Waddington's frameworks.

## SUMMARY

[0009]    In one aspect, the present disclosure includes a method of producing induced pluripotent stem cell comprising introducing a nucleic acid encoding Obox6 into a target cell to produce an induced pluripotent stem cell. In some embodiments, the methods further comprises introducing into the target cell at least one nucleic acid encoding a reprogramming factor selected from the group consisting of: Gdf9, Oct3/4, Sox2, Soxl, Sox3, Soxl5, Soxl7, Klf4, Klf2, c-Myc, N-Myc, L-Myc, Nanog, Lin28, Fbxl5, ERas, ECAT15-2, Tell, beta-catenin, Lin28b, Sall l, Sall4, Esrrb, Nr5a2, Tbx3, and Glisl. In some embodiments, the method further comprises introducing into the target cell at least one nucleic acid encoding a reprogramming factor selected from the group consisting of: Oct4, Klf4, Sox2 and Myc. In some embodiments, the nucleic acid encoding Obox6 is provided in a recombinant vector. In some embodiments, the vector is a lentivirus vector. In some embodiments, the nucleic acid encoding the reprogramming factor is provided in a recombinant vector. In some embodiments, the method further comprises a step of culturing the cells in reprogramming medium. In some embodiments, the method further comprises a step of culturing the cells in the presence of serum. In some embodiments,

the method further comprises a step of culturing the cells in the absence of serum. In some embodiments, the induced pluripotent stem cell expresses at least one of a surface marker selected from the group consisting of: Oct4, SOX2, KLf4, c-MYC, LIN28, Nanog, Glisl , TRA-160/TRA-1-81/TRA-2-54, SSEA1, SSEA4, Sal4, and Esrbbl. In some embodiments, the target cell is a mammalian cell. In some embodiments, the target cell is a human cell or a murine cell. In some embodiments, the target cell is a mouse embryonic fibroblast. In some embodiments, the target cell is selected from the group consisting of: fibroblasts, B cells, T cells, dendritic cells, keratinocytes, adipose cells, epithelial cells, epidermal cells, chondrocytes, cumulus cells, neural cells, glial cells, astrocytes, cardiac cells, esophageal cells, muscle cells, melanocytes, hematopoietic cells, pancreatic cells, hepatocytes, macrophages, monocytes, mononuclear cells, and gastric cells, including gastric epithelial cells.

[0010]     In another aspect, the present disclosure includes a method of producing an induced pluripotent stem cell comprising introducing at least one of Obox6, Spic, Zfp42, Sox2, Mybl2, Msc, Nanog, Hesxl and Esrrb into a target cell to produce an induced pluripotent stem cell.

[0011]     In another aspect, the present disclosure includes a method of producing an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6,  into a target cell to produce an induced pluripotent stem cell.

[0012]     In another aspect, the present disclosure includes a method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

[0013]     In another aspect, the present disclosure includes a method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell to produce an induced pluripotent stem cell.

[0014]     In another aspect, the present disclosure includes an isolated induced pluripotential stem cell produced by the methods disclosed herein.

[0015]     In another aspect, the present disclosure includes a method of treating a subject with a disease comprising administering to the subject a cell produced by differentiation of the induced pluripotent stem cell produced by the methods disclosed herein.

[0016]    In another aspect, the present disclosure includes a composition for producing an induced pluripotent stem cell comprising Obox6 in combination with reprogramming medium.

[0017]    In another aspect, the present disclosure includes a composition for producing an induced pluripotent stem cell comprising one or more of the factors identified in or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6 in combination with reprogramming medium.

[0018]    In another aspect, the present disclosure includes use of Obox6 for production of an induced pluripotent stem cell.

[0019]    In another aspect, the present disclosure includes use of a factor identified in or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6 for production of an induced pluripotent stem cell.

[0020]    In another aspect, the present disclosure includes a method of increasing the efficiency of reprogramming a cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

[0021]    In another aspect, the present disclosure includes a method of increasing the efficiency of reprogramming a cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5 and Table 6, into a target cell to produce an induced pluripotent stem cell.

[0022]    In another aspect, the present disclosure includes a computer-implemented method for mapping developmental trajectories of cells, comprising: generating, using one or more computing devices, optimal transport maps for a set of cells from single cell sequencing data obtained over a defined time course; determining, using one or more computing devices, cell regulatory models, and optionally identifying local biomarker enrichment, based on at least the generated optimal transport maps; defining, using the one or more computing devices, gene modules; and generating, using the one or more computing devices, a visualization of a developmental landscape of the set of cells.

[0023]    In some embodiments, determining cell regulatory models comprise sampling pairs of cells at a first time and a second time point according to transport probabilities. In some embodiments, the method further comprises using the expression levels of transcription factors at the earlier time point to predict non-transcription factor expression at the second time point. In

some embodiments, identifying local biomarker enrichment comprises identifying transcription factors enriched in cells having a defined percentage of descendants in a target cell population. In some embodiments, the defined percentage is at least 50% of mass. In some embodiments, defining gene modules comprises partitioning genes based on correlated gene expression across cells and clusters. In some embodiments, partitioning comprises partitioning cells based on graph clustering. In some embodiments, graph clustering further comprises dimensionality reduction using diffusion maps. In some embodiments, the visualization of the developmental landscape comprises high-dimensional gene expression data in two dimensions. In some embodiments, the visualization is generated using force-directed layout embedding (FLE). In some embodiments, the visualization provides one or more cell types, cell ancestors, cell descendants, cell trajectories, gene modules, and cell clusters from the single cell sequencing data.

[0024]    In another aspect, the present disclosure includes a computer program product, comprising: a non-transitory computer-executable storage device having computer-readable program instructions embodied thereon that when executed by a computer cause the computer to execute the methods disclosed herein.

[0025]    In another aspect, the present disclosure includes a system comprising: a storage device; and a processor communicatively coupled to the storage device, wherein the processor executes application code instructions that are stored in the storage device and that cause the system to executed the methods disclosed herein.

[0026]    In another aspect, the present disclosure includes a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Gdf9 into a target cell to produce an induced pluripotent stem cell.

[0027]    These and other aspects, objects, features, and advantages of the example embodiments will become apparent to those having ordinary skill in the art upon consideration of the following detailed description of illustrated example embodiments.


BRIEF DESCRIPTION OF THE DRAWINGS

[0028]    An understanding of the features and advantages of the present invention will be obtained by reference to the following detailed description that sets forth illustrative

embodiments, in which the principles of the invention may be utilized, and the accompanying drawings of which:

[0029]     **FIG. 1** - is a block diagram depicting a system for mapping developmental trajectories of cells, in accordance with certain example embodiments

[0030]     **FIG. 2** - is a block flow diagram depicting a method for mapping development trajectories of cells, in accordance with certain example embodiments.

[0031]     **FIG. 3** - is a diagram showing data $S_i$ from a generic branching developmental process. The x-axis represents the time and the y-axis represents expression.

[0032]     **FIG. 4** - provides a schematic of a regulatory vector file which gives rise to a time-dependent probability distribution.

[0033]     **FIGs. 5A-5G** - **(FIGs. 5A-5B)** Waddington's classical analogies of cells undergoing differentiation, initially (1936) illustrated by railroad cars on switching tracks **(FIG. 5A)** and later (1957) by marbles rolling in a landscape **(FIG. 5B),** with trajectories shaped by hills and valleys. **(FIGs. 5C-E)** Differentiation processes in which the ultimate fate of individual cells (filled dots) is **(C)** predetermined **(FIG. 5D)** not predetermined, or **(FIG. 5E)** progressively determined. Arrows indicate possible transitions, and color represents cell fate, with red and blue indicating distinct fates, light red and light blue indicating partially determined fates, and grey indicating undetermined fate. **(FIG. 5F)** Illustration of transported mass. A transport map, , describes how a point x at one stage (X) is redistributed across all points (denoted by "") at the subsequent stage (Y). **(FIG. 5G)** Transport maps computed from a time series of samples taken from a time-varying distribution. Between each pair of time points, a transport map redistributes the cells observed at time to match the distribution of cells observed at time.

[0034]     **FIGs. 6A-6C** - **(FIG. 6A)** Representation of reprogramming procedure and time points of sample collection. (Top) Mouse embryos (E13.5) were dissected to obtain secondary MEFs (2° MEF), which were reprogrammed into iPSCs. In Phase-1 of reprogramming (light blue; days 0-8), doxycycline (Dox) was added to the media to induce ectopic expression of reprogramming factors *(Oct4, Kl/4, Sox2,* and *Myc).* In Phase-2 (days 9-16), Dox was withdrawn from the media, and cells were grown either in the presence of 2i (light red) or serum (light green). Samples were also collected from established iPSC lines reprogrammed from the same 2° MEFs, maintained in either 2i or serum conditions (far right in each time course). Individual dots

along the time course indicate time points of scRNA-Seq collection, with two dots indicating biological replicates. **(FIG. 6B)** Number of scRNA-Seq profiles from each sample collection that passed quality control filters. **(FIG. 6C)** Bright field images of day 0 (Phase 1-(Dox)) and day 16 cells during reprogramming in (Phase-2(2i)) and (Phase-2(serum)) culture conditions.

**[0035]** **FIGs. 7A-7F** - scRNA-Seq profiles of all 65,781 cells were embedded in two-dimensional space using FLE, and annotated with indicated features. **(FIG. 7A)** Unannotated layout of all cells. Each dot represents one cell. **(FIGs. 7B-7C)** Annotation by time point (color) and biological feature, with Phase-2 points from either **(FIG. 7B)** 2i condition or **(FIG. 7C)** serum condition. Phase-1 points appear in both **(FIG. 7B)** and **(FIG. 7C).** Individual cells are colored by day of collection, with grey points (BC, background color) representing Phase-2 cells from serum (in **FIG. 7B)** or 2i (in **FIG. 7C). (FIG. 7D)** Annotation by cell cluster. Cells were clustered on the basis of similarity in gene expression. Each cell is colored by cluster membership (with clusters numbered 1-33). **(FIGs. 7E-7F)** Annotation by gene signature **(FIG. 7E)** and individual gene expression levels **(FIG. 7F).** Individual cells are colored by gene signature scores (in **FIG. 7E)** or normalized expression levels (in **FIG. 7F;** , where E is the number of transcripts of a gene per 10,000 total transcripts).

**[0036]** **FIGs. 8A-8F** - **(FIG. 8A)** Schematic representation of the major cluster-to-cluster transitions (see Table 10 for details[BC17] ). Individual arrows indicate transport from ancestral clusters to descendant clusters, with colors corresponding to the ancestral cluster. For each descendant cluster, arrows were drawn when at least 20% of the ancestral cells (at the previous time point) were contained within a given cluster (self-loops not shown). Arrow thickness indicates the proportion of ancestors arising from a given cluster. **(FIG. 8B)** Heatmap depiction of cluster descendants in 2i condition. In each row of the heatmap, color intensity indicates the number of descendant cells ("mass", normalized to a starting population of 100 cells) transported to each cluster at the subsequent time point (see Table 10 for details). Clusters with highly-proliferative cells *(e.g.,* cluster 4) transport more total mass than clusters with lowly-proliferative cells *(e.g.,* cluster 14). ((**FIG. 8C)** Depiction of divergent day 8 descendant distributions for two clusters of cells at day 2 (cluster 4 (left) and cluster 6 (right). Color intensity indicates the distribution of descendants at day 8, with bright teal indicating high probability fates and gray indicating low probability fates. **(FIG. 8D)** Enrichment of the ancestral distributions of iPSCs,

Valley of Stress, and alternative fates (neuron-like and placenta-like) in clusters of day 2 cells. The red horizontal dashed line indicates a null-enrichment, where a cluster contributes to the ancestral distribution in proportion to its size. Cluster 4 has a net positive enrichment because its descendants are highly proliferative, while cluster 6 has a net negative enrichment because its descendants are lowly proliferative. **(FIG. 8E)** and **(FIG. 8F)** Ancestral trajectories of indicated populations of cells at day 16 (iPSCs, placental, neural-like cells, *etc)* in serum **(FIG. 8E)** and 2i **(FIG. 8F).** Clusters used to define the indicated populations are shown in parentheses. Colors indicate time point. Sizes of points and intensity of colors indicate ancestral distribution probabilities by day (color bars, right; BC, background color, representing cells from the other culture condition).

[0037]    **FIGs. 9A-9D - (FIG. 9A)** Classification of genes into 14 groups based on similar temporal expression profiles along the trajectory to successful reprogramming. Averaged gene expression profiles for each group, in 2i and serum conditions (left). Heatmap for genes within each group, with intensity of color indicating log2-fold change in expression relative to day 0 (middle). Representative genes and top terms from gene-set enrichment analysis for each group (right). **(FIG. 9B)** Comparison of FACS and *in* silico sorting experiments. Scatterplot shows reprogramming efficiencies determined by FACS sort and growth experiments (blue triangles) *(16)* and our computationally inferred trajectories (red squares). The specific cell surface markers used for the *in silico* and experimental methods are indicated. Reprogramming efficiencies for these categories (calculated both experimentally and *in silico)* are normalized to the percentage of EGFP+ colonies in CD44$^-$ ICAM1$^+$ Nanog$^+$ condition (details found in *Appendix 5).* **(FIG. 9C)** Schematic of regulatory model in which TF expression in ancestral cells is predictive of gene expression in descendant cells. **(FIG. 9D)** Onset of iPSC-associated TFs in 2i (left) and serum (right). (Top) Mean expression levels weighted by iPSC ancestral distribution probabilities (Y axis) of *Nanog, Obox6,* and *Sox2* at each day (X axis). (Bottom) Normalized expression of TF modules "A" and "B" from our regulatory model (as in **FIG. 9B)** that were associated with gene expression in iPSCs.

[0038]    **FIGs. lOA-lOC - (FIGs. 10A-10B)** Bright field and fluorescence images of iPSC colonies generated by lentiviral overexpression of *Oct4, Kl/4, Sox2,* and *Myc* (OKSM) with either an empty control, *Zfp 42* or *Obox6* expression cassette, in either Phase-l(Dox)/Phase-2(2i)

**(FIG. 10A)** and Phase- l(Dox)/Phase-2(serum) **(FIG. 10B)** conditions (indicated). Cells were imaged at day 16 to measure Oct4-EGFP$^{+}$ cells. Bar plots representing average percentage of Oct4-EGFP$^{+}$ colonies in each condition on day 16 are included below the images. Shown are data from one of five independent experiments, with three biological replicates each. Error bars represent standard deviation for the three biological replicates. **(FIG. IOC)** Schematic of the overall reprogramming landscape highlighting: the progression of the successful reprogramming trajectory, alternative cell lineages, and specific transition states (Horn of Transformation). Also highlighted are transcription factors (orange) predicted to play a role in the induction and maintenance of indicated cellular states, and putative cell-cell interactions between contemporaneous cells in the reprogramming system.

[0039]    **FIGs. 11A-11D** - Single-cell RNA-Seq quality metrics. **(FIG. 11A)** Correlation between number of genes and tran- scripts per cell (loglO transformed). Cells with fewer than 1000 genes detected were filtered out. The color gradient represents cell density. **(FIG. 11B)** Variation in single cell data depicted by correlation between transcript levels (loglO transformed average transcript counts) detected in biological replicates generated from day 10 samples in 2i conditions. Pearson correlation coefficient (r) is given. The color gradient represents cell density. **(FIG. 11C)** Biological variation in single cell data depicted by correlation between tran- script levels (loglO transformed average transcript counts) detected in iPSCs and MEFs. Pearson corre- lation coefficient (r) is given. The color gradient represents cell density. **(FIG. 11D)** Correlogram visualizing correlation between single cell gene expression profiles between various time points and their biological replicates. In this plot, the correlation coefficients (circles) are colored according to their values, ranging from 0.75 (blue) to 1 (red). The size of the circles represents the magnitude of the coefficient. The replicates within the timepoints are denoted with suffixes 1 and 2.

[0040]    **FIGs. 12A-12C** - Comparison of various dimensionality reduction methods to visualize single cell RNA- Seq data. High-dimensional structure of single-cell expression data was embedded in low-dimensional space for visualization using **(FIG. 12A)** the Force-directed Layout Embedding algorithm (FLE) (directed graph approach) and the t-Distributed Stochastic Neighbor Embedding algorithm (t-SNE) with **(FIG. 12B)** principal components and **(FIG. 12C)** diffusion maps as input parameters.

**[0041]** **FIG. 13** - Visualization of gene modules across reprogramming time points. Expression profiles of all 65,781 cells studied were embedded in two-dimensional space, using force-directed layout embed- ding (FLE). The layouts were annotated by single-cell z-scores for 44 gene modules (details in *Table 1)*. The color gradient represents the distribution of z-scores across all cells for a given gene module.

**[0042]** **FIGs. 14A-14B** - Characterization of cell clusters. **(FIG. 14A)** Heatmap representing the enrichment of cells from the indicated samples at various time points and culture conditions across 33 different clusters. The color gradient represents the range of cell fractions from 0-0.25. **(FIG. 14B)** Heatmap depicting the enrichment of correlated gene modules within specific cell clusters. The color gradient represents the average gene module scores at the indicated cell clusters. Specific cell clusters that show highly correlated gene module scores were numerically labeled as shown

**[0043]** **FIG. 15** - Visualization of individual gene expression levels.Normalized expression levels [log2(E+l)] for indicated genes were used to annotate force-directed layout embedding (FLE) graphs generated from the expression profiles of 65,781 cells. E represents the number of transcripts of a gene per 10,000 total transcripts

**[0044]** **FIGs. 16A-16E** - Distribution of gene signatures. **(FIG. 16A)** Distribution of proliferation scores for cells at day 0 (solid black). Proliferation scores were calculated from combined expression levels of Gl/S and G2/M cell cycle genes (see *Appendix 5)*. Normal mixture modeling (dashed line) was used to classify the cells based on proliferation scores into non-cycling (red) and cycling (blue) cells (top). Visualization of the cycling and non-cycling of cells on FLE at day 0 (bottom). **(FIG. 16B)** Violin plots of single-cell scores for indicated gene signatures and Shisa8 expression levels in clusters 3, 4, 5, and 6. **(FIG. 16C)** Violin plots of single cell scores for indicated gene signatures in clusters 7, 8, and 18. **(FIG. 16D)** Bar plots of normalized expression levels [log2(E+l)] for indicated genes, where E is the number of transcripts of a gene per 10,000 total transcripts. **(FIG. 16E)** Single-cell scores for indicated gene signatures across all 33 cell clusters.

**[0045]** **FIGs. 17A-17C** - Heatmap depiction of origins and fates of cells inferred from optimal transport. Heatmap depiction of cluster descendants in **(FIG. 17A)** serum condition, and cluster ancestors in **(FIG. 17B)** 2i and **(FIG. 17C)** serum conditions. Each row of the heatmap in

11

**(FIG. 17A)** shows how the descendants of the cells in a particular cluster are distributed over all clusters. Color intensity indicates the number of descendant cells ("mass", normalized to a starting population of 100 cells) transported to each cluster at the next time point. Each column of the heatmaps in **(FIG. 17B, FIG. 17C)** shows how the ancestors of a particular cluster are distributed over all clusters. *Table 10* contains the specific numerical values.

[0046]     **FIGs. 18A-18F** - Potential cell-cell interactions across the reprogramming time course. **(FIG. 18A)** Temporal pattern of the net potential for paracrine signaling between contemporaneous cells. Each dot represents the aggregated interaction score across all ligand-receptor pairs for a given combination of clusters (all 149 detected ligands). The aggregate interaction score is defined as a sum of individual interaction scores. **(FIG. 18B)** As in A, but genes specific to SASP signature are considered (20 detected ligands). **(FIG. 18C)** Heatmap representing the aggregate interaction scores on day 16 cells in 2i condition for ligands specific to SASP signature. Rows correspond to clusters of cells expressing ligands. Columns correspond to clusters of cells expressing cognate receptors. Only clusters containing more than 1% of cells from day 16 (2i) are shown. **(FIGs. 18D-18F)** Potential ligand-receptor pairs ranked by their standardized interaction scores calculated from the permuted data (see *Appendix 5* for details). Ligand-receptor pairs between **(FIG. 18D)** valley of stress cells (clusters 11-17) and iPSCs (clusters 28-33) on day 16 (2i), **(FIG. 18E)** valley of stress cells and preneural/neural-like cells (clusters 23, 26, and 27) on day 16 (serum), and **(FIG. 18F)** placental-like cells (clusters 24 and 25) and valley of stress cells on day 12 (2i)

[0047]     **FIGs. 19A-19F** - Gene modules and associated transcription factors based on optimal transport. Using optimal transport trajectories, TF levels in cells at time t are used to predict the activity levels of gene modules in descendant cells at time t + 1. Gene modules are learned during model training to capture coherent expression programs. For five modules **(FIGs. 19A-19E),** bar plots depict the top 50 genes in the module (black), and the top 20 TFs each associated with positive (red) and negative (blue) module activity. **(FIGs. 19A- 19B)** Two modules that are active in cells with placental identity. **(FIG. 19C)** A module active in cells with neural identity. **(FIG. 19D-19E)** Two modules active in successfully reprogrammed cells. **(FIG. 19F)** Enrichment analysis of TFs in day 12 cells with high (>80%) vs. low (<20%) probability of successful reprogramming. Dot size and color represent percentage of day 12 cells expressing the

indicated TF in high- or low-probability cells. Bar heights indicate the fold enrichment in high-vs. low-probability cells.

[0048]    **FIGs. 20A-20C** - Effect of overexpression of *Obox6* and *Zpf42* on reprogramming efficiency. **(FIG. 20A)** Percentage of Oct4-EGFP+ cells at day 16 of reprogramming from secondary MEFs by lentiviral overexpression of *Oct4, Kl/4, Sox2,* and *Myc* (OKSM) combined with either *Zfp42*, *Obox6*, or an empty control, in either 2i or serum conditions. Oct4-EGFP+ cells were measured by flow cytometry. Plot includes the percentage of Oct4-EGFP+ cells in three biological replicates (for *Zfp42* and *Obox6* overexpression, or an empty control) from five independent experiments (Exp). **(FIG. 20B, FIG. 20C)** Number of Oct4-EGFP+ colonies at day 16 of reprogramming from primary MEFs by lentiviral overexpression of individual *Oct4, Kl/4, Sox2,* and *Myc* combined with either *Zfp42*, *Obox6*, or an empty control in **(FIG. 20B)** 2i and **(FIG. 20C)** serum conditions. Plot includes the number of Oct4-EGFP+ cells in three biological replicates (for *Zfp42* and *Obox6* overexpression, or an empty control) from two independent experiments (Exp).

[0049]    **FIGs. 21A-21E** - X-chromosome reactivation. **(FIGs. 21A**-21C) Boxplots showing X/Autosome expression ratio (left panel) and Xist expression $\log2(E+l)$ across individual cells by clusters (right panel): **(FIG. 21A)** all cells, **(FIG. 21B)** phase-l(Dox) and phase-2(2i) cells, **(FIG. 21C)** phase-l(Dox) and phase-2(serum) cells. **(FIGs. 21D-21F)** - X/Autosome expression ratio and A6, A7 activation pattern changes along the successful trajectory determined by optimal transport: Relative gene expression changes of individual genes from A6 **(FIG. 21D)** and A7 **(FIG. 21E)** activation patterns (gray solid lines). Black and blue solid lines correspond to average relative expression of genes and average X/Autosome expression ratios, respectively. **(FIG. 21F)** Comparison between activation of A6 and A7 programs (average relative expression) with X/Autosome expression ratio. Distribution of X/Autosome expression ratios **(FIG. 21G)** and A7 scores **(FIG. 21H)** across all cells. Dotted lines represent threshold values used in classification of cells that reactivated X-chromosome ($> 1.4$) and upregulated A7 genes ($> 0.25$).

[0050]    **FIGs. 22A-22C** - Single-cell expression levels were used to identify cells with aberrant expression in large chromosomal regions. **(FIG. 22A)** Whole chromosome aberrations were detected in 1% of all cells. Each dot represents one chromosome (X axis) in a single cell

with significant aberrations (FDR 10%), with violin plots capturing the distributions of dots. The net expression of these chromosomes relative to the average expression across all cells (Y axis) is 1.7-fold higher (median, left panel) and 2.2-fold lower (right panel), indicating whole chromosome gain and loss, respectively. The median relative expression levels are slightly higher (lower) than the 1.5-fold (2-fold) increase (decrease) that would be expected from a true chromosomal gain (loss) because our statistics are conservative in calling significant events but allow for a long tail of high (low) expression. **(FIG. 22B)** Visualization of cells with significant subchromosomal aberrations (red) in FLE. **(FIG. 22C)** Bar plots depict the fraction of cells in each cluster with significant subchromosomal (25-200Mbp) aberrations (FDR 10%).

[0051]     **FIGs. 23A-23F** - Modeling developmental processes with optimal transport. Waddington-OT: a probabilistic model for developmental processes. **(FIG. 23A)** A temporal progression of a time-varying distribution $\mathbf{P}_t$ (left) can be sampled to obtain finite empirical distributions of cells $\widehat{\mathbf{P}}_{t_i}$ at various time points $t_1, t_2, t_3$ (right). Over short time scales, the unknown true coupling, $Y_{t1,t_2}$, is assumed to be close to the optimal transport coupling, $\pi_{t_1,t_2}$, which can be approximated by $\pi_{t1,t2}$ computed from the empirical distributions $\widehat{\mathbf{P}}_{t_1}$ and $\widehat{\mathbf{P}}_{t_2}$. **(FIGs. 23B-23F)** Simulated data and analysis performed by Waddington-OT. **(FIG. 23B)** Single-cell profiles (individual dots) are embedded in two dimensions and colored by the time of collection. Optimal transport can be used to calculate the descendant trajectories **(FIG. 23C)** and ancestor trajectories **(FIG. 23D)** of any subpopulation of interest (cells highlighted in black; color indicates time). Ancestor distributions of distinct subpopulations can be compared to calculate their shared ancestry **(FIG. 23E)** (ancestors of each population shown in red and blue, shared ancestors in purple). **(FIG. 23F)** The expression of gene signatures (left; green, high expression; grey, low expression) can be predicted from the earlier expression of transcription factors (middle; black, high expression; grey, low expression) in a gene regulatory model by analyzing trends along ancestor trajectories. In the plot at right, at each time point, the height of the curve depicts the average expression in the ancestors of cells in the leftmost tip.

[0052]     **FIGs. 24A-24H** - A single cell RNA-Seq time course of iPSC reprogramming. **(FIG. 24A)** Representation of reprogramming procedure and time points of sample collection. (Top) Mouse embryos (E13.5) were dissected to obtain secondary MEFs (2° MEF), which were reprogrammed into iPSCs. In Phase-1 of reprogramming (light blue; days 0-8), doxycycline

(Dox) was added to the media to induce ectopic expression of reprogramming factors *(Oct4, Kl/4, Sox2,* and *Myc).* In Phase-2 (days 9-18), Dox was withdrawn from the media, and cells were grown either in the presence of 2i (light red) or serum (light green). Samples were also collected from established iPSC lines reprogrammed from the same 2° MEFs, maintained in either 2i or serum conditions (far right in each time course). Individual dots indicate time points of scRNA-Seq collection. **(FIGs. 24B-24E)** scRNA-Seq profiles of all 251,203 cells (individual dots) were embedded in two-dimensional space using FLE, and annotated with indicated features. **(FIG. 24B)** Unannotated layout of all cells, with the density of cells in each region indicated by intensity. **(FIG. 24C)** Cells colored by time point, with Phase-2 points from either 2i condition (left) or serum condition (right). Phase-1 points appear in both subplots. Grey points represent Phase-2 cells from the other condition. **(FIG. 24D)** In different regions of the FLE, cells have distinct expression patterns of six major gene signatures (average expression z-score of genes in a signature indicated by red color bar). Gene signature activity and trajectory analysis were used to define the major cell sets **(FIG. 24E)** and to establish the overall flow through the landscape **(FIG. 24F)** (schematic representation). **(FIG. 24G)** The relative abundance (y-axis) of each cell set (colored lines) is plotted over time (x-axis) in 2i (top) and serum (bottom). **(FIG. 24H)** Validation via geodesic interpolation in serum condition. Data at withheld timepoints (x-axis) are interpolated using data at the neighboring timepoints. Interpolation is done using a null estimator of independent coupling (blue) and the optimal transport coupling (red), with the distance between interpolated and withheld data indicated on the y-axis. The distance between two batches of withheld data at the same point is shown in green. Shaded regions indicate standard deviations over independent samples of the coupling map.

[0053]    **FIGs. 25A-25H** - In initial stages of reprogramming, cells progress toward stromal or MET fates. **(FIG. 25A)** Cells in the stromal region have higher expression of gene signatures (red color bar, average z-score) and individual genes (red color bar, log(TPM+l)) that are associated with stromal activity and senescence. Ancestors of day 18 stromal cells are visualized on the FLE **(FIG. 25B)** (colored by day, intensity indicates probability), and expression trends along this ancestor trajectory **(FIG. 25C)** are depicted for gene signatures (left) and individual transcription factors (TFs; right). The ancestors of day 8 MET cells **(FIG. 25D)** have a distinct trajectory and gene signature trends **(FIG. 25E),** and show differential expression of several TFs

**(FIG. 25F)** (dashed line, average TPM in stromal ancestors; solid line, average TPM in MET ancestors). **(FIG. 25G, FIG. 25H)** The MET and stromal fates are gradually specified from day 0 through 8. Color bar in **(FIG. 25G)** indicates log-likelihood of obtaining stromal vs. MET fate. **(FIG. 25H)** The extent to which the stromal ancestor distribution has diverged (y-axis) from all other fates at each point in time (x-axis). The divergence is quantified as ½ times the total variation distance between the ancestor distributions.

[0054]     **FIGs. 26A-26F**   - iPSCs emerge from cells in the MET Region. **(FIG. 26A)** Ancestors of day 18 iPSCs in 2i (left) and serum (right) are visualized on the FLE (colored by day, intensity indicates probability). Cells in the iPSC region express pluripotency marker genes **(FIG. 26B)** (red color bar, log(TPM+l))  and diverge from alternative fates also arising from the MET region (neural, epithelial, and trophoblast) from days 8-12 **(FIG. 26C)** (divergence between pairs of lineages indicated by individual lines; green line, divergence between iPSC and all others). **(FIG. 26D)** Expression trends along the ancestor trajectory in serum are depicted for gene signatures (left) and individual transcription factors (right). **(FIG. 26E)** A signature of X reactivation (left; red color bar, average z-score) and *Xist* expression (right; log(TPM + 1)) visualized on the FLE. **(FIG. 26F)** Trends in X-inactivation, X-reactivation and pluripotency along the iPSC trajectory in 2i. The values on the axis refer to average expression across early (black) and late (red) pluripotency activation genes, *Xist* average expression (log(TPM+l), orange) and X/Autosome expression ratio (blue) along the iPSC trajectory.

[0055]     **FIGs. 27A-27G** - Extra-embryonic and neural-like cells emerge during reprogramming. Subpopulations of trophoblast- **(FIGs. 27A-27C)** and neural-like **(FIGs. 27D-27G)** cells are found in the late stages of reprogramming. Ancestors of day 18 trophoblasts are visualized on the FLE **(FIG. 27A)** (colored by day, intensity indicates probability), and expression trends along the ancestor trajectory in serum **(FIG. 27B)** are depicted for gene signatures (left) and individual transcription factors (right). **(FIG. 27C)** Cells in the trophoblast cell set were re-embedded by FLE, and scored for signatures of trophoblast progenitors (TP), spiral artery trophoblast giant cells (SpA-TGC), and spongiotrophoblasts (SpTB). Colors indicate significant expression of TP, SpA-TGC, and SpTB signatures (4ogl0( FDR q-value)), or expression of labyrinthine trophoblast marker gene *Gcml* (red color bar, log(TPM + 1)). Ancestors of day 18 cells in the neural region are visualized on the FLE **(FIG. 27D)**

16

(colored by day, intensity indicates probability), and expression trends along the ancestor trajectory in serum **(FIG. 27E)** are depicted for gene signatures (left) and individual transcription factors (right). **(FIG. 27F)** Cells with radial glial (RG) and differentiated subtype signatures begin to appear around day 12 (x-axis, time; y-axis, relative abundance in serum). **(FIG. 27G)** All cells in the neural region we re-embedded by FLE, and scored for significant expression of differentiated signatures (OPC, astrocyte, cortical neurons; color, -loglO(FDR q-value)), or annotated by expression of markers of inhibitory and excitatory neurons (red color bars, log(TPM + 1)). OPC, oligodendrocyte precursor cells.

**[0056]**     **FIGs. 28A-28K** - Paracrine signaling and genomic aberrations. **(FIG. 28A)** Schematic of the paracrine signaling interaction scores. High potential interaction occurs between two groups of contemporaneous cells in which one group secretes a ligand and a second group expresses a cognate receptor. **(FIG. 28B)** Temporal pattern of the net potential for paracrine signaling between contemporaneous cells in serum condition. Each dot represents the aggregated interaction score across all ligand-receptor pairs for a given combination of clusters (Figure S5A, all 180 detected ligands). The aggregate interaction score is defined as a sum of individual interaction scores. **(FIGs. 28C-E)** Potential ligand-receptor pairs between ancestors of stromal cells and iPSCs **(FIG. 28C),** neural-like cells **(FIG. 28D),** and trophoblasts **(FIG. 28E),** ranked by their standardized interaction scores calculated from the permuted data (see STAR Methods for details). **(FIGs. 28F-H)** Individual cells on the FLE colored by the expression level (log(TPM+l))  of ligands (upper row) and receptors (lower row) for top interacting pairs between stromal cells and iPSCs **(FIG. 28F),** neural-like cells **(FIG. 28G),** and trophoblasts **(FIG. 28H).  (FIGs. 28I-28K)** Evidence for genomic aberrations was found at the level of whole chromosomes **(I)** and sub-chromosomal regions spanning 25 housekeeping genes **(FIGs. 28J, 28K). (FIG. 281)** Average expression of housekeeping genes on chromosomes (numbered on x-axis) in single cells (dots with violin plots) with evidence of genomic amplification (left panel) or loss (right panel), relative to all cells without evidence of aberrations (y-axis, relative expression). **(FIG. 28J)** Individual cells on the FLE are colored by statistical significance (-logl0( q-value ), colorbar ) of evidence for sub-chromosomal aberrations. **(FIG. 28K)** Average expression of genes on chromosome 15 in trophoblast-like cells with evidence of a recurrent sub-

chromosomal amplification (FDR 10%, region indicated by red lines), relative to trophoblast-like cells without evidence of amplification in this region (y-axis, relative expression).

[0057]      **FIGs. 29A-29D** - *Obox6* enhances reprogramming. **(FIG. 29A)** For cells (individual dots) at each timepoint (x-axis), the log-likelihood ratio of obtaining iPSCs fate vs non iPSCs fate in 2i is depicted on the y-axis. Cells expressing *Obox6* are highlighted in red. **(FIG. 29B)** Bright field and fluorescence images of iPSC colonies generated by lentiviral overexpression of *Oct4, Klf4, Sox2,* and *Myc* (OKSM) with either an empty control, *Zfp 42* or *Obox6* expression cassette, in Phase-l(Dox)/Phase-2(2i). **(FIG. 29C)** Bar plots representing average percentage of Oct4-EGFP$^+$ colonies in 2i on day 16. Data shown is one of five independent experiments, with three biological replicates each. Error bars represent standard deviation for the three biological replicates. **(FIG. 29D)** Schematic of the overall reprogramming landscape in serum highlighting: the progression of the successful reprogramming trajectory (represented in black), alternative cell lineages and subtypes within these lineages (Stromal in blue, trophoblast-like in red, neural in green and epithelial in orange), and specific transition states (MET in purple). Also highlighted are transcription factors predicted to play a role in the transition to indicated cellular states (as indicated by the specific color), and putative cell-cell interactions between contemporaneous cells in the reprogramming system. i and e Neurons refers to inhibitory and excitatory neurons respectively.

[0058]      **FIGs. 30A-30G**  - Related to **FIGs. 24A-24H:** Validation, stability, and comparison to pilot study. **(FIGs. 30A-30C)** Unbalanced transport can be used to tune growth rates. **(FIG. 30A)** When the unbalanced regularization parameter is large (=16), growth constraints are imposed strictly, and the input growth (x-axis; determined by gene signatures — see STAR Methods) is well-correlated to the output growth (y-axis; implicit growth rate determined from the transport map). **(FIG. 30B)** When the unbalanced parameter is small (=1), the growth constraints are only loosely imposed, allowing implicit growth rates to adjust and better fit the data. **(FIG. 30C)** The correlation of output vs input growth as a function of . **(FIG. 30D)** Validation by geodesic interpolation for 2i conditions. As in **FIG. 24H** (which shows serum), the red curve shows the performance of interpolating held-out time points with optimal transport. The green curve shows the batch-to-batch Wasserstein distance for the held-out time points, which is a measure of the baseline noise level. The blue curve shows the performance of a null

model (interpolating according to the independent coupling, including growth). **(FIGs. 30E-30F)** Comparison to pilot dataset. **(FIG. 30E)** Trends in signature scores along ancestor trajectories to iPSC, Stromal, Neural, and Trophoblast cell sets. Trends for the pilot dataset are shown with open circles and trends for the large dataset are shown with solid lines. **(FIG. 30F)** Shared ancestry results for pilot dataset (solid lines) and for the larger dataset (dashed lines). **(FIG. 30G)** Bright field images of day 2 (Phase l-(Dox)), day 4 (Phase l-(dox)) and day 18 cells during reprogramming in (Phase-2(2i)) and (Phase-2(serum)) culture conditions. BF (bright field). GFP (Oct4-GFP).

**[0059]**      **FIGs. 31A-31F** - Related to **FIGs. 25A-25H** Divergence of Stromal and MET fates during the initial stages of reprogramming. **(FIGs. 31A-31B)** Cells from the stromal region were re-embedded by FLE, and scored for signatures of long-term cultured MEFs (left) or stromal cells in the embryonic mesenchyme (right) found in the Mouse Cell Atlas **(FIG. 31A),** or from signatures derived from genes co-expressed (see STAR-Methods) with *Cxcll2, Ifltml,* or *Matn4* in the stromal cell set **(FIG. 31B)** (red color bars, average z-score of expression). **(FIG. 31C)** Ectopic OKSM expression levels are predictive of MET fate. The y-axis shows correlation between OKSM expression and the log-likelihood of obtaining MET fate. Color (red vs blue) distinguishes the two batches at each time point (x-axis). **(FIG. 31D)** *Fut9+* and *Shisa8+* expression patterns visualized in a fate-divergence layout. Each dot represents a single cell, colored by expression of either *Fut9* (left) or *Shisa8* (right). The x-axis shows time of collection and the y-axis shows the log-likelihood ratio of obtaining MET vs Stromal fate, as predicted by optimal transport. **(FIG. 31E)** The Stromal region is a terminal destination as evidenced by (1) the large flow of cells into the region around day 9 (green spike, first and second panels) and (2) essentially zero flow out of the region (blue curves, first and second panels). By contrast, the MET region is a transient state as evidenced by the blue curves in the right two panels showing significant transitions out of MET. **(FIG. 31F)** Day 0 MEFs (DO; black dots) we re-embedded together with cells from the stromal set (red dots) in a TSNE plot.

**[0060]**      **FIGs. 32A-32C** - Related to **FIGs. 26A-26F:** iPSCs. **(FIG. 32A)** Cells with significant expression of 2 cell (2C), 4 cell (4C), 8 cell (8C), 16 cell (16C) and 32cell (32C) signatures at an FDR of 10% on iPSC-specific FLE. **(FIG. 32B)** Overlap between different early embryonic stages. The horizontal bars show the number of cells identified as 2C, 4C, 8C, 16C, or

32C. The vertical bars indicate the number of cells in each possible combination of these cell sets (e.g. 2C and 4C). **(FIG. 32C)** Heatmap showing trends in expression of 1479 variable genes (STAR-Methods) along the ancestor trajectory to iPSCs. Color indicates fold-change in expression relative to day 0 (white). Each row shows the mean expression trend for a single gene, where the mean is computed with respect to the ancestor distribution. Genes are clustered into groups with similar trends. Terms on the right indicate significant gene set enrichment (GSEA, all adjusted p-values < 0.01) in one of several databases (M, MSigDB; BP, GO biological process; W, WikiPathways; C, chromosome; CC, GO cellular component).

[0061]     **FIGs. 33A-33E** - Related to **FIGs. 27A-27G:** Trophoblast and Neural subtypes. **(FIG. 33A)** Expression of individual marker genes (red color bars, log(TPM +1); see also Table S2) for each subtype on the trophoblast FLE (as in Figure 5C). TP, trophoblast progenitors; SpA-TGC, spiral artery trophoblast giant cells; SpTB, spongiotrophoblasts; LaTB, labyrinthine trophoblasts. **(FIG. 33B)** Cells with a gene signature of extra-embryonic endoderm (XEN) arise in a single batch on day 15.5 (red color bar, average z-score). **(FIGs. 33C-33E)** Cells in the neural region were re-embedded by tSNE and annotated with various features. **(FIG. 33C)** Marker gene expression (red color bar, log(TPM + 1)) of neural subtypes on the neural tSNE. **(FIG. 33D)** Cells with significant expression (black dots) of indicated signatures from the Allen Mouse Brain Atlas on the neural tSNE at an FDR of 10%. OPC refers to oligodendrocyte precursor cells. **(FIG. 33E)** Cells in the neural region present from days 12.5-14.5 (left) or days 17-18 (right).

[0062]     **FIGs. 34A-34E** - Related to **FIGs. 28A-28K:** Temporal patterns of paracrine signaling. **(FIG. 34A)** Cell clusters determined by Louvain-Jaccard community detection algorithm. **(FIG. 34B)** Temporal pattern of the net potential for paracrine signaling between contemporaneous cells in 2i condition. Each dot represents the aggregated interaction score across all ligand-receptor pairs for a given combination of clusters from **(FIG. 34A)** (see STAR Methods for details). **(FIGs. 34C-34E)** Changes in the standardized interaction scores for top ligand-receptor pairs between ancestors of stromal cells and ancestors of iPSCs **(FIG. 34C),** neural-like cells **(FIG. 34D),** and trophoblast cells **(FIG. 34E).**

[0063]     **FIGs. 35A-35B** - Related to **FIGs. 29A-29D:** Comparison with alternate methods. **(FIG. 35A)** Monocle2 computes a graph upon which each cell is embedded. The graph, which

consists of 5 segments, is visualized in the upper-left pane. The 5 segments are visualized on our FLE in the 5 remaining panels of **(FIG. 35A).** Segment 1 (green) consists of day 0 cells together with day 18 Stromal cells. Segments 2 and 3 consist of cells from day 2 - 8 that supposedly arise from Segment 1 cells. Segment 3 gives rise to Segments 4 (purple) and 5 (red). Segment 4 contains the cells we identify as on the MET region and Segment 5 contains the iPSCs, Trophoblasts, and Neural populations, which Monocle2 infers come directly from the non-proliferative cells in segment 3. **(FIG. 35B)** URD computes a graph representing random walks from a collection of tips to a root. This graph, which consists of **7** segments, is visualized in the upper-left pane. The **7** segments are visualized on our FLE in the remaining panels of **(FIG. 35B).** Segment 1 (magenta) contains the day 0 MEF cells. The first bifurcation occurs on day 0.5, where segment 2 (consisting of day 0.5 cells) splits off from segment 3 (consisting of day 12-18 Stromal cells). Segment 2 splits to give rise to Segment 4 (consisting of day 2 cells) and Segment 5 consisting of day 12-18 Trophoblasts and Epithelial cells. Segment 4 splits on day 3 to give rise to Segment 6 (consisting of a diverse population including day 3 cells and day 14-18 iPSCs) and Segment **7** (consisting of a diverse population including day 3 cells and day 12-18 Neural-like cells).

[0064]      **FIGs. 36A-36F** - Related to **FIGs. 29A-29D:** *Obox6 + Obox6* graphs. **(FIGs. 36A-36C)** Identical to **FIGs. 29A-29C** except here we show results for serum conditions. **(FIG. 36D)** Percentage of Oct4-EGFP+ cells at day 16 of reprogramming from secondary MEFs by lentiviral overexpression of *Oct4, Kl/4, Sox2,* and *Myc* (OKSM) combined with either *Zfp42, Obox6,* or an empty control, in either 2i or serum conditions. Oct4-EGFP+ cells were measured by flow cytometry. Plot includes the percentage of Oct4-EGFP+ cells in three biological replicates (for *Zfp 42* and *Obox6* overexpression, or an empty control) from five independent experiments (Exp). **(FIG. 36E, FIG. 36F)** Number of Oct4-EGFP+ colonies at day 16 of reprogramming from primary MEFs by lentiviral overexpression of individual *Oct4, Kl/4, Sox2,* and *Myc* combined with either *Zfp 42, Obox6,* or an empty control in **(FIG. 36E)** 2i and **(FIG. 36F)** serum conditions. Plot includes the number of Oct4-EGFP+ cells in three biological replicates (for *2fp42* and *Obox6* overexpression, or an empty control) from two independent experiments (Exp).

[0065]      **FIG. 37** - Effects of GDF9 on reprogramming efficiency.

[0066]      **FIG. 38** shows adding GDF9 to the medium resulted in more iPSCs.

## DETAILED DESCRIPTION OF THE EXAMPLE EMBODIMENTS

### General Definitions

[0067]     Unless defined otherwise, technical and scientific terms used herein have the same meaning as commonly understood by one of ordinary skill in the art to which this disclosure pertains.  Definitions of common terms and techniques in molecular biology may be found in Molecular Cloning: A Laboratory Manual, 2nd edition (1989) (Sambrook, Fritsch, and Maniatis); Molecular Cloning: A Laboratory Manual, 4th edition (2012) (Green and Sambrook); Current Protocols in Molecular Biology (1987) (F.M. Ausubel et al.  eds.); the series Methods in Enzymology (Academic Press, Inc.): PCR 2: A Practical Approach (1995) (M.J. MacPherson, B.D. Hames, and G.R. Taylor eds.): Antibodies, A Laboraotry Manual (1988) (Harlow and Lane, eds.): Antibodies A Laboraotry Manual, 2nd edition 2013 (E.A. Greenfield ed.); Animal Cell Culture (1987) (R.I. Freshney, ed.); Benjamin Lewin, Genes IX, published by Jones and Bartlet, 2008 (ISBN 0763752223); Kendrew *et al.* (eds.), The Encyclopedia of Molecular Biology, published by Blackwell Science Ltd., 1994 (ISBN 0632021829); Robert A. Meyers (ed.), Molecular Biology and Biotechnology: a Comprehensive Desk Reference, published by VCH Publishers, Inc., 1995 (ISBN 9780471 185710); Singleton *et al,* Dictionary of Microbiology and Molecular Biology 2nd ed., J. Wiley & Sons (New York, N.Y. 1994), March, Advanced Organic Chemistry Reactions, Mechanisms and Structure 4th ed., John Wiley & Sons (New York, N.Y. 1992); and  Marten H. Hofker and Jan van Deursen, Transgenic Mouse Methods and Protocols, 2nd edition (201 1) .

[0068]     As used herein, the singular forms "a", "an", and "the" include both singular and plural referents unless the context clearly dictates otherwise.

[0069]     The term "optional" or "optionally" means that the subsequent described event, circumstance or substituent may or may not occur, and that the description includes instances where the event or circumstance occurs and instances where it does not.

[0070]     The recitation of numerical ranges by endpoints includes all numbers and fractions subsumed within the respective ranges, as well as the recited endpoints.

[0071]     The terms "about" or "approximately" as used herein when referring to a measurable value such as a parameter, an amount, a temporal duration, and the like, are meant to encompass

variations of and from the specified value, such as variations of +/-10% or less, *+1-5%* or less, +/-1% or less, and +/-0.1% or less of and from the specified value, insofar such variations are appropriate to perform in the disclosed invention. It is to be understood that the value to which the modifier "about" or "approximately" refers is itself also specifically, and preferably, disclosed.

[0072]     Reference throughout this specification to "one embodiment", "an embodiment," "an example embodiment," means that a particular feature, structure or characteristic described in connection with the embodiment is included in at least one embodiment of the present invention. Thus, appearances of the phrases "in one embodiment," "in an embodiment," or "an example embodiment" in various places throughout this specification are not necessarily all referring to the same embodiment, but may. Furthermore, the particular features, structures or characteristics may be combined in any suitable manner, as would be apparent to a person skilled in the art from this disclosure, in one or more embodiments. Furthermore, while some embodiments described herein include some but not other features included in other embodiments, combinations of features of different embodiments are meant to be within the scope of the invention. For example, in the appended claims, any of the claimed embodiments can be used in any combination.

[0073]     All publications, published patent documents, and patent applications cited herein are hereby incorporated by reference to the same extent as though each individual publication, published patent document, or patent application was specifically and individually indicated as being incorporated by reference.

**Overview**

[0074]     Embodiments disclosed herein provide methods and systems intended to reflect Waddington's image of marbles rolling within a development landscape. It captures the notion that cells at any position in the landscape have a *distribution* of both probable origins and probable fates. It seeks to reconstruct both the landscape and probabilistic trajectories from scRNA-seq data at various points along a time course. Specifically, it uses time-course data to infer how the probability distribution of cells in gene-expression space evolves over time, by using the mathematical approach of Optimal Transport (OT). The utility of this method is demonstrated in the context of reprogramming of fibroblasts to induced pluripotent stem cells

(iPSCs). However, the same method may be applied to other cell development and biological context where an understanding of cell orgins, trajectories, and fates is needed. For ease of reference, the methods disclosed herein and in their various embodiments may be referred to collectively as "Waddington-OT." As demonstrated herein, Waddington-OT readily rediscovers known biological features of reprogramming, including that successfully reprogrammed cells exhibit an early loss of fibroblast identity, maintain high levels of proliferation, and undergo a mesenchymal-to-epithelial transition before adopting an iPSC-like state *(12)*. In addition, by exploiting single-cell resolution and the new model, it also extends these results by (1) identifying alternative cell fates, including senescence, apoptosis, neural identity, and placental identity; (2) quantifying the portion of cells in each state at each time point; (3) inferring the probable origin(s) and fate(s) of each cell and cell class at each time point; (4) identifying early molecular markers associated with eventual fates; and (5) using trajectory information to identify transcription factors (TFs) associated with the activation of different expression programs. In particular, TFs that are putative regulators of neural identity, placental identity, and pluripotency during reprogramming, and we experimentally demonstrate that one such TF, *Obox6,* enhances reprogramming efficiency are provided. Together, the data provide a high-resolution resource for studying the roadmap of reprogramming, and the methods provide a general approach for studying cellular differentiation in natural or induced settings.

[0075]     Prior to describing implementation of the methods in detail, the following overview and definitions utilized in execution of the method are defined.

[0076]     scRNA-seq may be obtained from cells using standard techniques known in the art. A collection of mRNA levels for a single cell is called an *expression profile* and is often represented mathematically by a vector in *gene expression space.* This is a vector space that has a dimension corresponding to each gene, with the value of the ith coordinate of an expression profile vector representing the number of copies of mRNA for the ith gene. Note that real cells only occupy an integer lattice in gene expression space (because the number of copies of mRNA is an integer), but it is assumed herein that cells can move continuously through a real-valued G dimensional vector space.

[0077]     As an individual cell changes the genes it expresses over time, it moves in gene expression space and describes a trajectory. As a population of cells develops and grows, a

distribution on gene expression space evolves over time. When a single cell from such a population is measured with single cell RNA sequencing, a noisy estimate of the number of molecules of mRNA for each gene is obtained. The measured expression profile of this single cell is represented as a sample from a probability distribution on gene expression space. This sampling captures both (a) the randomness in the single cell RNA sequencing measurement process (due to sub-sampling reads, technical issues, etc.) and (b) the random selection of a cell from a population. This probability distribution is treated as nonparametric in the sense that it is not specified by any finite list of parameters.

[0078]     A precise mathematical notion for a *developmental process* as a generalization of a stochastic process is provided below. A goal of the methods disclosed herein is to infer the ancestors and descendants of subpopulations evolving according to an unknown developmental process. While not bound by a particular theory, this may be possible over short time scales because it is reasonable to assume that cells don't change too much and therefore it can be inferred which cells go where.

[0079]     In certain example embodiments, the following definitions to define a precise notion of the developmental trajectory of an individual cell and its descendants are used. It is a continuous path in gene expression that bifurcates with every cell division.

Formally, consider a cell $x(0) \in \mathbb{R}^G$. *Let* $k(t) \geq 0$ *specify the number of descendants at time* t, *where* $k(0) = 1$. *A single cell developmental trajectory is a continuous function*

$$x : [0, T) \to \underbrace{\mathbb{R}^G \text{ x } \mathbb{R}^G \text{ x } \ldots \text{ x } \mathbb{R}^G}_{k(t) \text{ times}}.$$

*This means that x(t) is a k(t)-tuple of cells, each represented by a vector* $\mathbb{R}^G$:

$$x(t) = \left( x_1(t), \ldots, x_{k(t)}(t) \right).$$

*Cells xi(t),  ...., $x_{k(t)}(t)$ as the descendants of x(0).*

[0080]     $\mathbb{R}^G$ and $R^G$ are used interchangeably.

[0081]     Note that the temporal dynamics of an individual cell cannot be directly measured because scRNA-Seq is a destructive measurement process: scRNA-Seq lyses cells so it is only possible to measure the expression profile of a cell at a single point in time. As a result, it is not possible to directly measure the descendants of that cell, and it is (usually) not possible to directly measure which cells share a common ancestor with ordinary scRNA-Seq. Therefore the

full trajectory of a specific cell is unobservable. However, one can learn something about the probable trajectories of individual cells by measuring snapshots from an evolving population.

[0082]    Published methods typically represent the aggregate trajectory of a population of cells with a graph. While this recapitulates the branching path traveled by the descendants of an individual cell, it may over-simplify the stochastic nature of developmental processes. Individual cells have the potential to travel through different paths, but in reality any given cell travels one and only one such path. The methods disclosed herein help to describe this potential, which might not be a represented by a graph as a union of one dimensional paths.

[0083]    Instead, a developmental process is defined to be a time-varying distribution on gene expression space. The word distribution is used to refer to an object that assigns mass to regions of $\mathbb{R}^G$. Note that a distinction is made between distribution and *probability* distribution, which necessarily has total mass 1. Distributions are formally defined as generalized functions (such as the delta function $\delta_\chi$) that act on test functions. A used herein a "distribution" is the same as a measure. One simple example of a distribution of cells is that a set of cells $x_1, \ldots, x_n$ can be represented by the distribution

$$\mathbb{P} = \sum_{i=1}^{n} \delta_{x_i}.$$

Similarly, a set of single cell trajectories may be represented $xj(t), \ldots, x_n(t)$ with a distribution over trajectories. A developmental process $\mathbb{P}_*$ is a time-varying distribution on gene expression space. A developmental process generalizes the definition of *stochastic process.* A developmental process with total mass 1 for all time is a (continuous time) stochastic process, *i.e.* an ordered set of random variables with a particular dependence structure. Recall that a stochastic process is determined by its temporal dependence structure, *i.e.* the coupling between random variables at different time points. The coupling of a pair of random variables refers to the structure of their joint distribution. The notion of coupling for developmental processes is the same as for stochastic processes, except with general distributions replacing probability distributions.

[0084]    A coupling of a pair of distributions P, Q on $R^G$ is a distribution $\pi$ on $R^G \times R^G$ with the property that $\pi$ has P and Q as its two marginals. A coupling is also called a transport map.

**[0085]**     As a distribution on the product space $R^G$ x $R^G$, a transport map π assigns a number π(A, B) to any pair of sets A,B $\subset$ $R^G$ .

$$\pi(A, B) = \int_{x \in A} \int_{y \in B} \pi(x, y) dx dy.$$

When π is the coupling of a developmental process, this number π(A, B) represents the mass transported from A to B by the developmental process. This is the amount of mass coming from A and going to B. When a particular destination is note specified, the quantity π(A, ·) specifies the full distribution of mass coming from A. This action may be referred to as *pushing* A through the transport map π. More generally, we can also push a *distribution* μ forward through the transport map π via integration

$$\mu \mapsto \int \pi(x, \cdot) d\mu(x).$$

The reverse operation is referred to as pulling a set B back through π. The resulting distribution π(·, B) encodes the mass ending up at B. Distributions μ can also be pulled back through π in a similar way:

$$\mu \mapsto \int \pi(\cdot, y) d\mu(y).$$

This may also be referred as *back-propagating* the distribution μ (and to pushing μ forward as *forward propagation).*

**[0086]**     Recall that a stochastic process is Markov if the future is independent of the past, given the present. Equivalently, it is fully specified by its couplings between pairs of time points. A general stochastic process can be specified by further higher order couplings. Markov developmental processes, which are defined in the same way:

**[0087]**     *A Markov developmental process* $P_t$ *is a time-varying distribution on* $R^G$ *that is completely specified by couplings between pairs of time points.* It is an interesting question to what extent developmental processes are Markov. On gene expression space, they are likely not Markov because, for example, the history of gene expression can influence chromatin modifications, which may not themselves be reflected in the observed expression profile but could still influence the subsequent evolution of the process. However, it is possible that developmental processes could be considered Markov on some augmented space.

**[0088]**    A definition of descendants and ancestors of subgroups of cells evolving according to a Markov developmental process is now provided. The earlier definition of descendants is extended as follows:  *Consider a set of cells* $S \subset R^G$*, which live at time* $x_1$ *are part of a population of cells evolving according to a Markov developmental process* $P^{\wedge}$*. Let* $\pi$ *denote the transport map for* $V^{\wedge}$ *from  time* $\backslash\backslash$ *to time* $x^{\wedge}$*. The descendants of* $S$ *at time* $x_2$ *are obtained by pushing* $S$ *through the transport map* $\pi$. Note that if a developmental process is not Markov, then the descendants of S are not well defined. The descendants would depend on the cells that gave rise to $S$, which we refer to as the *ancestors* of $S$.

**[0089]**    Definition 6 (ancestors in a Markov developmental process). *Consider a set of cells* $S \subset R^G$*, which live at time* $t_2$ *and are part of a population of cells evolving according to a Markov developmental process* $P^{\wedge}$*. Let* $\pi$ *denote the transport map for* $V^{\wedge}$ *from  time* $x_2$ *to time* $X\backslash$*. The ancestors of $S$  at time* $t_1$ *are obtained by pushing* $S$ *through the transport map* $\pi$.

**Empirical developmental processes**

**[0090]**    In certain aspects, a goal of the embodiments disclosed herein is to track the evolution of a developmental process from a scRNA-Seq time course. Suppose we are given input data consisting of a sequence of sets of single cell expression profiles, collected at T different time slices of development. Mathematically, this time series of expression profiles is a sequence of sets $S_1,...,S_T \subset R^G$ collected at times $\ddot{i}_1,...,\ddot{i}_T \in R$.

**[0091]**     Developmental time series. *A developmental time series is a sequence of samples from a developmental process* $P_t$ *on* $R^G$*. This is a sequence of sets* $S_1, \ldots, S_N \subset R^G$*. Each* $S_i$ *is a set of expression profiles in* $R^G$ *drawn i.i.d from  the probability distribution obtained by normalizing the distribution* $P_{t_i}$ *to have total mass X.* From this input data, we form an empirical version of the developmental process. Specifically, at each time point tj we form the empirical probability distribution supported on the data $x \in S_i$ is formed. This is summarized inin the following definition:

**[0092]**    Empirical developmental process. *An empirical developmental process* $\hat{P}_t$ *is a time vary-ing distribution constructed from  a developmental time course* $S_1, \ldots, S_N$ *:*

$$\hat{\mathbf{P}}_{t_i} = \frac{1}{|S_i|} \sum_{x \in S_i} \delta_x.$$

*he empirical developmental process is undefined for* $t \in / \{t_1, \ldots, t_N\}$.

**[0093]** Our goal is to recover information about a true, unknown developmental process $P_t$ from the empirical developmental process $\hat{P}_t$. The measurement process of single cell RNA-Seq destroys the coupling, and the observed empirical developmental process does not come with an informative coupling between successive time points. Over short time scales, it is reasonable to assume that cells do not change too much and therefore inferences regarding which cells go where and estimate the coupling.

**[0094]** This may be done with *optimal transport:* the transport map $\pi$ that minimizes the total work required for redistributing $\hat{P}_{t_i}$ to $\hat{P}_{t_{i+1}}$. is selected. One motivation for minimizing this objective, is a deep relationship between optimal transport and dynamical systems that provides a direct connection to Waddington's landscape: the optimal transport problem can formulated as a *least-action advection* of one distribution into another according to an unknown velocity field (see Theorem 1 in Section 6 below). At a high level, differentiation follows a velocity field on gene expression space, and the potential inducing this velocity field is in direct correspondence with Waddington's landscape^.

**Optimal transport for scRNA-Seq time series**

**[0095]** A process for how to compute probabilistic flows from a time series of single cell gene expression profiles by using optimal transport (SI) is provided. The embodiments disclosed herein show how to compute an optimal *coupling* of adjacent time points by solving a convex optimization problem.

**[0096]** Optimal transport defines a metric between probability distributions; it measures the total distance that mass must be transported to transform one distribution into another. For two measures $P$ and $Q$ on $R^G$, a *transport plan* is a measure on the product space $R^G \times R^G$ that has marginals $P$ and $Q$. In probability theory, this is also called a *coupling*. Intuitively, a transport plan $\pi$ can be interpreted as follows: if one picks a point mass at position x, then $\pi(\chi, \cdot)$ gives the distribution over points where x might end up.

[0097]    If c(x, y) denotes the cost$^2$ of transporting a unit mass from x to y, then the expected cost under a transport plan π is given by

$$\iint c(x,y)\pi(x,y)dxdy.$$

The optimal transport plan minimizes the expected cost subject to marginal constraints:

$$\underset{\pi}{\text{minimize}} \quad \iint c(x,y)\pi(x,y)dxdy$$

$$\text{subject to} \quad \int \pi(x,\text{-})dx = \mathbb{Q}$$

$$\int \pi(\cdot,y)dy = \mathbb{P}.$$

[0098]    Note that this is a linear program in the variable π because the objective and constraints are both linear in π. Note that the optimal objective value defines the *transport distance* between P and Q (it is also called the Earthmover's distance or Wasserstein distance). Unlike most other ways to compare distributions (such as KL-divergence or total variation), optimal transport takes the geometry of the underlying space into account. For example, the KL-Divergence is infinite for any two distributions with disjoint support, but the transport distance between two unit masses depends on their separation.

[0099]    When the measures P and Q are supported on finite subsets of $R^G$, the transport plan is a matrix whose entries give transport probabilities and the linear program above is finite dimensional. In this context, *empirical distributions* are formed from the sets of samples $S_1, \ldots,$ $S_T$ :

$$\hat{\mathbb{P}}_{t_i} = \frac{1}{|S_i|} \sum_{x \in S_i} \delta_x,$$

were $\delta_\chi$ denotes the Dirac delta function centered at $x \in R^G$. These empirical distributions $\hat{P}_{t_i}$ are definitely supported, and so it is possible solve the linear program[l] with $P=\hat{P}_{t_j}$ and $Q=\hat{P}_{t_{i+1}}$.

[00100]    However, the classical formulation [1] does not allow cells to grow (or die) during transportation (because it was designed to move piles of dirt and conserve mass). When the classical formulation is applied to a time series with two distinct subpopulations proliferating at

different rates[3], the transport map will artificially transport mass between the subpopulations to account for the relative proliferation. Therefore, we modify the classical formulation of optimal transport in equation [1] is modified to allow cells to grow at different rates.

**[00101]** Is it assumed that a cell's measured expression profile x determines its growth rate g(x). This is reasonable because many genes are involved in cell proliferation (e.g. cell cycle genes). It is further assumed g(x) is a known function (based on knowledge of gene expression) representing the exponential increase in mass per unit time, but also note that the growth rate can be allowed to be miss-specified by leveraging techniques from *unbalanced transport* (S2). In practice, g(x) is defined in terms of the expression levels of genes involved in cell proliferation.

**[00102]** **Derivation of transport with growth:** For any cell x $\in$ Sj-j, let r(x, y) be the fraction of x that transitions towards y. Then the amount of probability mass from x that ends up at y (after proliferation) is

$$r(x,\ y)g(x)^{\Delta t},$$

where $A_t = t_i+1 - t_i$. The total amount of mass that comes from x can be written two ways:

$$\sum_{y \in S_{i+1}} r(x,\ y)g(x)^{\Delta t} \approx g(x)^{\Delta t} d\hat{\mathbb{P}}_{t_i}(x)\ .$$

This gives us a first constraint. Similarly, there is also the constraint that *the total mass observed at y is equal to the sum of masses coming from each* x *and ending up at y*. In symbols,

$$d\hat{\mathbb{P}}_{t_{i+1}}(y) \sum_{x \in S_i} g(x)^{\Delta t} \sim \sum_{x \in S_i} r(x, V)9(x)^{\Delta t} \qquad \text{for each } y \in S_{i+1}\ .$$

The factor $_{x \in gj} g(x)^{\Delta t}$ on the left hand side accounts for the overall proliferation of all the cells from $S_i$. Note that this factor is required so that the constraints are consistent: when one sums up both sides of the first constraint over x, this must equal the result of summing up both sides of the second constraint over y. Finally, for convenience these constraints are rewritten in terms of the optimization variable

$$\pi(x,\ y) = r(x,\ y)g(x)^{\Delta t}.$$

Therefore, to compute the transport map between the empirical distributions of expression profiles observed at time $t_i$ and $t_i+1$, the following linear program is set up:

$$\underset{\pi}{\text{minimize}} \quad \sum_{x \in S_i} \sum_{y \in S_{i+1}} c(x,y)\pi(x,y)$$

$$\text{subject to} \quad \sum_{x \in S_i} \pi(x,y) \text{ sa } d\hat{\mathbb{P}}_{t_{i+1}}(\mathbf{y}) \sum_{x \in S_i} g(x)^{\Delta_t}$$

$$\sum_{y \in S_{i+1}} \pi(x,y) \approx d\hat{\mathbb{P}}_{t_i}(x)g(x)^{\Delta_t}$$

[00103]   **Regularization and algorithmic considerations:** Fast algorithms have been recently developed to solve an entropically regularized version of the transport linear program (S3). Entropic regularization means adding the entropy $H(\pi) = E_\pi \log \pi$ to the objective function, which penalizes deterministic transport plans (a purely deterministic transport plan would have only one nonzero entry in each row). Entropic regularization speeds up the computations because it makes the optimization problem strongly convex, and gradient ascent on the dual can be realized by successive diagonal matrix scalings (S3). These are very fast operations. This scaling algorithm has also been extended to work in the setting of *unbalanced transport,* where equality constraints are relaxed to bounds on KL-divergence (S2). This allows the growth rate function g(x) to be misspecified to some extent.

[00104]   Both entropic regularization and unbalanced transport may be used. To compute the transport map between the empirical distributions of expression profiles observed at time $t_i$ and $t_{i+1}$, the embodiments disclosed herein solve the following optimization problem:

$$\underset{\pi}{\text{minimize}} \quad \sum_{x \in S_i} \sum_{y \in S_{i+1}} c(x,y)\pi(x,y) - \epsilon H\{\pi\}$$

$$\text{subject to} \quad KL\left[\sum_{x \in S_i} \pi(x,y) \middle\| d\hat{\mathbb{P}}_{t_{i+1}}(y) \sum_{x \in S_i} g(x)^{\Delta_t}\right] \leq \frac{1}{\lambda_1}$$

$$KL\left[\sum_{y \in S_{i+1}} \pi(x,y) \middle\| d\hat{\mathbb{P}}_{t_i}(x)g(x)^{\Delta_t}\right] \leq \frac{1}{\lambda_2}$$

where $\epsilon$, $\lambda_1$ and $\lambda_2$ are regularization parameters. This is a convex optimization problem in the matrix variable $\pi \in R^{Ni \times Ni+1}$, where $N_I = |\mathbf{g_I}| j_s$ the num ber of cells sequenced at time tj. It takes about 5 seconds to solve this unbalanced transport problem using the scaling algorithm of Chizat et al. 2016 (S2) on a standard laptop with Nj $\approx$ 5000. Note that the densities (on the discrete set Sj) of the empirical distributions specified in equation [2] are simply $d\hat{P}_t(x) = 1$ .

However, in principle one could use nonuniform empirical distributions (e.g. i Ni if one wanted to include information about cell quality).

[00105]    To summarize: given a sequence of expression profiles $S_i$, . . . , $S_T$ , the optimization problem [5] for each successive pair of time points $S_i$, $S_i+1$ is solved. This gives us a sequence of transport maps as illustrated in **FIG. 3.**

[00106]    To make this more precise, consider a single cell y $\in$ Sj. The column π(·, y) of the transport map π from tj_i to $t_i$ describes the contributions to y of the cells in $S_i-1$. This is the origin of y at the time point tj_i . Similarly, the row r(y, ·) of the transition map from tj to $t_i+1$ describes the probabilities y would transition to cells in $S_i+1$. These are the fates of y, i.e. the descendants of y.

[00107]    The origin of y further back in time may be computed via matrix multiplication: the contributions to y of cells in $S_i-2$ are given by a column of the matrix

$$\tilde{K}_{[i-2,i]} = \pi_{[i-2,i-1]}\pi_{[i-1,i]}.$$

[00108]    This matrix $\tilde{\pi}_{[i-2,i]}$ represents the inferred transport from time point tj_2 to $t_i$, and note it with a tilde to distinguish it from the maps computed directly from adjacent time points. Note that, in principle, the transport between any non-consecutive pairs of time points $S_i$, Sj, may be directly computed but it is not anticipated that the principle of optimal transport to be as reliable over long time gaps.

[00109]    Finally, note that expression profiles can be interpolated between pairs of time points by averaging a cell's expression profile at time $t_i$ with its fated expression profiles at time $t_i+1$.

**Transport maps encode regulatory information**

[00110]    Transport maps can encode regulatory information, and provided herein are methods on how to set up a regression to fit a regulatory function to our sequence of transport maps. It is assumed that a cell's trajectory is cell-autonomous and, in fact, depends only on its own internal gene expression. We know this is wrong as it ignores paracrine signaling between cells, and we return to discuss models that include cell-cell communication at the end of this section. However, this assumption is powerful because it exposes the time-dependence of the stochastic process $P_t$ as arising from pushing an initial measure through a differential equation:

$$\dot{x} = /(*).$$

**[00111]** Here f is a vector field that prescribes the flow of a particle x (see fig. 3 for a cartoon illustration of a distribution flowing according to a vector field). Our biological motivation for estimating such a function f is that it encodes information about the regulatory networks that create the equations of motion in gene-expression space.

**[00112]** We propose to set up a regression to learn a regulatory function f that models the fate of a cell at time $t_{i+1}$ as a function of its expression profile at time ¾. For motivation that the transport maps might contain information about the underlying regulatory dynamics, we appeal to a classical theorem establishing a dynamical formulation of optimal transport.

**[00113]** Theorem 1 (Benamou and Brenier, 2001). *The optimal objective value of the transport problem* [1] *is equal to the optimal objective value of the following optimization problem:*

$$\underset{\rho,v}{\text{minimize}} \quad \int_0^1 \int_{\mathbb{R}^G} \|v(t,x)\|^2 \rho(t,x)dtdx$$
$$\text{subject to} \quad \rho(0,\cdot) = \mathbb{P}, \quad \rho(1,\cdot) = \mathbb{Q}$$
$$\nabla \cdot (\rho v) = \frac{\partial \rho}{\partial t}$$

**[00114]** In this theorem, v is a vector-valued velocity field that advects4 the distribution p from P to Q, and the objective value to be minimized is the kinetic energy of the flow (mass **x** squared velocity). Intuitively, the theorem shows that a transport map $\pi$ can be seen as a point-to-point summary of a least-action continuous time flow, according to an unknown velocity field. While the optimization problem [8] can be reformulated as a convex optimization problem, and modified to allow for variable growth rates, it is inherently infinite dimensional and therefore difficult to solve numerically.

**[00115]** We therefore propose a tractable approach to learn a static regulatory function f from our sequence of transport maps. Our approach involves sampling pairs of points using the couplings from optimal transport, and solving a regression to learn a regulatory function that predicts the fate of a cell at time $t_{i+1}$ as a function of its expression profile at time $t_i$:

**[00116]** **Regulatory network regression:** For each pair of time points ¾,¾+!, we consider the pair of random variables $X_t$ ,$X_t$ jointly distributed according to $r_{[t,t]}$, (which we obtained from

34

the i i+1 i i+1 transport map $\pi_{[t_i, t_{i+1}]}$ by removing the effect of proliferation as in equation [3]). We set up the following optimization problem over regulatory functions f:

$$\min_{f \in \mathcal{F}} \quad \mathbb{E}_r \left\| \frac{X_{t_i} - X_{t_{i+1}}}{\Delta_t} - f(X_{t_i}) \right\|^2 .$$

Here F specifies a parametric function class to optimize over.

[00117]    **Cell non-autonomous processes:** We conclude our treatment of gene regulatory networks by discussing an approach to cell-cell communication. Note that the gradient flow [8] only makes sense for cell autonomous processes. Otherwise, the rate of change in expression $\dot{x}$ is not just a function of a cell's own expression vector x(t), but also of other expression vectors from other cells. We can accommodate cell non-autonomous processes by allowing f to also depend on the full distribution $P_t$

$$\frac{dx}{dt} = f(x, \mathbb{P}_t).$$

4. Extensions to continuous time.

[00118]    In this section we discuss how our method could be improved by going beyond pairs of time points to track the continuous evolution of $P_t$. We begin by pointing out a peculiar behavior of our method: whenever we have a time point with few sampled cells, our method is forced through an information bottleneck. As an extreme example - suppose we had a time point with only one cell. Everything would transition through that single cell, which is absurd! In this extreme case, we would be better off ignoring the time point. We therefore propose a smoothed approach that shares information between time slices and gracefully improves as data is added.

[00119]    Our continuous-time formulation is based on locally-weighted averaging, an elementary interpolation technique. Recall that given noisy function evaluations $y_j \sim f(x_j)$, one can interpolate f by averaging the $y_i$ for all $x_i$ close to a point of interest x:

$$f(x) \approx \sum_i \alpha_i f(x_i),$$

where $\alpha_i$ are weights that give more influence to nearby points

.........

**[00120]** In our setup, we seek to interpolate a distribution-valued function $P_t$ from the collections of i.i.d. samples $S_i, \ldots, S_T$. We can interpolate a distribution-valued function by computing the *barycenter* (or centroid) of nearby time points with respect to the optimal transport metric. The *transport barycenter* of

$$\underset{\mathbb{Q}}{\text{minimize}} \quad \sum_{i=l}^{T} \alpha_i W^2(\mathbb{P}_i, \mathbb{Q}),$$

where W (P, Q) denotes the transport distance (or Wasserstein distance) between P and Q. The transport distance is defined by the optimal value of the transport problem [1]. The weights $\alpha_i$ can be chosen to interpolate about time point t by setting, for example,

$$\underset{\mathbb{Q}}{\text{minimize}} \quad \sum_{i=l}^{T} \alpha_i G^2(\hat{\mathbb{Y}}_{t-}, \mathbb{Q}),$$

where G(P, Q) denotes our modified transport distance from equation [5]. To solve this optimization problem, we can fix the support of Q to the samples observed at all time points $U_{i=\backslash}^{T} S_i$. Then we can applythescalingalgorithmforunbalancedbarycentersduetoChizatetal.    (S2).

**[00121]** However, fixing the support of the barycenter ahead of time may not be completely satisfactory, and this motivates further research in the computation of transport barycenters: can we design an algorithm to solve for the barycenter Q without fixing the support in advance? Is there a dynamic formulation for barycenters analogous to the Brenier Benamou formula of Theorem 1, and can we leverage it to better learn gene regulatory networks?

**[00122]** Finally, we conclude this section with the observation that this continuous-time approach could pro-vide a principled approach to sequential experimental design. We can identify optimal time points for further data collection by examining the loss function (fit of barycenter) across time, and adding data where the fit is poor. Moreover, we could also use this continuous time approach to test the principle of optimal transport by withholding some time points and testing the quality of the interpolation against the held-out truth.

**Example System Architectures**

**[00123]** **FIG. 1** is a block diagram depicting a system for mapping developmental trajectories of cells using single cell sequencing data, in accordance with certain example embodiments. As depicted in **FIG. 1,** the system 100 includes network devices 110, 115, and 120, that are

configured to communicate with one another via one or more networks 105. In some embodiments, a user associated with the user device 115, may have to install an application and/or make a feature selection to obtain the benefits of the techniques described herein.

[00124]    Each network 105 includes a wired or wireless telecommunication means by which network devices (including devices 110, 135 and 140) can exchange data. For example, each network 105 can include a local area network ("LAN"), a wide area network ("WAN"), an intranet, an Internet, a mobile telephone network, or any combination thereof. Throughout the discussion of example embodiments, it should be understood that the terms "data" and "information" are used interchangeably herein to refer to text, images, audio, video, or any other form of information that can exist in a computer-based environment.

[00125]    Each network device 110, 135 and 140 includes a device having a communication module capable of transmitting and receiving data over the network 105. For example, each network device 110, 135 and 140 can include a server, desktop computer, laptop computer, tablet computer, a television with one or more processors embedded therein and / or coupled thereto, smart phone, handheld computer, personal digital assistant ("PDA"), or any other wired or wireless, processor-driven device. In the example embodiment depicted in **FIG. 1,** the network devices (including systems 110, 115 and 120) are operated by end-users or consumers, merchant operators (not depicted), and feedback system operators (not depicted), respectively.

[00126]    A user can use the application 112, such as a web browser application or a stand-alone application, to view, download, upload, or otherwise access documents or web pages via a distributed network 105. The network 105 includes a wired or wireless telecommunication system or device by which network devices (including devices 110, 115 and 120) can exchange data. For example, the network 105 can include a local area network ("LAN"), a wide area network ("WAN"), an intranet, an Internet, storage area network (SAN), personal area network (PAN), a metropolitan area network (MAN), a wireless local area network (WLAN), a virtual private network (VPN), a cellular or other mobile communication network, Bluetooth, NFC, or any combination thereof or any other appropriate architecture or system that facilitates the communication of signals, data, and/or messages. Throughout the discussion of example embodiments, it should be understood that the terms "data" and "information" are used

interchangeably herein to refer to text, images, audio, video, or any other form of information that can exist in a computer based environment.

[00127]    The communication application 112 can interact with web servers or other computing devices connected to the network 105, including the single cell sequencing system 110 and optimal transport system 120.

[00128]    It will be appreciated that the network connections shown are example and other means of establishing a communications link between the computers and devices can be used. Moreover, those having ordinary skill in the art having the benefit of the present disclosure will appreciate that the single cell sequencing system 110, user device 115, and optimal transport system 120 illustrated in **FIG. 1** can have any of several other suitable computer system configurations.  For example a user device 115 embodied as a mobile phone or handheld computer may not include all the components described above

**Example Processes**

[00129]    The example methods illustrated in **FIG. 2** are described hereinafter with respect to the components of the example operating environment **100.** The example methods of **FIG. 2** may also be performed with other systems and in other environments

[00130]    **FIG. 2** is a block flow diagram depicting a method **200** to determine developmental trajectories of cells, in accordance with certain example embodiments.

[00131]    Method **200** begins at block **205,** where the optimal transport module **125** performs optimal transport analysis on single cell RNA-seq data (scRNA-seq) from a time course, by calculating optimal transport maps and using them to find ancestors, descendants and trajectories for any set of cells. Given a subpopulation of cells, the sequence of ancestors coming before it and descendants coming after it are referred to as its developmental trajectory. Further example of how development trajectories may be computed in block **205** is described in Example 1 below. Briefly, transport maps are calculated, as described above, between consecutive time points, with cells allowed to grow according to a gene-expression signature of cell proliferation. From these transport maps, the forward and backward transport possibilities can be calculated between any two classes of cells at any time points. For example, a successfully reprogrammed cell at day 16 and use back-propagation to infer the distribution over their precursors at day 12. This can then be further propagated back to day 11, and so one to obtain the ancestor

distributions at all previous time points. From this trend in gene expression over time may be plotted. See **FIGs. 9A-9D.**

**[00132]** In certain example embodiments, an expression matrix may be computed by the optimal transport module 125 from the scRNA-Seq data. Sequence reads may be aligned to obtain a matrix $U$ of UMI counts, with a row for each gene and column for each cell. To reduce variation due to fluctuations in the total number of transcripts per cell, we divide the UMI vector for each cell by the total number of transcripts in that cell. Thus we define the expression matrix E in terms of the UMI matrix U via:

$$E_{ij} = \frac{U_{ij}}{\sum_{i=1}^{G} U_{ij}} \times 10^4.$$

**[00133]** Two variance-stabilizing transforms of the expression matrix E may be used for further analysis. In particular

1. $\tilde{E}$ to be the log-normalized expression matrix. The entries of $\tilde{E}$ are obtained via

$$\tilde{E}_{ij} = \log{(E_{ij} + 1)}.$$

2. E˜to be the truncated expression matrix. The entries of E¯are obtained by capping the entries of E at the 99.5% quantile.

**[00134]** At block **210,** the optimal transport module 125 determines cell regulatory models based on the optimal transport maps. In certain example embodiments, the optimal transport module 125 determines cell regulatory models based at least in part on the optimal transport maps. In certain example embodiments, the optimal transport module 125 may further identify local biomarker enrichment based at least in part on the optimal transport maps. An example implementation is described in further detail in Example 1 below. Transcription factors (TFs) that appear to play important roles along trajectories to key desitnations are identified by two approaches. The first approach involves constructing a global regulatory model. Pairs of cells at consecutive time points are sampled according to their transport probabilities; expression levels of Tfs in the cell at time $t$ are used to predict expression levels of all non-TFs in the paired cell at time $t + 1$., under the assumption that the regulatory rules are constant across cells and time points. TFs may be excluded from the predicted set to avoid cases of spurious self-regulation). The second approach involves enrichment analysis. TFs are identified based on enrichment in

cells at an earlier time point with a high probability (e.g. >80%) of transitioning to a given state vs. those with a low probability (e.g. <20%).

[00135]     At block 215, the optimal transport module 125 may further define gene modules. In certain example embodiments, this step is optional. Cells may be clustered based on their gene-expression profiles, after performing two rounds of dimensionality reduction to increase statistical power in subsequent analyses. For the reprogramming data disclosed herein, the analysis partitioned 16,339 detected genes into 44 gene modules, which were then analyzed for enrichment of gene sets (signatures) related to specific pathways, cells types, and conditions. **(FIG. 13, Table 1).** Based on the expression profiles in each cell, signature scores were calculated (defined by curated gene sets) for relevant features including MEF identity, pluripotency, proliferation, apoptosis, senescence, X-reactivation, neural identity, placental identity and genomic copy-number variation.

**Table 1**

| Clusters | Gene Modules | ID (Term) | q-Value | Database |
|---|---|---|---|---|
| 1 | GM4 | GO:0036211 (protein modification process) | 7.0 10-3 | BP |
|  | GM10 | GO:001604 (cellular component organization) |  | BP |
|  |  | GO:0036211 (protein modification process) |  | BP |
|  |  | GO:0006325 (chromain organization) |  | BP |
|  |  | GO:0016570 (histone modification) |  | BP |
| 2 | GM5 | GO:0007049 (cell cycle) | 9.6 10-123 | BP |
|  |  | GO:0000278 (mitotic cell cycle) | 6.7 10-110 | BP |
|  |  | GO:0006260 (DNA replication) | 6.7 10-55 | BP |
| 3 | GM33 | IPR001400 (Somatotropin) | 9.0 10-06 | 1 |
|  |  | GO:0005179 (hormone activity) | 3.3 10-09 | MF |
|  |  | R-MMU-1170546 (Prolactin receptor signaling) | 7.0 10-15 | R |
|  |  | R-MMU-982772 (Growth hormone receptor signaling) | 1.1 10-13 | R |
|  | GM40 | GO:0045664 (regulation of neuron differentiation) |  | BP |
| 4 | GM8 | GO:0030855 (epithelial cell differentiation) | 2.6 10-11 | BP |
|  |  | GO:0060429 (epithelium development) | 1.5 10-07 | BP |
|  |  | mmu04530 (Tight junction) | 2.7 10-08 | K |
|  | GM14 | GO:0001890 (placenta development) | 2.5 10-5 | BP |
|  | GM42 | GO:0016126 (sterol biosynthetic process) | 4.8 10-38 | BP |
|  |  | Hallmark cholesterol homeostasis | 8.0 10-29 | M |
| 5 | GM2 | GO:0009653 (anatomical structure morphogenesis) | 5.8 10-29 | BP |
|  |  | GO:0050793 (regulation of developmental process) | 1.6 10-25 | BO |

| | | | | |
|---|---|---|---|---|
| | | GO:0031012 (extracellular matrix) | 1.6 10-17 | CC |
| | GM6 | Lee Bmp2 Targets up | 2.3 10-16 | M |
| | GM7 | GO:0034976 (response to endoplasmic reticulum stress) | 3.8 10-16 | BP |
| | GM9 | GO:0072331 (signal transduction by p53 class mediator) | 6.5 10-06 | BP |
| | | mmu04115 (p53 signaling pathway) | 2.9 10-10 | K |
| | | HALLMARK_P53_PATHWAY | 2.1 10-26 | M |
| | GM23 | GO:0043568 (positive regulation of insulin-like growth factor receptor signaling pathway) | 1.0 10-4 | BP |
| | | GO:0005520 (insulin-like growth factor binding) | 3.1 10-5 | MF |
| | GM27 | GO:0031012 (extracellular matrix) | 2.9 10-3 | CC |
| | GM32 | GO:0006749 (glutathione metabolic process) | 1.5 10-3 | BP |
| | | MOUSE_PWY-4061 (glutathione-mediated detoxification) | 1.7 10-2 | BI |
| | GM34 | GO:0035456 (response to interferon-beta) | 2.5 10-13 | BP |
| | | GO:0006952 (defense response) | 8.0 10-11 | BP |
| | GM35 | GO:0006952 (defense response) | 6.6 10-08 | BP |
| | | GO:0006958 (complement activation, classical pathway) | 1.7 10-5 | BP |
| | GM37 | GO:0034097 (response to cytokine) | 5.0 10-11 | BP |
| | | mmu04668 (TNF signaling pathway) | 4.8 10-11 | K |
| | GM43 | Hallmark Tgf beta signaling | 2.0 10-3 | M |
| | GM44 | GO:0009952 (ranterior/posterior pattern specification) | 2.9 10 15 | BP |
| | | GO:0001501 (skeletal system development) | 1.2 10-12 | BP |
| 6 | GM 13 | Pasini Suzl2 Targets up | 3.0 10-20 | M |
| | | WP1763 PluriNetWork | 3.6 10-06 | W |
| | GM 18 | Mikkelsen Pluripotent State up | 2.2 10-3 | M |
| | GM25 | mouse chrX\|X | 1.1 10-3 | C |
| 7 | GM22 | GO:0007399 (nervous system development) | 4.64 10-5 | BP |
| | | GO:0097458 (neuron part) | 2.4 10-5 | CC |

[00136]    In certain example embodiments, dimensionality reduction may be used to increase robustness. As a first step towards dimensionality reduction, genes that do not show significant variation are removed. The resulting variable-gene expression matrix may be denoted $E_{var}$.

[0100]    A second round of dimensionality reduction may comprise non-linear mapping such as Laplacian embedding, or diffusion component embedding. While principal component analysis (PCA) is a traditional approach to reduce dimensionality, it is only typically appropriate for preserving linear structures. To accommodate nonlinear shapes in high-dimensional gene expression space, diffusion components which are a generalization of principal components were used.

**[0101]** The diffusion components defined in terms of a similarity function k : RG **x** RG → [0, ∞). For a pair (x, y) of G-dimensional gene-expression profiles, the similarity function — or *kernel* function — k(x, y) measures the similarity between x and y. We use the Gaussian kernel function

$$k(x,y) = e^{-\frac{\|\tilde{x} - \tilde{y}\|^2}{2\sigma^2}}.$$

Where *x* and *y* are log-transformed expression profiles (*i.e.* columns of $\tilde{E}^{'}$ )

**[0102]** The diffusion components are defined as the top eigenvectors of a certain matrix constructed by evaluating the kernel function for all pairs of expression profiles $x_i, \ldots, x_N$. Specifically, the kernel matrix $K$ is formed with entries

$$Ki_j = k(x_i, x_j),$$

and then the Laplacian matrix $L$ is formed by multiplying $K$ on the left and the right by $D^{-1/2}$, where $D$ is a diagonal matrix with entries

$$D_{ii} = \sum_{j=1}^{N} k(x_i, x_j).$$

The Laplacian matrix $L$ is given by

$$L = D^{-\frac{1}{2}} K D^{-\frac{1}{2}}.$$

The diffusion components are the eigenvectors $v_i, \ldots, v_N$ of L, sorted by eigenvalue. We embed the data in d dimensional diffusion component space by selecting the top d diffusion components $v_1, \ldots, v_d$, and sending data point $x_i$ to the vector obtained by selecting the ith entry of $v_1, \ldots, v_{20}$. The diffusion component embedding of an expression profile x may be denoted by $\Phi_d(x)$. The top 20 diffusion components were enriched for gene signatures related to biological processes, and therefore were elected to use the top 20 diffusion components to represent data (see below for details).

**[00137]** At block 215, the visualization module 130 generates a visualization of a developmental landscape of the set of cells. To visualize the developmental landscape, the dimensionality of the data is reduced with diffusion components (such as those described above), and then the data is embedded in two dimension with force-directed graph visualization. While alternative visualization methods, such as t-distributed Stochastic Neighbor Embedding (t-SNE), are well suited for identifying clusters, they do not preserve global structures by including

42

repulsive forces between dissimilar points. In particular, these repulsive forces seem to do a good job of splaying out the spikes present in the diffusion map embedding. **FIGs. 7A-7F.**

[0103]     The invention is further described in the following examples, which do not limit the scope of the invention described in the claims.

**Methods for Inducing Pluripotent Stems cell**

[0104]     The invention provides for a method of producing an induced pluripotent stem cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell. In one embodiment, a nucleic acid encoding Obox6 is introduced into a target cell. The method may include a step of introducing into the target cell at least one nucleic acid encoding a reprogramming factor selected from the group consisting of: Oct3/4, Sox2, Soxl, Sox3, Soxl5, Soxl7, Klf4, Klf2, c-Myc, N-Myc, L-Myc, Nanog, Lin28, Fbxl5, ERas, ECAT15-2, Tell, beta-catenin, Lin28b, Sail 1, Sall4, Esrrb, Nr5a2, Tbx3, and Glisl, or selected from the group consisting of: Oct4, Klf4, Sox2 and Myc.

[0105]     In one embodiment, the nucleic acid encoding Obox6 is provided in a recombinant vector, for example, a lentivirus vector. In another embodiment, the nucleic acid encoding the reprogramming factor is provided in a recombinant vector. The nucleic acid may be incorporated into the genome of the cell. The nucleic may not be incorporated into the genome of the cell.

[0106]     The method may include a step of culturing the cells in reprogramming medium as defined herein. The method may also include a step of culturing the cells in the presence of serum or the absence of serum, for example, after a culturing step in reprogramming medium.

[0107]     The induced pluripotent stem cell produced according to the methods of the invention can express at least one of a surface marker selected from the group consisting of: Oct4, SOX2, KLf4, c-MYC, LIN28, Nanog, Glisl , TRA-160/TRA-1-81/TRA-2-54, SSEA1, SSEA4, Sal4 and Esrbb 1.

[0108]     The method can be performed with a target cell that is a mammalian cell, including but not limited to a human, murine, porcine or canine cell. The target cell can be a primary or secondary mouse embryonic fibroblast (MEF).The target cell can be any one of the following: fibroblasts, B cells, T cells, dendritic cells, keratinocytes, adipose cells, epithelial cells, epidermal cells, chondrocytes, cumulus cells, neural cells, glial cells, astrocytes, cardiac cells,

esophageal cells, muscle cells, melanocytes, hematopoietic cells, pancreatic cells, hepatocytes, macrophages, monocytes, mononuclear cells, and gastric cells, including gastric epithelial cells.

[0109]     The target cell can be embryonic, or adult somatic cells, differentiated cells, cells with an intact nuclear membrane, non-dividing cells, quiescent cells, terminally differentiated primary cells, and the like.

[0110]     The invention also provides for a method of producing an induced pluripotent stem cell comprising introducing at least one of Obox6, Spic, Zfp42, Sox2, Mybl2, Msc, Nanog, Hesxl and Esrrb into a target cell to produce an induced pluripotent stem cell. In one embodiment, a nucleic acid encoding Obox6, Spic, Zfp42, Sox2, Mybl2, Msc, Nanog, Hesxl or Esrrb is introduced into a target cell.

[0111]     The invention also provides a method of producing an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5 or Table 6 into a target cell to produce an induced pluripotent stem cell. . In one embodiment, a nucleic acid encoding a transcription factor identified in Table 2, Table 3, Table 4, Table 5 or Table 6 is introduced into a target cell.

**Table 2**

| Genes detected  in less than  1% of cells in clusters 1-27 |
| --- |
| Rhox2a |
| Myolf |
| Xlr3c |
| Stra8 |
| Smtnl l |
| Tspo2 |
| Aurkc |
| Dazl |
| Rhoxl |
| Crxos |
| Rbakdn |
| Smclb |
| Tuba3a |
| Sycp3 |
| Apobec2 |
| Obox6 |
| Patl2 |
| Platr3 |

| |
|---|
| Gpx6 |
| 1700013H16Rik |
| Lncencl |
| Tell |
| Spic |
| Hsf2bp |
| Fkbp6 |
| Arl14epl |
| Pacsinl |
| Faml83b |
| Dpys |
| Fmrlnb |
| Gm9732 |
| Dppa4 |
| Fam25c |
| Dppa2 |
| Lrrc34 |
| Trpm l |
| Khdc3 |
| Col9a2 |
| Magebl6 |
| Hesxl |
| Myl7 |
| Ly6g6e |
| Gm9 |
| Gml3580 |
| Aard |
| Zfp42 |
| Gm7325 |

## Table 3

| TF | frequency in high / frequency in low | frequency in high | frequency in low |
|---|---|---|---|
| Spic | 15.63 | 38.5% | 2.4% |
| Zfp42 | 17.41 | 33.4% | 1.9% |
| Obox6 | 61.90 | 9.3% | 0.1% |
| Sox2 | 11.68 | 33.5% | 2.9% |
| Mybl2 | 22.55 | 17.2% | 0.7% |
| Msc | 20.37 | 16.9% | 0.8% |
| Nanog | 6.08 | 51.3% | 8.4% |

| | | | |
|---|---|---|---|
| **Hesxl** | **8.68** | **35.5%** | **4.1%** |
| **Esrrb** | **17.00** | **16.4%** | **1.0%** |

**Bold: Intersection between global regulatory network and enrichment analysis**

**Table 4**

**Late pluripotency markers unique to successful trajectory**

Genes detected in less than 1% of cells in clusters 1-27
Rhox2a
Myolf
Xlr3c
Stra8
Smtnll
Tspo2
Aurkc
Dazl
Rhoxl
Crxos
Rbakdn
Smclb
Tuba3a
Sycp3
Apobec2
Obox6
Patl2
Platr3
Gpx6
1700013H16Rik
Lncencl
Tell
Spic
Hsf2bp
Fkbp6
ArlHepl
Pacsinl
Faml83b
Dpys
Fmrlnb

Gm9732
Dppa4
Fam25c
Dppa2
Lrrc34
Trpml
Khdc3
Col9a2
Magebl6
Hesxl
Myl7
Ly6g6e
Gm9
Gml3580
Aard
Zfp42
Gm7325

## Table 5

| TF | frequency in high / frequency in low | frequency in high | frequency in low |
|---|---|---|---|
| **Spic** | 15.63 | 38.5% | 2.4% |
| **Zfp42** | 17.41 | 33.4% | 1.9% |
| **Obox6** | 61.90 | 9.3% | 0.1% |
| **Sox2** | 11.68 | 33.5% | 2.9% |
| **Mybl2** | 22.55 | 17.2% | 0.7% |
| **Msc** | 20.37 | 16.9% | 0.8% |
| **Nanog** | 6.08 | 51.3% | 8.4% |
| **Hesxl** | 8.68 | 35.5% | 4.1% |
| **Esrrb** | 17.00 | 16.4% | 1.0% |

**Bold: Intersection between global regulatory network and enrichment analysis**

## Table 6

Candidate Transcription Factors

| Gene | Description | Reference |
|---|---|---|
| Spic | Spi-C transcription factor (Spi-l/PU.l related) | Roderick TH, Chromosomal inversions in studies of mammalian mutagenesis. Genetics. 1979 May;92(l Pt 1 Suppl):sl21-6 |
| Zfp42 | zinc finger protein 42 | Hosier BA, et al., Expression of REX-1, a gene containing zinc finger motifs, is rapidly reduced by retinoic acid in F9 teratocarcinoma cells. Mol Cell Biol. 1989 Dec;9(12):5623-9 |
| Obox6 | oocyte specific homeobox 6 | Ko MS, et al., Large-scale cDNA analysis reveals phased gene expression patterns during preimplantation mouse development. Development. 2000 Apr; 127(8): 1737-49 |
| Sox2 | SRY (sex determining region Y)-box 2 | Lyon MF, et al., Dose-response curves for radiation-induced gene mutations in mouse oocytes and their interpretation. Mutat Res. 1979 Nov;63(l): 161-73 |
| Mybl2 | myeloblastosis oncogene-like 2 | Lam EW, et al., Characterization and cell cycle-regulated expression of mouse B-myb. Oncogene. 1992 Sep;7(9): 1885-90 |
| Msc | musculin | Robb L, et al., musculin: a murine basic helix-loop-helix transcription factor gene expressed in embryonic skeletal muscle. Mech Dev. 1998 Aug;76(l-2): 197-201 |
| Nanog | Nanog homeobox | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Hesxl | homeobox gene expressed in ES cells | Thomas PQ, et al., F£ES-1, a novel homeobox gene expressed by murine embryonic stem cells, identifies a new class of homeobox genes. Nucleic Acids Res. 1992 Nov 11;20(21):5840 |
| Esrrb | estrogen related receptor, beta | Pettersson K, et al., Expression of a novel member of estrogen response element-binding nuclear receptors is restricted to the early stages of chorion formation during mouse embryogenesis. Mech Dev. 1996 Feb;54(2):21 1-23 |
| Rhox2a | reproductive homeobox 2A | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Myolf | myosin IF | Hasson T, et al., Mapping of |

| | | |
|---|---|---|
| | | unconventional myosins in mouse and human. Genomics. 1996 Sep 15;36(3):431-9 |
| Xlr3c | X-linked lymphocyte-regulated 3C | Bergsagel PL, et al., Sequence and expression of murine cDNAs encoding Xlr3a and Xlr3b, defining a new X-linked lymphocyte-regulated Xlr gene subfamily. Gene. 1994 Dec 15;150(2):345-50 |
| Stra8 | stimulated by retinoic acid gene 8 | Bouillet P, et al., Efficient cloning of cDNAs of retinoic acid-responsive genes in P19 embryonal carcinoma cells and characterization of a novel mouse gene, Stra1 (mouse LERK-2/Eplg2). Dev Biol. 1995 Aug;170(2):420-33 |
| Smtnll | smoothelin-like 1 | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Tspo2 | translocator protein 2 | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Aurkc | aurora kinase C | Tseng TC, et al., Protein kinase profile of sperm and eggs: cloning and characterization of two novel testis-specific protein kinases (AIE1, AIE2) related to yeast and fly chromosome segregation regulators. DNA Cell Biol. 1998 Oct; 17(10): 823-33 |
| Dazl | deleted in azoospermia-like | Kasahara M, et al., Genetic mapping of a male germ cell-expressed gene Tpx-2 to mouse chromosome 17. Immunogenetics. 1991;34(2):132-5 |
| Rhoxl | reproductive homeobox 1 | Maclean JA 2nd, et al., Rhox: a new homeobox gene cluster. Cell. 2005 Feb 11;120(3):369-82 |
| Crxos | cone-rod homeobox, opposite strand | Ko MS, et al., Large-scale cDNA analysis reveals phased gene expression patterns during preimplantation mouse development. Development. 2000 Apr; 127(8): 1737-49 |
| Rbakdn | RB-associated KRAB zinc finger downstream neighbor (non-protein coding) | MGD Nomenclature Committee, 2/14/1995; |
| Smclb | structural maintenance of chromosomes IB | Biswas U, et al., Distinct Roles of Meiosis-Specific Cohesin Complexes in Mammalian Spermatogenesis. PLoS |

| | | |
|---|---|---|
| | | Genet. 2016 Oct; 12(1 0):el 0063 89 |
| | | Villasante A, et al., Six mouse alpha-tubulin mRNAs encode five distinct isotypes: testis-specific expression of two sister genes. Mol Cell Biol. 1986 |
| Tuba3a | tubulin, alpha 3A | Jul;6(7):2409-19 |
| | | Roderick TH, Chromosomal inversions in studies of mammalian mutagenesis. |
| | synaptonemal complex protein | Genetics. 1979 May;92(l Pt 1 |
| Sycp3 | 3 | Suppl):sl21-6 |
| | | Hirano K, et al., Targeted disruption of the mouse apobec-1 gene abolishes |
| | apolipoprotein B mRNA | apolipoprotein B mRNA editing and |
| | editing enzyme, catalytic | eliminates apolipoprotein B48. J Biol |
| Apobec2 | polypeptide 2 | Chem. 1996 Apr 26;271(17):9887-90 |
| | | Ko MS, et al., Large-scale cDNA analysis reveals phased gene expression patterns during preimplantation mouse development. Development. 2000 |
| Obox6 | oocyte specific homeobox 6 | Apr; 127(8): 173 7-49 |
| | | Marnef A, et al., Distinct functions of maternal and somatic Patl protein |
| | protein associated with | paralogs. RNA. 2010 Nov; 16(1 1):2094- |
| Patl2 | topoisomerase II homolog 2 | 107 |
| | pluripotency associated | Leo D, et al., Transgenic mouse models for |
| Platr3 | transcript 3 | ADHD. Cell Tissue Res. 2013 May 17 |
| | | Roderick TH, Producing and detecting paracentric chromosomal inversions in |
| Gpx6 | glutathione peroxidase 6 | mice. Mutat Res. 1971 Jan;l l(l):59-69 |
| | | Kawai J, et al., Functional annotation of a |
| | RIKEN cDNA 1700013H16 | full-length mouse cDNA collection. |
| 1700013H16Rik | gene | Nature. 2001 Feb 8;409(6821):685-90 |
| | | Lai KM, et al., Diverse Phenotypes and |
| | long non-coding RNA, | Specific Transcription Patterns in Twenty |
| | embryonic stem cells | Mouse Lines with Ablated LincRNAs. |
| Lncencl | expressed 1 | PLoS One. 201 5; 10(4):e0 125522 |
| | | Narducci MG, et al., The murine Tell oncogene: embryonic and lymphoid cell expression. Oncogene. 1997 Aug |
| Tell | T cell lymphoma breakpoint 1 | 18;15(8):919-26 |
| | | Roderick TH, Chromosomal inversions in studies of mammalian mutagenesis. |
| | Spi-C transcription factor | Genetics. 1979 May;92(l Pt 1 |
| Spic | (Spi-l/PU.l related) | Suppl):sl21-6 |
| Hsf2bp | heat shock transcription factor | Kawai J, et al., Functional annotation of a |

| | | |
|---|---|---|
| | 2 binding protein | full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Fkbp6 | FK506 binding protein 6 | Coss MC, et al., Molecular cloning, DNA sequence analysis, and biochemical characterization of a novel 65-kDa FK506-binding protein (FKBP65). J Biol Chem. 1995 Dec 8;270(49):29336-41 |
| Arll4epl | ADP-ribosylation factor-like 14 effector protein-like | Zambrowicz BP, et al., Wnkl kinase deficiency lowers blood pressure in mice: a gene-trap screen to identify potential targets for therapeutic intervention. Proc Natl Acad Sci U S A. 2003 Nov 25;100(24):14109-14 |
| Pacsinl | protein kinase C and casein kinase substrate in neurons 1 | Plomann M, et al., PACSIN, a brain protein that is upregulated upon differentiation into neuronal cells. Eur J Biochem. 1998 Aug 15;256(1):201-1 1 |
| Faml83b | family with sequence similarity 183, member B | Roderick TH, Chromosomal inversions in studies of mammalian mutagenesis. Genetics. 1979 May;92(l Pt 1 Suppl):sl21-6 |
| Dpys | dihydropyrimidinase | Skarnes WC, et al., A conditional knockout resource for the genome-wide study of mouse gene function. Nature. 201 1 Jun 16;474(7351):337-42 |
| Fmrlnb | fragile X mental retardation 1 neighbor | Skarnes WC, et al., A conditional knockout resource for the genome-wide study of mouse gene function. Nature. 201 1 Jun 16;474(7351):337-42 |
| Gm9732 | predicted gene 9732 | Roderick TH, Using inversions to detect and study recessive lethals and detrimentals in mice, in Utilization of Mammalian Specific Locus Studies in Hazard Evaluation and Estimation of Genetic Risk. 1983:135-67. |
| Dppa4 | developmental pluripotency associated 4 | Ko MS, et al., Large-scale cDNA analysis reveals phased gene expression patterns during preimplantation mouse development. Development. 2000 Apr; 127(8): 173 7-49 |
| Fam25c | family with sequence similarity 25, member C | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Dppa2 | developmental pluripotency associated 2 | Ko MS, et al., Large-scale cDNA analysis reveals phased gene expression patterns |

| | | |
|---|---|---|
| | | during preimplantation mouse development. Development. 2000 Apr; 127(8): 1737-49 |
| Lrrc34 | leucine rich repeat containing 34 | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Trpm1 | transient receptor potential cation channel, subfamily M, member 1 | Dickinson ME, et al., High-throughput discovery of novel developmental phenotypes. Nature. 2016 Sep 14;537(7621):508-514 |
| Khdc3 | KH domain containing 3, subcortical maternal complex member | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Col9a2 | collagen, type IX, alpha 2 | Dickinson ME, et al., High-throughput discovery of novel developmental phenotypes. Nature. 2016 Sep 14;537(7621):508-514 |
| Mageb16 | melanoma antigen family B, 16 | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Hesx1 | homeobox gene expressed in ES cells | Thomas PQ, et al., HES-1, a novel homeobox gene expressed by murine embryonic stem cells, identifies a new class of homeobox genes. Nucleic Acids Res. 1992 Nov 11;20(21):5840 |
| Myl7 | myosin, light polypeptide 7, regulatory | Lowey S, et al., Light chains from fast and slow muscle myosins. Nature. 1971 Nov 12;234(5324):81-5 |
| Ly6g6e | lymphocyte antigen 6 complex, locus G6E | Kawai J, et al., Functional annotation of a full-length mouse cDNA collection. Nature. 2001 Feb 8;409(6821):685-90 |
| Gm9 | predicted gene 9 | The FANTOM Consortium and RIKEN Genome Exploration Research Group and Genome Science Group (Genome Network Project Core Group), The Transcriptional Landscape of the Mammalian Genome. Science. 2005;309(5740): 1559-1563 |
| Gm13580 | predicted gene 13580 | Zambrowicz BP, et al., Wnkl kinase deficiency lowers blood pressure in mice: a gene-trap screen to identify potential targets for therapeutic intervention. Proc Natl Acad Sci U S A. 2003 Nov 25;100(24):14109-14 |
| Aard | alanine and arginine rich domain containing protein | Roderick TH, et al., Nineteen paracentric chromosomal inversions in mice. Genetics. |

| | | 1974 Jan;76(l): 109-17 |
| | | Hosier BA, et al., Expression of REX-1, a gene containing zinc finger motifs, is rapidly reduced by retinoic acid in F9 teratocarcinoma cells. Mol Cell Biol. 1989 |
| Zfp42 | zinc finger protein 42 | Dec;9(12):5623-9 |
| | | Hansen J, et al., A large-scale, gene-driven mutagenesis approach for the functional analysis of the mouse genome. Proc Natl |
| | myomixer, myoblast fusion | Acad Sci U S A. 2003 Aug |
| Gm7325 | factor | 19;100(17):9918-22 |

[0112]    The invention also provides a method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

[0113]    The invention also provides a method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5 or Table 6 into a target cell to produce an induced pluripotent stem cell.

[0114]    The invention also provides a method of increasing the efficiency of reprogramming of a cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

[0115]    The invention also provides a method of increasing the efficiency of reprogramming a cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5 or Table 6 into a target cell to produce an induced pluripotent stem cell.

[0116]

[0117]    The invention also provides for an isolated induced pluripotent stem cell produced by the methods of the invention.

[0118]    The invention also provides a method of treating a subject with a disease comprising administering to the subject a cell produced by differentiation of the induced pluripotent stem cell produced by the methods of the invention.

[0119]    The invention also provides for a composition for producing an induced pluripotent stem cell comprising Obox6 or any of the factors identified in Table 2, Table 3, Table 4, Table 5 or Table 6 in combination with reprogramming media.

**[0120]**    The invention also provides for use of Obox6 or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5 or Table 6 for production of an induced pluripotent stem cell.

Definitions

**[0121]**    As used herein, "pluripotent" as it refers to a "pluripotent stem cell" means a cell with the developmental potential, under different conditions, to differentiate to cell types characteristic of all three germ cell layers, i.e., endoderm (e.g., gut tissue), mesoderm (including blood, muscle, and vessels), and ectoderm (such as skin and nerve). Pluripotent cell as used herein, includes a cell that can form a teratoma which includes tissues or cells of all three embryonic germ layers, or that resemble normal derivatives of all three embryonic germ layers (i.e., ectoderm, mesoderm, and endoderm). A pluripotent cell of the invention also means a cell that can form an embryoid body (EB) and express markers for all three germ layers including but not limited to the following: endoderm markers-AFP, FOXA2, GATA4; mesoderm markers-CD34, CDH2 (N-cadherin), COL2A1, GATA2, HAND1, PECAMI, RUNX1, RUNX2; and Ectoderm markers-ALDHlAl,   COL1A1, NCAM1, PAX6, TUBB3 (Tujl).

**[0122]**    A pluripotent cell of the invention also means a human cell that expresses at least one of the following markers: SSEA3, SSEA4, Tra-1-81, Tra-1-60, Rexl,  Oct4, Nanog, Sox2 as detected using methods known in the art. A pluripotent stem cell of the invention includes a cell that stains positive with alkaline phosphatase or Hoechst Stain.

**[0123]**    In some embodiments, a pluripotent cell is termed an "undifferentiated cell." Accordingly, the terms "pluripotency" or a "pluripotent state" as used herein refer to the developmental potential of a cell that provides the ability of the cell to differentiate into all three embryonic germ layers (endoderm, mesoderm and ectoderm). Those of skill in the art are aware of the embryonic germ layer or lineage that gives rise to a given cell type. A cell in a pluripotent state typically has the potential to divide in vitro for a long period of time, e.g., greater than one year or more than 30 passages.

**[0124]**    As used herein, the term "induced pluripotent stem cells (iPSCs or "iPS cells)" refers to cells having similar properties to those of ES cells. In particular, an "iPSC" or "iPS cell" as used herein, includes an undifferentiated cell which is reprogrammed from somatic cells and have pluripotency and proliferation potency. However, this term is not to be construed as

limiting in any sense, and should be construed to have its broadest meaning. As used herein, the term "pluripotent stem cell", as it refers to the cell produced by the claimed methods is synonymous with the term "iPS".

[0125]    Obox6 and any of the other factors described herein can be used to generate induced pluripotent stem cells from differentiated adult somatic cells. In the preparation of induced pluripotent stem cells by using the factors of the present invention, types of cells to be reprogrammed are not particularly limited, and any kind of cells may be used. For example, matured somatic cells may be used, as well as somatic cells of an embryonic period. Other examples of cells capable of being generated into iPS cells and/or encompassed by the present invention include mammalian cells such as fibroblasts, mouse embryonic fibroblasts, B cells, T cells, dendritic cells, keratinocytes, adipose cells, epithelial cells, epidermal cells, chondrocytes, cumulus cells, neural cells, glial cells, astrocytes, cardiac cells, esophageal cells, muscle cells, melanocytes, hematopoietic cells, pancreatic cells, hepatocytes, macrophages, monocytes, mononuclear cells, and gastric cells, including gastric epithelial cells. The cells can be embryonic, or adult somatic cells, differentiated cells, cells with an intact nuclear membrane, non-dividing cells, quiescent cells, terminally differentiated primary cells, and the like. The pluripotent or multipotent cells of the present invention possess the ability to differentiate into cells that have characteristic attributes and specialized functions, such as hair follicle cells, blood cells, heart cells, eye cells, skin cells, placental cells, pancreatic cells, or nerve cells. In particular, pluripotent cells of the invention can differentiate into multiple cell types including but not limited to: cells derived from the endoderm, mesoderm or ectoderm, including but not limited to cardiac cells, neural cells (for example, astrocytes and oligodendrocytes), hepatic cells (for example, pancreatic islet cells), osteogentic, muscle cells, epithelial cells, chondrocytes, adipocytes, placental cells, dendritic cells and, haematopoietic and retinal pigment epithelial (RPE) cells.

[0126]    Induced pluripotent stem cells may express any number of pluripotent cell markers, including: alkaline phosphatase (AP); ABCG2; stage specific embryonic antigen-1 (SSEA-1); SSEA-3; SSEA-4; TRA-1-60; TRA-1-81; Tra-2-49/6E; ERas/ECAT5, E-cadherin; βIII-tubulin; α-smooth muscle actin (α-SMA); fibroblast growth factor 4 (Fgf4), Cripto, Daxl; zinc finger protein 296 (Zfp296); N-acetyltransferase-1 (Natl); (ES cell associated transcript 1 (ECAT1);

ESG1/DPPA5/ECAT2;    ECAT3;  ECAT6;  ECAT7;  ECAT8;  ECAT9;  ECAT10;  ECAT15-1;
ECAT15-2;  Fthll7;  Sall4;  undifferentiated  embryonic  cell  transcription  factor  (Utfl);  Rexl;
p53;  G3PDH;  telomerase,  including  TERT;  silent  X  chromosome  genes;  Dnmt3a;  Dnmt3b;
TRIM28;  F-box  containing  protein  15  (Fbxl5);  Nanog/ECAT4;  Oct3/4;  Sox2;  Klf4;  c-Myc;
Esrrb;  TDGF1;  GABRB3;  Zfp42,  FoxD3;  GDF3;  CYP25A1;  developmental  pluripotency-
associated  2  (DPPA2);  T-cell  lymphoma  breakpoint  1  (Tell);  DPPA3/Stella;  DPPA4;  other
general  markers  for  pluripotency,  etc.  Other  markers  can  include  Dnmt3L;  Soxl5;  Stat3;  Grb2;
SV40  Large  T  Antigen;  HPV16  E6;  HPV16  E7,    -catenin,  and  Bmil.  Such  cells  can  also  be
characterized  by  the  down-regulation  of  markers  characteristic  of  the  differentiated  cell  from
which  the  iPS  cell  is  induced.  For  example,  iPS  cells  derived  from  fibroblasts  may  be
characterized  by  down-regulation  of  the  fibroblast  cell  marker  Thyl  and/or  up-regulation  of
SSEA-1.  It  is  understood  that  the  present  invention  is  not  limited  to  those  markers  listed  herein,
and  encompasses  markers  such  as  cell  surface  markers,  antigens,  and  other  gene  products
including  ESTs,  RNA  (including  microRNAs  and  antisense  RNA),  DNA  (including  genes  and
cDNAs),  and  portions  thereof.

[0127]    As  used  herein,  "increases  the  efficiency"  as  it  refers  to  the  production  of  induced
pluripotent  stem  cells,  means  an  increase  in  the  number  of  induced  pluripotent  stem  cells  that  are
produced,  for  example  in  the  presence  of  Obox6  or  one  or  more  of  the  factors  identified  in  Table
2,  3,  4,  5  or  6,  as  compared  to  the  number  of  cells  produced  in  the  absence  of  Obox6  or  one  or
more  of  the  factors  identified  in  Table  2,  3,  4,  5  or  6  under  identical  conditions.  An  increase  in
the  number  of  induced  pluripotent  cells  means  an  increase  of  at  least  5%,  for  example,  5%,  10%,
15%,  20%,  25%,  30%,  35%,  40%,  45%,  50%,  55%,  60%,  65%,  70%,  75%,  80%,  85%,  90%,
95%,  or  100%  or  more.  An  increase  also  means  at  least  5-fold  more,  for  example,  5-fold,  -fold,
20-fold,  30-fold,  40-fold,  50-fold,  60-fold,  70-fold,  80-fold,  90-fold,  100-fold,  500-fold,  1000-
fold  or  more.  Increases  the  efficiency  also  means  decreasing  the  time  required  to  produce  an
induced  pluripotent  stem  cell,  for  example  in  the  presence  of  Obox6  or  one  or  more  of  the  factors
identified  in  Table  6,  7,  8,  9  or  10,  as  compared  to  the  number  of  cells  produced  in  the  absence  of
Obox6  or  one  or  more  of  the  factors  identified  in  Table  2,  Table  3,  Table  4,  Table  5  and  Table  6.
In  the  presence  of  Obox6  or  any  one  of  the  factors  identified  in  Table  2,  Table  3,  Table  4,  Table
5  and  Table  6,  an  iPSC  can  be  formed  between  5  and  30  days,  between  5  and  20  days,  between

10 and 20 days, for example 10 days, 11 days, 12 days, 13 days, 14 days, 15 days, 16 days, 17 days, 18 days, 19 days or 20 days after the addition of Obox6 or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5 and Table 6or following induction of expression of Obox6 or or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5 and Table 6.

[0128]    Candidate transcriptional regulators to augment reprogramming efficiency include but are not limited to the transcription regulators presented in Tables 2, 3, 4, 5 and 6.


EXPERIMENTAL METHODS

## 1. Derivation of MEFs

[0129]    Mouse embryonic fibroblasts (MEFs) were derived from E13.5 embryos with a mixed B6;129 background. The cell line used in this study was homozygous for ROSA26-M2rtTA, homozygous for a polycistronic cassette carrying *Pou5fl, Kl/4, Sox2,* and *Myc* at the *Collal* locus (18), and homozygous for an EGFP reporter under the control of the *Pou5fl* promoter. Briefly, MEFs were isolated from E13.5 embryos resulting from timed-matings by removing the head, limbs, and internal organs under a dissecting microscope. The remaining tissue was finely minced using scalpels and dissociated by incubation at 37°C for 10 minutes in trypsin-EDTA (Thermo Fisher Scientific). Dissociated cells were then plated in MEF medium containing DMEM (Thermo Fisher Scientific), supplemented with 10% fetal bovine serum (GE Healthcare Life Sciences), non-essential amino acids (Thermo Fisher Scientific), and GlutaMAX (Thermo Fisher Scientific). MEFs were cultured at 37°C and 4% $CO_2$ and passaged until confluent. All procedures, including maintenance of animals, were performed according to a mouse protocol (2006N000104) approved by the MGH Subcommittee on Research Animal Care.

## 2. Reprogramming assay

[0130]    For the reprogramming assay, 20,000 low passage MEFs (no greater than 3-4 passages from isolation) were seeded in a 6-well plate. These cells were cultured at 37°C and 5% CO2 in reprogramming medium containing KnockOut DMEM (GIBCO), 10% knockout serum replacement (KSR, GIBCO), 10% fetal bovine serum (FBS, GIBCO), 1% GlutaMAX (Invitrogen), 1% nonessential amino acids (NEAA, Invitrogen), 0.055 mM 2-mercaptoethanol (Sigma), 1%, penicillin-streptomycin (Invitrogen) and 1,000 U/ml leukemia inhibitory factor

(LIF, Millipore). Day 0 medium was supplemented with 2 g/mL doxycycline Phase-l(Dox) to induce the polycistronic OKSM expression cassette. Medium was refreshed every other day. At day 8, doxycycline was withdrawn, and cells were transferred to either serum-free 2i medium containing 3 μM CHIR99021, 1 μM PD0325901, and LIF (Phase-2(2i)) (25) or maintained in reprogramming medium (Phase-2(serum)). Fresh medium was added every other day until the final time point on day 16. Oct4-EGFP positive iPSC colonies should start to appear on day 10, indicative of successful reprogramming of the endogenous *Oct4* locus.

## 3. Sample collection

**[0131]**    A total of 66,000 cells were collected from twelve time points over a period of 16 days in two different culture conditions. Single or duplicate samples were collected at day 0 (before and after Dox addition), 2, 4, 6, and 8 in Phase-l(Dox); day 9, 10, 11, 12, 16 in Phase-2(2i); and day 10, 12, 16 in Phase-2(serum). Cells were also collected from established iPSCs cell lines reprogrammed from the same MEFs, maintained either in Phase-2(2i) conditions or in Phase-2(serum) medium. For all time points, selected wells were trypsinized for 5 mins followed by inactivation of trypsin by addition of MEF medium. Cells were subsequently spun down and washed with IX PBS supplemented with 0.1% bovine serum albumin. The cells were then passed through a 40 micron filter to remove cell debris and large clumps. Cell count was determined using Neubauer chamber hemocytometer to a final concentration of 1000 cells/ l.

## 4. Single-cell RNA sequencing

**[0132]**    Single-cell RNA-Seq libraries were generated from each time point using the 10X Genomics Chromium Controller Instrument (10X Genomics, Pleasanton, CA) and Chromium™ Single Cell 3' Reagent Kits v1 (PN-120230, PN-120231, PN-120232) according to manufacturer's instructions. Reverse transcription and sample indexing were performed using the CI000 Touch Thermal cycler with 96-Deep Well Reaction Module. Briefly, the suspended cells were loaded on a Chromium controller Single-Cell Instrument to first generate single-cell Gel Bead-In-Emulsions (GEMs). After breaking the GEMs, the barcoded cDNA was then purified and amplified. The amplified barcoded cDNA was fragmented, Atailed and ligated with adaptors. Finally, PCR amplification was performed to enable sample indexing and enrichment of the 3' RNA-Seq libraries. The final libraries were quantified using Thermo Fisher Qubit dsDNA HS Assay kit (Q32851) and the fragment size distribution of the libraries were

determined using the Agilent 2100 BioAnalyzer High Sensitivity DNA kit (5067-4626). Pooled libraries were then sequenced using Illumina Sequencing By Synthesis (SBS) chemistry.

## 5. Lentivirus Vector Construction and Particle Production

[0133]    To test whether transcription factors (TFs) improve late-stage reprogramming efficiency, lentiviral constructs for the top candidates *Zfp42,* and *Obox6* were generated. cDNA for these factors were ordered from Origene (Zfp42-MG203929, and Obox6-MR215428) were cloned into the FUW Tet-On vector (Addgene, Plasmid #20323) using the Gibson Assembly (NEB, E261 1S). Briefly, the cDNA for each TF was amplified and cloned into the backbone generated by removing Oct4 from the FUW-Teto-Oct4 vector. All vectors were verified by Sanger sequencing analysis. For lentivirus production, FIEK293T cells were plated at a density of $2.6 \times 10^6$ cells/well in a 10cm dish. The cells were transfected with the lentiviral packaging vector and a TF-expressing vector at 70-80% growth confluency using the Fugene FID reagent (Promega E231 1) according to the manufacturer's protocols. At 48 hours after transfection, the viral supernatant was collected, filtered and stored at -80°C for future use.

## 6. Reprogramming efficiency of secondary MEFS together with individual TFs

[0134]    We sought to determine the ability of the candidate TFs to augment reprogramming efficiency in secondary MEFs; the use of secondary MEFs for reprogramming overcomes limitations associated with random lentiviral integration events at variable genomic locations. Briefly, secondary MEFs were plated at a concentration of 20,000 cells per well of a 6-well plate. Cells were infected with virus containing *2fp42, Obox6,* or an empty vector and maintained in reprogramming medium as described above. At day 8 after induction, cells were switched to either Phase-2(2i) or Phase-2(serum). On day 16, reprogramming efficiency was quantified by measuring the levels of the EGFP reporter driven by the endogenous *Oct4* promoter. FACS analyses was performed using the Beckman Coulter CytoFLEX S, and the percentage of Oct4-EGFP+ cells was determined. Triplicates were used to determine average and standard deviation **(FIG. 10B).**

## 7. Reprogramming efficiency of primary MEFS with individual TFs and OKSM

[0135]    In addition to demonstrating the ability of a TF to increase reprogramming efficiency in secondary MEFs, the performance of the TFs were independently tested in primary MEFs. To this end, lentiviral particles were generated from four distinct FUW-Teto vectors, containing

*Oct4, Sox2, Kl/4,* and *Myc,*. MEFs from the background strain *B6.Cg-Gt(ROSA)26Sortml(rtTA\*M2)Jae/J x* B6;129S4-Pou5fltm2Jae/J were infected with these lentiviral particles, together with a lentivirus expressing tetracycline-inducible *Zfp42, Obox6* or no insert. Infected cells were then induced with 2μg/mL doxycycline in ESC reprogramming medium (day 0). At day 8 after induction, cells were switched to either Phase-2(2i) or Phase-2(serum). On day 16, the number of Oct4-EGFP+ colonies were counted using a fluorescence microscope. Triplicates for each condition used to determine average values and standard deviation.

## EXAMPLES

**Example 1**

[0136]    Computing trajectories with optimal transport

[0137]    As noted above, for any pair of time points we compute a transport plan that minimizes the expected cost of redistributing mass, subject to constraints involving a proliferation score (see Appendix 1 for a precise statement of the optimization problem). To compute these transport matrices, we need to specify a cost function, a proliferation function, and numerical values for the regularization parameters.

[0138]    Cost functions: We tried several different cost functions based on squared Euclidean distance in different input spaces. Specifically, for cells with expression profiles x and y, given by two columns of the expression matrix E, we specify a cost function c(x, y)

Expression space
$$c_1(x, y) = \| \chi^- - y^- \|^2$$

100 dimensional diffusion component space
$$c_2(\chi, y) = \| A\phi_{100}(x) - A\phi_{1Qo}(y) \|^2$$

20 dimensional diffusion component space
$$c_3(x, y) = \| A\phi_{20}(\chi) - A\Phi_{20}(y) \|^2$$

[0139]    The bar above $x^-$, $y^\sim$ denotes that we apply the truncation transform from section 2, and $\Phi_d$ is the Laplacian embedding from section 3. Note that $\Phi_d$ has the log transform $x \rightarrow \tilde{x}$ built-in. In the equations above, A is a diagonal matrix containing the eigenvalues of the Laplacian matrix, raised to the power 8. Hence $c_2$ and $c_3$ are both truncated versions of the diffusion distance $D_4(x, y)$ from (S5).

**[0140]**     The cost function c3 was used to report the numerical values in the main text, and we computed separate transport maps for 2i and serum. Note that all the cost functions cl, c2, c3 give largely similar results.

**[0141]**     Proliferation function: We estimate the relative growth rate for every cell using the proliferation signature displayed in **FIG.** 7D in the main text. To transform the proliferation score into an estimate of the growth rate (in doublings per day), we first observed that the proliferation score is bimodally distributed over the dataset. We transformed the proliferation score so that the two modes were mapped to a growth ratio of 2.5 per day (this means that over 1 day, a cell in the more proliferative group is expected to produce 2.5 times as many offspring as a cell in the non-proliferative group). However, note that we allow for some laxity in the prescribed growth rate (see supplemental figure on input vs implied proliferation).

**[0142]**     Regularization parameters: We employed the following strategy to select the regularization pa- rameters $\lambda$ and $\varepsilon$. The entropy parameter $\varepsilon$ controls the entropy of the transport map. An extremely large entropy parameter will give a maximally entropic transport map, and an extremely small entropy parameter will give a nearly deterministic transport map (but could also lead to numerical instability in the algorithm). We adjusted the entropy parameter until each cell transitions to between 10 and 50 percent of cells in the next time point, as measured by the Shannon diversity of the rows of the transport map.

**[0143]**     The regularization parameter $\lambda$ controls the fidelity of the constraints: as $\lambda$ gets larger, the constraints become more stringent. We selected $\lambda$ so that the marginals of the transport map are 95% correlated with the prescribed proliferation score.

**[0144]**     Implementation: The scaling algorithm for unbalanced transport (S2) was implemented to compute optimal transport maps. This algorithm performs gradient ascent steps on the dual optimization problem. Because of the entropic regularization, these gradient ascent steps can be performed via diagonal matrix scalings. We implemented versions of the solver in both R and Python.

**[0145]**     Experiments: Computational experiments were performed to evaluate the stability of our results to choice of cost function, regularization parameters, and subsampling the dataset.

[0146]    The cluster-to-cluster origin were compared and fate tables for the different cost functions listed above, and consistent results were found. Moreover, the transport probabilities described above are all robust to choice of cost function.

[0147]    A bootstrap analysis was performed on a batch of 100 subsamples consisting of 50% of the data from each time point. The variance in the cluster-to-cluster origin and fate tables is extremely small (see Table 7).

## Table 7

| MEF.identity | Pluripotency | GI.S | G2.M | Cell.cycle | ER.stress | Epithelial.identity | ECM.rearrangement | Apoptosis | SASP | Neural.identity | Placental.identity | X.reactivation |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gm5571 | Rhox5 | Cdca7 | Cbx5 | Mcm4 | Nck2 | Cdhl | Sulfl | Ercc5 | 116 | Vtn | 4933433pl4rik | Gm21950 |
| Rbfox2 | Tdgfl | Mcm4 | Aurkb | Smc4 | Ankzf1 | Tgml | Coll9al | Serpinb5 | 117 | Ednrb | Esxl | Gm21364 |
| Btbdl9 | Utfl | Mcm2 | Cksl b | Gtse1 | Dnajb2 | Cldn3 | Col3al | Inhbb | Ilia | Sox21 | Afapl | Gml4346 |
| Actnl | Mkrnl | Rfc2 | Cks2 | Ttk | Rhbdl | Cldn4 | Col5a2 | SteaP3 | IIIb | Zeb2 | Zfyve21 | Gml4345 |
| Gatad2a | Dppa5a | Ung | Hn1 | Rangap1 | Bcl2 | Cldn7 | Fnl | Btg2 | 1113 | Hes5 | Erv3 | Gml4351 |
| Med6 | Uppl | Mcm6 | Hmgb2 | Ccnb2 | Ubxn4 | Cldnll | Ihh | Phlda3 | 1115 | Fabp7 | Atgl2 | Gm3701 |
| Mex3a | Chchd10 | Rrm1 | Anp32e | Cenpa | Yodl | Ocln | Col4a4 | Tnni1 | Cxcl15 | Soxl | Lasll | Gm3706 |
| Ccdc80 | Klf2 | Slbp | Lbr | Cenpe | Ppplrl5b | Epcam | Col4a3 | Rgsl6 | Cxcl1 | Neurod l | Rbpl | Gml4347 |
| Mex3c | Trapla | Pena | Tmpo | Cdca8 | Faml29a | Crb3 | Serpinb5 | ier5 | Cxcl2 | Pax3 | Prl2bl | Gml0921 |
| Sdpr | Mylpf | Atad2 | Top2a | CkaP2 | Edem3 | Krt8 | Fmod | Slcl9a2 | Cxcl3 | Pax6 | Prl3dl | Gml0922 |
| Pcdhb2 | 1700013H16Rik | Tipin | Tacc3 | Rad51 | Atf6 | Krtl9 | Elf3 | Adck3 | | Cdh2 | Rnf2 | Gm3750 |
| Triml6 | AA467197 | Mcm5 | Tub b4b | Pena | Ufcl | Pkp3 | Lamcl | Ephxl | Ccl8 | Sox9 | Set | Gm3763 |
| Obsll | Dhxl6 | Uhrf1 | Ncapd2 | Ube2c | Atf3 | Dsp | Tnr | Ptpn14 | Ccl3 | Sox2 | Mrgprg | Mycs |
| Ephal | Mt2 | Rpa2 | Rangap1 | Lbr | Manlbl | Pkpl | Dpt | Atf3 | Ccl20 | Id2 | Aa763515 | Gml4374 |
| Stxlb | Ube2a | Dtl | Cdk1 | CenPf | Toria | | Ddr2 | Notehi | Cell6 | Hoxbl | Tfpi | Nudtll |

| Staul | Khdc3 | Priml | Smc4 | Birc5 | Hspa5 | | Olfml2b | Rxra | Ccl26 | Msxl | Etosl | AU022751 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Serpinel | Pycard | Fen1 | Kif2Ob | Dtl | Dab2ip | | Tgfb2 | Ralgds | Csf2 | Msil | Slc5a6 | NudtlO |
| Aa881470 | Hsp90aal | Hells | Cdca8 | Dscc1 | Nfe2l2 | | Itga8 | Akl | Csf3 | Msi2 | 1600025ml7rik | Bmpl5 |
| Coll2a1 | Prrcl | Gmnn | CkaP2 | Cbx5 | Dnajc10 | | Adamtsl2 | Stom | Ifng | Atohl | Gm9 | Shroom4 |
| 2010300fl7rik | Hatl | Pold3 | Ndc80 | Usp1 | Psmc3 | | Col5al | Ddb2 | Mif | Rbfox3 | Creb3l2 | Dgkk |
| CcdclO2a | Calcoco2 | NasP | Digap5 | Hmmr | Creb311 | | Pomtl | Cd82 | Areg | Map2 | Bbx | Ccnb3 |
| Nradd | Impa2 | Chaflb | Hjurp | Wdr76 | Thbs1 | | Eng | Ilia | Ereg | Tubb3 | Prl3cl | Akap4 |
| Pard6g | Saa3 | Gins2 | CkaP5 | Ung | Eif2ak4 | | Lmxlb | Pcna | Nrg1 | | Mta3 | Clcn5 |
| Ntn4 | Ooep | Pola1 | Bub1 | Hnl | Chac1 | | Gsn | Bmp2 | Egf | | Prl2al | Usp27x |
| 5730471hl9rik | Bnip3 | Msh2 | Ckap2l | Cks2 | Pdia3 | | Olfml2a | Trib3 | Fgf2 | | Gm9112 | Ppplr3f |
| Sepnl | Mtl | Casp8aP2 | Ect2 | Kif2Ob | Bcl2l11 | | Creb3ll | Procr | Hgf | | Afapll2 | Ppplr3fos |
| Pegl2 | Asns | Cdc6 | Kifl1 | Cdk1 | Ddrgkl | | Hsdl7bl2 | BlcaP | Fgf7 | | Erlin2 | Foxp3 |
| Dpysl3 | Aldoa | Ubr7 | Birc5 | Slbp | Tmx4 | | Wtl | Ada | Vegfa | | Pard3 | Ccdc22 |
| 1110012d08rik | Tdh | Ccne2 | Cdca2 | Aurkb | Trib3 | | Greml | Fgfl3 | Ang | | Aifll | Cacnalf |
| Aktl | Gjb3 | Wdr76 | Nuf2 | Kifl1 | H13 | | Spintl | Irak1 | Kitl | | Dmrtcla | Syp |
| Zfp286 | Rbpms2 | Tyms | Cdca3 | Cksl b | Edem2 | | Cst3 | Tspy12 | Cxcl12 | | 4932442l08rik | Gml4703 |
| Ubap2l | Prpsl | Cdc45 | Nusapl | Blm | Cebpb | | Fkbpla | Satl | Pigf | | Gjb2 | Prickle3 |
| Samd4 | Fam25c | Clspn | Ttk | Msh2 | Ptpn1 | | Mmp9 | Zmat3 | IgfbP2 | | Gjb5 | Plp2 |
| Phc2 | Eif2s2 | Rrm2 | Aurka | Gas213 | Vapb | | Sulf2 | Hspa4l | IgfbP3 | | Slco5al | Magix |
| Mcam | Cenpm | Dscc1 | Mki67 | Tyms | Srpx | | Atp7a | Slc7all | Igfbp4 | | Wdr61 | Gpkow |
| Pla2g4c | Nanog | Rad51 | Fam64a | HjurP | Aifm1 | | Noxl | Tm4sfl | Igfbp6 | | Kitl | Wdr45 |
| Fzd7 | Ndufa412 | Usp1 | Ccnb2 | Hells | Ubqln2 | | Col4a6 | Rap2b | IgfbP7 | | 9430027b09rik | RP23-109E24.10 |

63

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pappa | Syce2 | Exo1 | Tpx2 | Priml | Mbtps2 | | Prdx4 | Fbxw7 | Mm Pi | | Tfrc | Praf2 |
| Ptk7 | Gml3251 | Blm | Hjurp | Uhrf1 | Uspl3 | | Gpm6b | S100a4 | Mm P3 | | Slc6a2 | Ccdcl20 |
| Nuakl | Taf7 | Rad51a Pi | Anln | Ndc80 | Ufml | | Egfl6 | S10Oal0 | Mm plO | | Wdr45 | Tfe3 |
| Ill7rd | Nudt4 | Mlflip | Kif2c | Mcm6 | Serp1 | | Postn | TxniP | Mmpl2 | | Zxda | Gripap1 |
| Ptk2 | Cox5a | E2f8 | Cenpe | Rrm1 | Creb314 | | Rxfpl | Nhlh2 | Mmpl3 | | Prdx4 | Kcndl |
| Ehd2 | Sod2 | Brip1 | Gtsel | Mlflip | Tmem67 | | Sfrp2 | Dnttip2 | Mmpl4 | | Faml22b | Otud5 |
| Lats2 | SlOOa13 | | Kif23 | Top2a | Ufll | | Hapln2 | Clca2 | TimP2 | | Zxdb | Pim2 |
| Hspg2 | Fkbp6 | | Cdc20 | Hmgb2 | Ube2jl | | Ctss | WwPi | Serpine1 | | Zxdc | Slc35a2 |
| 493045 6gl4rik | Rhox9 | | Ube2c | Cence2 | Vcp | | Adamtsl4 | Klf4 | Serpinb2 | | Pip5kla | Pqbpl |
| 493029 b21rik | Gdf3 | | CenPf | G2e3 | Creb3 | | St7l | Ikbkap | Plat | | Placl | Timml7b |
| Rps20 | 2700094K13Rik | | Cenpa | Tmp0 | Sec6lb | | Colllal | Cdkn2a | Plau | | Igf2as | Gml0491 |
| Vgll3 | Fmrlnb | | Hmmr | Nusapl | Erp44 | | Npnt | Cdkn2b | Ctsb | | Usp9x | Gml0490 |
| Prrl5 | Hmgn2 | | Ctcf | Ncapd2 | AI314180 | | Cyr61 | Jun | leaml | | Psg28 | Pcskln |
| Fbxl7 | Ubald2 | | Psrcl | Mcm2 | Jun | | B4galtl | Slc35dl | leam3 | | Bmp8b | Eras |
| Maged2 | Lactb2 | | Cdc25c | Kif2c | Casp9 | | Reck | Plk3 | Tnfrsfllb | | Fnl | Hdac6 |
| Galntl4 | Folrl | | Nek2 | Cdca2 | Fbxo6 | | Tgfbrl | Rnfl9b | Tnfrsfla | | Psg23 | Gatal |
| Pdgfc | Gm7325 | | Gas213 | NasP | Fbxo2 | | Col27al | Sfn | Tnfrsflb | | Bmp8a | Glod5 |
| Tmtc4 | Agtrap | | G2e3 | Gmnn | Ube4b | | P3hl | Fuca1 | TnfrsflOb | | Psg21 | Gml4820 |
| Tmtc3 | Sppl | | | Cdc6 | Ube2j2 | | Hspg2 | Epha2 | Fas | | Dusp9 | Suv39h1 |
| Lpar4 | Hells | | | Pold3 | Psmc2 | | Vwal | Wrap73 | Plaur | | H19 | Was |
| Pcdhl9 | Dppa4 | | | Ckap2l | Tmubl | | Dnajb6 | Mxd4 | Il6st | | Tmem37 | Wdrl3 |
| Eda2r | Gabarapl2 | | | Fam64a | Tmeml2 | | Emilinl | Rchy1 | Egfr | | Mmpl5 | Rbm3 |

| | | | | | 9 | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Pcdhl8 | Rhox6 | | | Ubr7 | Wfsl | | Mpvl7 | Iscu | Fnl | | FamlOlb | Rbm3os |
| Gprl76 | Rhoxl | | | Fen1 | Ube2k | | Apbb2 | TriaPi | | | Phfl6 | Tbcld25 |
| LoclOO503471 | Cdc5l | | | Bub1 | Tbl2 | | Pdgfra | Prkabl | | | 4930422n03rik | Ebp |
| Mical2 | Texl9.1 | | | Brip1 | Get4 | | Ambn | Trafdl | | | Ada | Porcn |
| Dzipll | Trim28 | | | Atad2 | Bhlhal5 | | Dmpl | Pom121 | | | Mmpla | Ftsjl |
| Hoxc6 | Atp5gl | | | Psrc1 | Creb312 | | Ibsp | Pdgfa | | | Gprl26 | Slc38a5 |
| Hoxc5 | Sox2 | | | Rrm2 | Pdia4 | | Tfipll | Gadd45a | | | Arf2 | SsxblO |
| Mettl4-psl | Jam2 | | | Tipin | Eif2ak3 | | Eln | Vamp8 | | | Tinagll | Ssxb9 |
| Sec63 | Fkbp3 | | | Casp8aP2 | Rnfl03 | | Plod3 | Retsat | | | Mfi2 | Ssxbl |
| Ikbip | Cox7b | | | Tubb4b | Aupl | | Colla2 | Tprkb | | | Rpn2 | Ssxb2 |
| Tsc22d2 | Ash2l | | | Kif23 | Itprl | | Ndnf | Tgfa | | | Abhd2 | Gml4459 |
| 231007 6g05rik | Dut | | | Exo1 | Edem l | | Vhl | Mxd1 | | | Hrctl | Ssxb6 |
| Anxa6 | Dtymk | | | Rfc2 | Bbc3 | | Mfap5 | Sec6lal | | | Adm | Ssxb3 |
| Nfatc4 | Gpx4 | | | Pola1 | Psmc4 | | Ercc2 | Xpc | | | Abhd6 | Ssxb8 |
| Fnl | Eif4ebPi | | | Mki67 | Bax | | Bcl3 | Ccnd2 | | | Slc7al | Ssx9 |
| Wnt9a | Morel | | | Tpx2 | Pppl rl5a | | Tgfbl | H2afj | | | Tead4 | Ssxb5 |
| Sorcs2 | Fabp3 | | | Aurka | Vimp | | Mia | Ldhb | | | Mbnl3 | Gm6592 |
| Tmeffl | Zfp428 | | | Anln | Rnfl21 | | Spint2 | LrmP | | | Gprl | Gm5751 |
| C79491 | Aqp3 | | | Chafl b | Anks4b | | Aplpl | Tm7sf3 | | | 2900057el5rik | B630019K06Rik |
| Crlfl | Grhpr | | | HjurP | Ern2 | | Hpn | Tgfb1 | | | Ldocl | Fthll7b |
| 2610034e01rik | Higdla | | | Tacc3 | Atp2al | | Klk4 | Sertad3 | | | Adaml9 | Fthll7c |
| Gjd4 | Rpp25 | | | Mc | Brsk2 | | Acan | Ceb | | | Rybp | Fthll7 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | m5 | | | | pa | | | | d |
| Ccngl | Rbpms | | | Anp32e | Ins2 | | Serpinhl | Klk8 | | | Col4al | Fthll7e |
| Gprl24 | Mmp3 | | | DlgaP5 | Ccnd1 | | Apbbl | Bax | | | Fndc3c1 | Fthll7f |
| Fibin | Apobec3 | | | Ect2 | Map3k5 | | Ilk | Ppplrl5a | | | Col4a2 | 4930402K13Rik |
| 8030476ll9rik | Spc24 | | | Nuf2 | Nrbf2 | | Ric8 | Rpll8 | | | 4930502el8rik | Lancl3 |
| Ddr2 | Xlr3a | | | Cdc45 | Derl3 | | Muc5ac | Aen | | | Pkn2 | Gml4862 |
| Arf4 | Recll4 | | | CkaP5 | Ube2g2 | | Ctgf | Rrp8 | | | Rlim | Xk |
| Ptprs | Mtf2 | | | Ctcf | Tmem259 | | Nr2el | Ccp110 | | | 160001SilOrik | 1700012L04Rik |
| Sprr2k | Snrpn | | | Clspn | Creb313 | | Nepn | Nuprl | | | Afp | Gml4501 |
| Adm | Gml3580 | | | Cdca7 | Hsp9Obi | | P4hal | Ptpre | | | Tmeml40 | Cybb |
| A830029e22rik | Gmnn | | | Cdca3 | Apaf1 | | Spock2 | Hras | | | Fstl3 | Gm5132 |
| 9230114kl4rik | Chmp4c | | | Rpa2 | Ifng | | Adamtsl4 | Eps812 | | | Ing4 | Dynlt3 |
| Extl3 | Hsf2bp | | | Gins2 | Os9 | | Mmpll | Ctsd | | | Taf7l | Hypm |
| Mecom | Polr2e | | | E2f8 | Ddit3 | | Coll8al | Cd81 | | | Sultlel | 4930557A04Rik |
| Qsoxl | Blvrb | | | Cdc25c | Erlin2 | | Myf5 | Perp | | | Olrl | Sytl5 |
| Teadl | Ldhb | | | Nek2 | Ppp2cb | | Col4al | Rps12 | | | 2610019f03rik | Srpx |
| Snx7 | Apocl | | | Cdc20 | Ubxn8 | | Csgalnact1 | Tpd5211 | | | Fll | Rpgr |
| Cdkl4 | Syngrl | | | Rad51aPi | Casp3 | | Comp | Sesn1 | | | Fbxw8 | Otc |
| Cdkn2a | Bexl | | | | Pik3r2 | | Gfod2 | Foxo3 | | | Sema4c | Tspan7 |
| Cdkn2b | Nr2c2aP | | | | Amfr | | Has3 | Ddit4 | | | Ctnnbip1 | Gml0489 |
| Ccnyll | | | | | Herpudl | | Atxnll | Zfp365 | | | Tfpi2 | Midlip1 |
| Tubb2a-ps2 | | | | | Aars | | Crispld2 | Prmt2 | | | ZbtblO | Gml4493 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Aen | | | | | Selk | | Foxfl | Mknk2 | | | Mitf | Gm14483 |
| Farpl | | | | | Eroll | | Foxc2 | Draml | | | Gpr50 | Gm14474 |
| 4930402h24rik | | | | | Psmc6 | | Agt | Apaf1 | | | Hic2 | Gm14477 |
| Sh3rf3 | | | | | Trim13 | | Exoc8 | Btgl | | | Tpbpb | Gm14476 |
| Adaml9 | | | | | Dnajc3 | | Eroll | Mdm2 | | | Slc9a6 | Gm14484 |
| Ddbl | | | | | Casp4 | | Lgals3 | Ddit3 | | | Prl7dl | Gm14479 |
| Cttn | | | | | Casp12 | | Ripk3 | Gls2 | | | Tpbpa | Gm14482 |
| 9230112e08rik | | | | | ScamP5 | | Loxl2 | Dgka | | | Slco2al | Gm14478 |
| Dbnl | | | | | Pml | | Lcpl | Cdkn2aiP | | | Pkp2 | Gm14475 |
| Fyttdl | | | | | Parp16 | | Mmpl3 | Hmoxl | | | 9630050el6rik | Gm4906 |
| Lrrcl5 | | | | | Nckl | | Mmp20 | Rrad | | | Pvrl2 | Bcor |
| Fkbpl0 | | | | | Uba5 | | Col5a3 | Cdh13 | | | Zfp568 | Gm14635 |
| Trubl | | | | | Uspl9 | | Smarca4 | Osginl | | | Vtcnl | Atp6ap2 |
| Zdhhc20 | | | | | Stt3b | | Aplp2 | Cgrrfl | | | Il6ra | 1810030O07Rik |
| Stonl | | | | | Rnfl85 | | Mpzl3 | Abhd4 | | | Foxo4 | Med14 |
| Hoxdl3 | | | | | Xbpl | | Thsd4 | Kifl3b | | | Hsp90b1 | Usp9x |
| Nudt6 | | | | | Erlec1 | | Anxa2 | Rbl | | | Prl7cl | 2010308F09Rik |
| Hoxdl2 | | | | | Stc2 | | Myole | Nudtl5 | | | Prl6al | Ddx3x |
| Prss23 | | | | | Trp53 | | Nphp3 | Tsc22dl | | | Cdh5 | Nyx |
| 9430030nl7rik | | | | | Aloxl5 | | Dagl | Casp1 | | | Fgd6 | Cask |
| Arntl2 | | | | | Derl2 | | Lamb2 | Stl4 | | | Cysltr2 | Gpr34 |
| Sh3rfl | | | | | Trim25 | | Kif9 | Ei24 | | | Rhox6 | Gpr82 |
| Mrc2 | | | | | Cdk5rap3 | | Sh3pxd2b | Vwa5a | | | Cdh3 | Gm5382 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Mdhl | | | | | Cede 47 | | Adamts2 | Zbtb 16 | | | Spp2 | Gml45 05 |
| Rictor | | | | | Psmc 5 | | Wnt3a | Rps 271 | | | Ziml | Drrl |
| Map4k 5 | | | | | Ernl | | Mfap4 | Map kapk 3 | | | Flnb | Cyptl |
| Plcll | | | | | Nplo c4 | | Serpinf2 | Ip6k 2 | | | Rbbp7 | Maoa |
| Septll | | | | | P4hb | | Vtn | Tcn2 | | | Map3k7 | Maob |
| Ryk | | | | | Txnd c5 | | Nfl | Lif | | | Rhox9 | Ndp |
| Tgfb3 | | | | | Faf2 | | Collal | Upp 1 | | | Whscll 1 | Efhc2 |
| Ube2i | | | | | Ubql nl | | Ramp2 | Ceng 1 | | | Slc38al | Fundcl |
| Tgfb2 | | | | | Atgl 0 | | Gfap | Cyfi P2 | | | 160001 2pl7rik | Dusp2 1 |
| Zfp319 | | | | | Thbs 4 | | Sox9 | Gnb 211 | | | Adra2b | Kdm6a |
| GmlO 399 | | | | | Col4a 3bp | | Erollb | Hint 1 | | | Pgf | 493057 8C19Ri k |
| Fbxol 7 | | | | | Pik3r 1 | | Nidi | Gm2 a | | | 120000 9i06rik | Gm266 52 |
| Wnt5a | | | | | Pdia6 | | Foxf2 | Hist 3h2 a | | | Mfsd7c | BC049 702 |
| Criml | | | | | Dnaj b9 | | Foxcl | Alox 8 | | | Esam | Chst7 |
| Midi | | | | | Tmxl | | Ripkl | Trp5 3 | | | Gprl07 | Slc9a7 |
| Displ | | | | | Jkam P | | Tfap2a | Taxi bp3 | | | Au0157 91 | Rp2 |
| Ubox5 | | | | | Selll | | Ecm2 | Traf 4 | | | Arhgap 8 | Jade3 |
| St7l | | | | | Psmc 1 | | B4galt7 | Cdk 5rl | | | Ankrdl 7 | Rgn |
| Col5a2 | | | | | Atxn 3 | | Tgfbi | Ppm l d | | | Cul7 | Ndufbl 1 |
| Axl | | | | | Derll | | Pxdn | Rad 51c | | | 231006 7p03rik | RbmlO |
| Col5al | | | | | Rnfl 39 | | Smocl | Tob 1 | | | Irs3 | Ubal |
| Zyx | | | | | Foxre d2 | | Ltbp2 | Krtl 7 | | | Prl5al | Cdkl6 |
| Ror2 | | | | | Pla2g 6 | | Flrt2 | Hexi m l | | | Fntb | Uspll |
| Wdfy3 | | | | | Atf4 | | Fbln5 | Fdxr | | | Tceanc | Araf |
| Amotl | | | | | Ep30 | | Egflam | Itgb | | | Lepr | Synl |

|  |  |  |  |  |  |  |  |  |  |  |  |  |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 |  |  |  |  | 0 |  |  | 4 |  |  |  |  |
| **Yapl** |  |  |  |  | Tmbim6 |  | Tnfrsfllb | Sphk1 |  |  | Tnfrsf9 | Timp1 |
| Phldb2 |  |  |  |  | Txndell |  | Col14a1 | Rhbdf2 |  |  | Papola | Cfp |
| 6330562c20rik |  |  |  |  | Sdf2l1 |  | Has2 | BaiaP2 |  |  | Srd5al | Elk1 |
| Ctnnd1 |  |  |  |  | Ufdll |  | Ptk2 | Dcxr |  |  | Clqtnfl | Uxt |
| Rock2 |  |  |  |  | Eif2b5 |  | Sex | Histlhlc |  |  | Slc38a4 | Zfp182 |
| Maspl |  |  |  |  | Nrros |  | Fblnl | Ninj1 |  |  | Angpt4 | Spaca5 |
| Pvtl |  |  |  |  | Pdia5 |  | Adamts20 | Nol8 |  |  | Ctla2a | Zfp300 |
| Tnc |  |  |  |  | Gsk3b |  | Col2al | F2r |  |  | 9930012kllrik | Ssxa1 |
| Fbln2 |  |  |  |  | Park2 |  | Myhll | Ankra2 |  |  | Mical3 | Gm21876 |
| Hdlbp |  |  |  |  | Stub1 |  | Ccdc80 | Plk2 |  |  | Apoa4 | 4930453H23Rik |
| AtplOa |  |  |  |  | Pdia2 |  | Abi3bp | Sdcl |  |  | Cul4b | Gm6938 |
| Loxll |  |  |  |  | Crebrf |  | App | Gpx2 |  |  | 3632454l22rik | Gm26593 |
| Loxl2 |  |  |  |  | Bakl |  | Seracl | $Zfp3_{611}$ |  |  | Psg-psl | Agtr2 |
| Fbln5 |  |  |  |  | Rnf5 |  | Pig | Fos |  |  | Lcor | Slc6a14 |
| Ctgf |  |  |  |  | Atf6b |  | Smoc2 | Ccnk |  |  | Tnfrsf22 | Gm28269 |
| Efnb2 |  |  |  |  | Bag6 |  | Hasl | Jag2 |  |  | Tnfrsf23 | Gm28268 |
| Rxra |  |  |  |  | Flotl |  | Noxol | **Ndrgl** |  |  | Sosl | Klhl13 |
| Ccnd2 |  |  |  |  | Eif2ak2 |  | Collla2 | Pmml |  |  | Dlx3 | Wdr44 |
| Gpc2 |  |  |  |  | PmaiPi |  | Tnxb | Plxnb2 |  |  | Ippk | Gm4907 |
| Ntf3 |  |  |  |  | Tmx3 |  | Tnf | Vdr |  |  | Htr2b | Gm4985 |
| Kif5b |  |  |  |  | Syvn1 |  | 2300002M23Rik | CsrnP2 |  |  | Duspl6 | Gm27192 |
| Slit2 |  |  |  |  | Erlin1 |  | Flotl | Acvrlb |  |  | Cdc73 | Gm5934 |
| Tpml |  |  |  |  |  |  | Hsp90ab1 | Spl |  |  | 1700025g04rik | Gm4297 |
| Gpc4 |  |  |  |  |  |  | Washl | Abat |  |  | Prl4al | Gm593 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | 5 |
| Flnb | | | | | | | Vit | Socs1 | | | Zfp655 | Gm5169 |
| 4930555bllrik | | | | | | | Cyplbl | Abcc5 | | | Slcl3a4 | Gml993 |
| Fine | | | | | | | Fshr | Trp63 | | | Ceacam14 | E330010L02Rik |
| C76332 | | | | | | | Mkx | Fam162a | | | Ceacam15 | Gm5168 |
| Capn2 | | | | | | | Lox | App | | | Trapla | Gm2012 |
| Phlda3 | | | | | | | Hpse2 | Rab40c | | | Ceacam12 | Gm2030 |
| Map3k7 | | | | | | | Kazaldl | Bak1 | | | Gml6515 | Six |
| Myhl0 | | | | | | | Nfkb2 | Def6 | | | Ceacam13 | Gml4525 |
| D18ertd653e | | | | | | | | Cdknla | | | 4930447f24rik | Gm6121 |
| Stox2 | | | | | | | | Tap1 | | | Gzmd | Gml0230 |
| Igf2r | | | | | | | | Ier3 | | | Foxj2 | Gm2101 |
| D15ertd621e | | | | | | | | Polh | | | Fbxll9 | GmlOO58 |
| Arid5b | | | | | | | | Ccnd3 | | | Gzmc | Gm2117 |
| Tnfrsfl0b | | | | | | | | Hbegf | | | Gzmf | Gm4836 |
| 26100lle03rik | | | | | | | | Hdac3 | | | Gzme | GmlOl47 |
| Ckap4 | | | | | | | | Rad9a | | | Gzmg | Gm2165 |
| Efna2 | | | | | | | | Ctsf | | | Patl2 | GmlOO96 |
| Picalm | | | | | | | | Slc3a2 | | | 3830417al3rik | Gm2200 |
| Cdhl0 | | | | | | | | Fas | | | Tspanl4 | Gm26818 |
| Ddahl | | | | | | | | | | | Handl | Gm3669 |
| Uba3 | | | | | | | | | | | AtxnlO | Gml0488 |
| 0610038b21rik | | | | | | | | | | | Mgat4a | E330016L19Rik |

70

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Gemin7 | | | | | | | | | | | | Unc50 | Gm14632 |
| Uba1 | | | | | | | | | | | | Il2rb | Gm7437 |
| Fbn1 | | | | | | | | | | | | Ceacam11 | Gm14974 |
| Lhx9 | | | | | | | | | | | | Plekhg1 | Gm10487 |
| Eif4g2 | | | | | | | | | | | | Prl3b1 | Gm21447 |
| Vcl | | | | | | | | | | | | Folr1 | Spin2f |
| Bcl2l2 | | | | | | | | | | | | A830080Odolrik | Gm2784 |
| Cd276 | | | | | | | | | | | | Blzf1 | Gm2777 |
| Lrrc58 | | | | | | | | | | | | Zfp667 | Gm21883 |
| Wwc2 | | | | | | | | | | | | Flt1 | Spin2e |
| Lpp | | | | | | | | | | | | Usp27x | Gm21608 |
| Aril | | | | | | | | | | | | Hdac4 | Gm21637 |
| Ltbp1 | | | | | | | | | | | | Itgb3 | Gm21645 |
| Ltbp2 | | | | | | | | | | | | Sri | Gm27999 |
| Wisp1 | | | | | | | | | | | | Sema3f | Gmcl1l |
| Igf1r | | | | | | | | | | | | Prl3a1 | Gm5926 |
| Rhobtb3 | | | | | | | | | | | | Bahd1 | Gm21951 |
| Fam198b | | | | | | | | | | | | Sin3b | Gm21657 |
| Cnn2 | | | | | | | | | | | | Gm2a | Gm21789 |
| Glipr2 | | | | | | | | | | | | Serpinb9g | Gm2825 |
| Sydel | | | | | | | | | | | | Bend4 | Spin2-ps6 |
| Hhat | | | | | | | | | | | | Bend5 | Gm2863 |
| Zmat3 | | | | | | | | | | | | Serpinb9b | Gm2854 |
| Cald1 | | | | | | | | | | | | Serpinb9c | Gm2913 |
| Pmepa1 | | | | | | | | | | | | Serpinb9d | Gm2927 |
| E130112l23rik | | | | | | | | | | | | Plekhh1 | Gm2933 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Bag2 | | | | | | | | | | | 221001lc24rik | Gm2964 |
| Zfp583 | | | | | | | | | | | Cd320 | Gm21870 |
| Pibfl | | | | | | | | | | | Ccnjl | Gm21681 |
| Pmaip1 | | | | | | | | | | | Entpd2 | Spin2g |
| A130022jl5rik | | | | | | | | | | | Illr2 | Gm21699 |
| Bcl9l | | | | | | | | | | | Sfmbt2 | Gm14552 |
| Cpa6 | | | | | | | | | | | 170001lm02rik | Gm10486 |
| D13ertd787e | | | | | | | | | | | Plekha7 | Gm2309 |
| Pabpc4 1 | | | | | | | | | | | Sfrp5 | Gm14553 |
| Zfhx3 | | | | | | | | | | | Ppplr3f | Gm14819 |
| Itga5 | | | | | | | | | | | Obsll | Dock11 |
| Txnrdl | | | | | | | | | | | Slc23a3 | Il13ra1 |
| Htrlb | | | | | | | | | | | Tmem87b | Zcchc12 |
| Hmga2 | | | | | | | | | | | Epasl | Lonrf3 |
| Sept2 | | | | | | | | | | | Ccdc68 | Gm6268 |
| Lambl | | | | | | | | | | | Kdelr2 | Gm14569 |
| Zfp518b | | | | | | | | | | | Pramef12 | Pgrmc1 |
| Parva | | | | | | | | | | | Lrp8 | Akap17b |
| Gulpl | | | | | | | | | | | Pard6b | Slc25a43 |
| Shank1 | | | | | | | | | | | PeglO | Slc25a5 |
| Bmpl | | | | | | | | | | | N4bp2 | Gm14549 |
| Aktlsl | | | | | | | | | | | Pla2g4e | 2310010G23Rik |
| Itga9 | | | | | | | | | | | Fam78b | C330007P06Rik |
| Abccl | | | | | | | | | | | Arrdc3 | Ube2a |
| Eda | | | | | | | | | | | Pla2g4d | Nkrf |
| B4galt | | | | | | | | | | | Rassf8 | Gm150 |

72

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2 | | | | | | | | | | | | 08 |
| Nidi | | | | | | | | | | | Au0158 36 | Sept6 |
| Ncaml | | | | | | | | | | | Csnkle | Sowah d |
| Shc2 | | | | | | | | | | | Stagl | Rpl39 |
| Uba6 | | | | | | | | | | | Vnnl | Upf3b |
| Tradd | | | | | | | | | | | Tchhll | Nkap |
| Rtell | | | | | | | | | | | Plala | Akapl 4 |
| Bicd2 | | | | | | | | | | | Slc45a4 | Ndufal |
| Adamt sl2 | | | | | | | | | | | Tex264 | Rnfll3 al |
| Hs2stl | | | | | | | | | | | Pcdhl2 | Gm9 |
| DlOert d610e | | | | | | | | | | | Ctr9 | Rhoxl |
| Cyr61 | | | | | | | | | | | Ccrlll | Rhox2a |
| Gtf3cl | | | | | | | | | | | Htatsfl | Rhox3a |
| Lbh | | | | | | | | | | | 903040 9gllrik | Rhox4a |
| Krt33b | | | | | | | | | | | Tspan9 | Rhox3a 2 |
| Gm66 07 | | | | | | | | | | | Rassf6 | Rhox4a 2 |
| D3wsu 167e | | | | | | | | | | | 463140 2f24rik | Rhox2 b |
| Zc3h7 b | | | | | | | | | | | A2m | Rhox4 b |
| 76304 03g23r ik | | | | | | | | | | | Rimklb | Rhox2c |
| Tnpo2 | | | | | | | | | | | LoclOO 504569 | Rhox3c |
| Cepl7 0 | | | | | | | | | | | Apob | Rhox4c |
| Pdlim5 | | | | | | | | | | | Tmeml 50a | Rhox2 d |
| Pdlim7 | | | | | | | | | | | 913040 4d08rik | Rhox4 d |
| Cad | | | | | | | | | | | Prl8a6 | Rhox2e |
| Unc5b | | | | | | | | | | | Cts6 | Rhox3e |
| 24100 18ll3ri k | | | | | | | | | | | Prl8a8 | Rhox4e |
| LoclOO 21634 3 | | | | | | | | | | | Prl8a9 | Rhox2f |
| Glrx3 | | | | | | | | | | | Cts3 | Rhox3f |

| Kctd5 | | | | | | | | | | Krtl8 | Rhox4f |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Loc269 472 | | | | | | | | | | Nrnll | Rhox3g |
| Myolc | | | | | | | | | | Sfil | Rhox2g |
| 49305 62cl5r ik | | | | | | | | | | Tlr5 | Rhox4g |
| Till | | | | | | | | | | Rhou | Rhox3 h |
| Sema3 a | | | | | | | | | | Arhgef6 | Rhox2 h |
| Itgbl | | | | | | | | | | Tmeml 85b | Rhox5 |
| Nxn | | | | | | | | | | Tram2 | Rhox6 |
| Tmem 41b | | | | | | | | | | Citedl | Rhox7a |
| Sec23a | | | | | | | | | | Cited2 | Rhox8 |
| Gm22 | | | | | | | | | | Zfand2a | Rhox7 b |
| Itgb5 | | | | | | | | | | Krt25 | Rhox9 |
| Dysf | | | | | | | | | | Klk4 | Btgl- psl |
| Thbsl | | | | | | | | | | Tnfrsfl l b | Btgl- ps2 |
| Bc022 687 | | | | | | | | | | 201020 4kl3rik | RhoxlO |
| Dnm3 os | | | | | | | | | | Torlaip 2 | Rhoxll |
| Rnd3 | | | | | | | | | | Fmrlnb | Rhoxl2 |
| Pik3c2 a | | | | | | | | | | Ctsr | Rhoxl3 |
| 28100 08m24 rik | | | | | | | | | | Ctsq | Zbtb33 |
| Spred3 | | | | | | | | | | Prl8a2 | Tmem 255a |
| Senp5 | | | | | | | | | | Ctsm | Atplb4 |
| Arll3b | | | | | | | | | | Prl8al | Lamp2 |
| Polr2e | | | | | | | | | | Ctsj | Gm759 8 |
| Itgav | | | | | | | | | | Mpzll | Cul4b |
| Igf2bp 3 | | | | | | | | | | Stra6 | Mctsl |
| | | | | | | | | | | Bcap31 | Clgalt lcl |
| | | | | | | | | | | Cregl | Gml45 65 |
| | | | | | | | | | | Tcfap2c | 603049 |

| | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | | 8E09Rik |
| | | | | | | | | | | | | | Prl7bl | Cyptl5 |
| | | | | | | | | | | | | | Ghrh | Cyptl4 |
| | | | | | | | | | | | | | 4930486l24rik | Gria3 |
| | | | | | | | | | | | | | Neurog2 | Thoc2 |
| | | | | | | | | | | | | | 5430425jl2rik | Xiap |
| | | | | | | | | | | | | | Prl7al | Stag2 |
| | | | | | | | | | | | | | Prl7a2 | Gm43337 |
| | | | | | | | | | | | | | Mirll99 | Sh2dla |
| | | | | | | | | | | | | | Tbcldl0a | Tenml |
| | | | | | | | | | | | | | Ralbpl | Gm362 |
| | | | | | | | | | | | | | Pdgfra | Dcafl212 |
| | | | | | | | | | | | | | Morc4 | Dcafl211 |
| | | | | | | | | | | | | | Rarres2 | Prr32 |
| | | | | | | | | | | | | | Arid3a | 4930515L19Rik |
| | | | | | | | | | | | | | Lifr | Actrtl |
| | | | | | | | | | | | | | Shisa3 | Gm29242 |
| | | | | | | | | | | | | | Uevld | Smarca1 |
| | | | | | | | | | | | | | Scnnlb | Ocrl |
| | | | | | | | | | | | | | Dnajbl2 | Apln |
| | | | | | | | | | | | | | Brwd3 | Xpnpep2 |
| | | | | | | | | | | | | | Hhipll | Sash3 |
| | | | | | | | | | | | | | Fbln7 | Zdhhc9 |
| | | | | | | | | | | | | | Maspl | Utpl4a |
| | | | | | | | | | | | | | Nrk | 9530027J09Rik |
| | | | | | | | | | | | | | Pvr | Bcorll |
| | | | | | | | | | | | | | Atp2cl | Elf4 |
| | | | | | | | | | | | | | Amot | Aifml |
| | | | | | | | | | | | | | 1600014k23rik | Rab33a |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | Tbrgl | Zfp280c |
| | | | | | | | | | | | Slitl | Slc25a14 |
| | | | | | | | | | | | A730090h04rik | Gprll9 |
| | | | | | | | | | | | 4931406pl6rik | Rbmx2 |
| | | | | | | | | | | | Opn3 | Gm595 |
| | | | | | | | | | | | Pdia4 | Enox2 |
| | | | | | | | | | | | B930054o08 | Gml4696 |
| | | | | | | | | | | | 1700031f05rik | Gml4697 |
| | | | | | | | | | | | Inhba | Arhgap36 |
| | | | | | | | | | | | Inhbb | Olfrl320 |
| | | | | | | | | | | | Helz | Olfrl321 |
| | | | | | | | | | | | Sele | Igsfl |
| | | | | | | | | | | | Pdia6 | Olfrl322 |
| | | | | | | | | | | | Pdia5 | Olfrl323 |
| | | | | | | | | | | | Creb3 | Olfrl324 |
| | | | | | | | | | | | Efnal | Stk26 |
| | | | | | | | | | | | Dlg5 | Frmd7 |
| | | | | | | | | | | | Procr | Rap2c |
| | | | | | | | | | | | Fgfrl | Mbnl3 |
| | | | | | | | | | | | Gnb4 | Hs6st2 |
| | | | | | | | | | | | 2310030g06rik | Usp26 |
| | | | | | | | | | | | Gcml | 1700080O16Rik |
| | | | | | | | | | | | Psgl8 | Gpc4 |
| | | | | | | | | | | | Goltlb | Gpc3 |
| | | | | | | | | | | | Psgl9 | Gml4582 |
| | | | | | | | | | | | Psgl6 | A630012P03Rik |
| | | | | | | | | | | | Slc2al | Ccdcl60 |
| | | | | | | | | | | | Psgl7 | Phf6 |

| | | | | | | | | | | | Htra3 | Hprt |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | Klhll3 | Gm287 30 |
| | | | | | | | | | | | Ets2 | Placl |
| | | | | | | | | | | | Nppc | Faml2 2b |
| | | | | | | | | | | | Tgml | Faml2 2c |
| | | | | | | | | | | | Tmeml 08 | Mospd 1 |
| | | | | | | | | | | | Usp53 | Etd |
| | | | | | | | | | | | Mark3 | Gml45 97 |
| | | | | | | | | | | | Cbx8 | Cxxlc |
| | | | | | | | | | | | Hspa5 | Cxxla |
| | | | | | | | | | | | Spats2 | Cxxlb |
| | | | | | | | | | | | Limk2 | 493050 2E18Ri k |
| | | | | | | | | | | | Mkl2 | 170001 3H16Ri k |
| | | | | | | | | | | | Shroom 4 | Zfp36l3 |
| | | | | | | | | | | | Shroom 1 | Xlr |
| | | | | | | | | | | | Pou2f3 | Gml64 05 |
| | | | | | | | | | | | Acvr2b | Gml64 30 |
| | | | | | | | | | | | Rbms2 | Slxll |
| | | | | | | | | | | | Atg4b | 383040 3N18Ri k |
| | | | | | | | | | | | Pappa2 | Gm773 |
| | | | | | | | | | | | Rbm25 | 160002 5M 17R ik |
| | | | | | | | | | | | Gm479 3 | Zfp449 |
| | | | | | | | | | | | Nidi | Gm215 5 |
| | | | | | | | | | | | Uba6 | Smiml 012a |
| | | | | | | | | | | | Lamcl | Gm217 4 |
| | | | | | | | | | | | Slc40al | Ddx26 b |

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Hapln3 | Gml04 77 |
| | | | | | | | | | | | | Faml76 a | Gm648 |
| | | | | | | | | | | | | Pdliml | Mmgtl |
| | | | | | | | | | | | | Ube2q2 | Slc9a6 |
| | | | | | | | | | | | | Au0180 91 | Fhll |
| | | | | | | | | | | | | Bdkrb2 | Mtap7 d3 |
| | | | | | | | | | | | | E13020 3bl4rik | Adgrg4 |
| | | | | | | | | | | | | SlOOg | Brs3 |
| | | | | | | | | | | | | 493340 2el3rik | Htatsfl |
| | | | | | | | | | | | | Dapk2 | Vglll |
| | | | | | | | | | | | | Gmll9 85 | Gml47 18 |
| | | | | | | | | | | | | Fndc3b | Cd40lg |
| | | | | | | | | | | | | Twsgl | Arhgef 6 |
| | | | | | | | | | | | | Aldhla3 | Rbmx |
| | | | | | | | | | | | | Lnx2 | Gm364 |
| | | | | | | | | | | | | Taf7 | GprlOl |
| | | | | | | | | | | | | Ai84486 9 | Zic3 |
| | | | | | | | | | | | | Clecl2b | 493055 0L24Ri k |
| | | | | | | | | | | | | Prkcsh | Fgfl3 |
| | | | | | | | | | | | | Lama5 | F9 |
| | | | | | | | | | | | | Tchh | Mcf2 |
| | | | | | | | | | | | | Lamal | Atpllc |
| | | | | | | | | | | | | Rps6ka 6 | Gm707 3 |
| | | | | | | | | | | | | Vhl | Gml46 61 |
| | | | | | | | | | | | | Eps8l2 | Sox3 |
| | | | | | | | | | | | | Polg | Gml46 62 |
| | | | | | | | | | | | | | Gml46 64 |
| | | | | | | | | | | | | | Cdrl |
| | | | | | | | | | | | | | Ldocl |
| | | | | | | | | | | | | | 493340 2E13Ri k |

| | | | | | | | | | | | | 4931400O07Rik |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | 1700019B21Rik |
| | | | | | | | | | | | | Gm6760 |
| | | | | | | | | | | | | 3830417A13Rik |
| | | | | | | | | | | | | Slitrk4 |
| | | | | | | | | | | | | Ctag2 |
| | | | | | | | | | | | | 4930447F04Rik |
| | | | | | | | | | | | | Slitrk2 |
| | | | | | | | | | | | | 1700036O09Rik |
| | | | | | | | | | | | | Gmll40 |
| | | | | | | | | | | | | Gml4692 |
| | | | | | | | | | | | | 4933436l01Rik |
| | | | | | | | | | | | | Fmrlos |
| | | | | | | | | | | | | Fmrl |
| | | | | | | | | | | | | Fmrlnb |
| | | | | | | | | | | | | Gml4698 |
| | | | | | | | | | | | | Gm6812 |
| | | | | | | | | | | | | Gml4705 |
| | | | | | | | | | | | | Aff2 |
| | | | | | | | | | | | | 1700111N16Rik |
| | | | | | | | | | | | | 1700020N15Rik |
| | | | | | | | | | | | | Ids |
| | | | | | | | | | | | | 1110012L19Rik |
| | | | | | | | | | | | | 4930567H17Ri |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | k |
| | | | | | | | | | | | | BC023 829 |
| | | | | | | | | | | | | Mamld 1 |
| | | | | | | | | | | | | Mtml |
| | | | | | | | | | | | | Mtmrl |
| | | | | | | | | | | | | Cd99l2 |
| | | | | | | | | | | | | Gml61 89 |
| | | | | | | | | | | | | Hmgb3 |
| | | | | | | | | | | | | Gpr50 |
| | | | | | | | | | | | | Vma21 |
| | | | | | | | | | | | | Gmll4 1 |
| | | | | | | | | | | | | Prrg3 |
| | | | | | | | | | | | | Fatel |
| | | | | | | | | | | | | Cnga2 |
| | | | | | | | | | | | | Magea 4 |
| | | | | | | | | | | | | Gabre |
| | | | | | | | | | | | | Magea 10 |
| | | | | | | | | | | | | Gabra3 |
| | | | | | | | | | | | | Gabrq |
| | | | | | | | | | | | | Cetn2 |
| | | | | | | | | | | | | Nsdhl |
| | | | | | | | | | | | | Gml46 84 |
| | | | | | | | | | | | | Zfpl85 |
| | | | | | | | | | | | | Pnma5 |
| | | | | | | | | | | | | Pnma3 |
| | | | | | | | | | | | | Xlr4a |
| | | | | | | | | | | | | Xlr3a |
| | | | | | | | | | | | | Xlr5a |
| | | | | | | | | | | | | Gml46 85 |
| | | | | | | | | | | | | DXBay 18 |
| | | | | | | | | | | | | Xlr5b |
| | | | | | | | | | | | | Spin2d |
| | | | | | | | | | | | | Xlr3b |
| | | | | | | | | | | | | Xlr4b |
| | | | | | | | | | | | | F8a |

| | | | | | | | | | | | | Xlr4c |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Xlr3c |
| | | | | | | | | | | | | Xlr5c |
| | | | | | | | | | | | | RP23-95K12.13 |
| | | | | | | | | | | | | Zfp275 |
| | | | | | | | | | | | | Gml8336 |
| | | | | | | | | | | | | Gm26726 |
| | | | | | | | | | | | | Zfp92 |
| | | | | | | | | | | | | Trex2 |
| | | | | | | | | | | | | Haus7 |
| | | | | | | | | | | | | Bgn |
| | | | | | | | | | | | | Atp2b3 |
| | | | | | | | | | | | | Dusp9 |
| | | | | | | | | | | | | Pnck |
| | | | | | | | | | | | | Slc6a8 |
| | | | | | | | | | | | | Bcap31 |
| | | | | | | | | | | | | Abcdl |
| | | | | | | | | | | | | Plxnb3 |
| | | | | | | | | | | | | Srpk3 |
| | | | | | | | | | | | | Idh3g |
| | | | | | | | | | | | | Ssr4 |
| | | | | | | | | | | | | Pdzd4 |
| | | | | | | | | | | | | Llcam |
| | | | | | | | | | | | | Arhgap4 |
| | | | | | | | | | | | | Avpr2 |
| | | | | | | | | | | | | NaalO |
| | | | | | | | | | | | | Renbp |
| | | | | | | | | | | | | Hcfcl |
| | | | | | | | | | | | | Iraki |
| | | | | | | | | | | | | Mecp2 |
| | | | | | | | | | | | | Opnlmw |
| | | | | | | | | | | | | Tex28 |
| | | | | | | | | | | | | Tktll |
| | | | | | | | | | | | | Flna |
| | | | | | | | | | | | | Emd |
| | | | | | | | | | | | | RpllO |
| | | | | | | | | | | | | Dnasel |

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | 11 |
| | | | | | | | | | | | | | Taz |
| | | | | | | | | | | | | | Atp6ap1 |
| | | | | | | | | | | | | | Gdil |
| | | | | | | | | | | | | | Fam50a |
| | | | | | | | | | | | | | Plxna3 |
| | | | | | | | | | | | | | Lage3 |
| | | | | | | | | | | | | | Ubl4a |
| | | | | | | | | | | | | | SlclOa3 |
| | | | | | | | | | | | | | Fam3a |
| | | | | | | | | | | | | | Ikbkg |
| | | | | | | | | | | | | | G6pdx |
| | | | | | | | | | | | | | Gm6880 |
| | | | | | | | | | | | | | Olfrl326-psl |
| | | | | | | | | | | | | | Olfrl325 |
| | | | | | | | | | | | | | Gm5640 |
| | | | | | | | | | | | | | Gm6890 |
| | | | | | | | | | | | | | Gm5936 |
| | | | | | | | | | | | | | Gab3 |
| | | | | | | | | | | | | | Dkcl |
| | | | | | | | | | | | | | Mppl |
| | | | | | | | | | | | | | Smim9 |
| | | | | | | | | | | | | | F8 |
| | | | | | | | | | | | | | Fundc2 |
| | | | | | | | | | | | | | Cmc4 |
| | | | | | | | | | | | | | Mtcpl |
| | | | | | | | | | | | | | Brcc3 |
| | | | | | | | | | | | | | Vbpl |
| | | | | | | | | | | | | | Gml5384 |
| | | | | | | | | | | | | | Rab39b |
| | | | | | | | | | | | | | Gml5063 |
| | | | | | | | | | | | | | Pls3 |
| | | | | | | | | | | | | | Gml47 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | 15 |
| | | | | | | | | | | | | Gml47 07 |
| | | | | | | | | | | | | Gml47 17 |
| | | | | | | | | | | | | Cldn34 b3 |
| | | | | | | | | | | | | Cldn34 b4 |
| | | | | | | | | | | | | Cldn34 d |
| | | | | | | | | | | | | Tbllx |
| | | | | | | | | | | | | Prkx |
| | | | | | | | | | | | | Gml47 42 |
| | | | | | | | | | | | | Pbsn |
| | | | | | | | | | | | | Gml47 44 |
| | | | | | | | | | | | | 543040 2E10Ri k |
| | | | | | | | | | | | | Obpla |
| | | | | | | | | | | | | Gm593 8 |
| | | | | | | | | | | | | Obplb |
| | | | | | | | | | | | | Gml47 43 |
| | | | | | | | | | | | | 493048 OEllRi k |
| | | | | | | | | | | | | Prrgl |
| | | | | | | | | | | | | Fam47 c |
| | | | | | | | | | | | | Gm717 3 |
| | | | | | | | | | | | | Mageb 16 |
| | | | | | | | | | | | | Gm267 75 |
| | | | | | | | | | | | | Tmem 47 |
| | | | | | | | | | | | | 493059 5M18R ik |
| | | | | | | | | | | | | Dmd |
| | | | | | | | | | | | | Tsga8 |
| | | | | | | | | | | | | Fthll7a |
| | | | | | | | | | | | | Tab3 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Gk |
| | | | | | | | | | | | | Gml4764 |
| | | | | | | | | | | | | Gml4762 |
| | | | | | | | | | | | | 5430427019Rik |
| | | | | | | | | | | | | Samt3 |
| | | | | | | | | | | | | NrObl |
| | | | | | | | | | | | | Mageb4 |
| | | | | | | | | | | | | Illrapl1 |
| | | | | | | | | | | | | Gm27000 |
| | | | | | | | | | | | | Pet2 |
| | | | | | | | | | | | | 4932429P05Rik |
| | | | | | | | | | | | | 4930415L06Rik |
| | | | | | | | | | | | | Gm44 |
| | | | | | | | | | | | | Gml4773 |
| | | | | | | | | | | | | Mageb2 |
| | | | | | | | | | | | | Gm5072 |
| | | | | | | | | | | | | Gm8914 |
| | | | | | | | | | | | | 1700084M14Rik |
| | | | | | | | | | | | | Gml4781 |
| | | | | | | | | | | | | Mageb5 |
| | | | | | | | | | | | | Mageb1 |
| | | | | | | | | | | | | Mageb18 |
| | | | | | | | | | | | | Gm5941 |
| | | | | | | | | | | | | 1700003E24Rik |
| | | | | | | | | | | | | BC061 |

| | | | | | | | | | | | | | 195 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Arx |
| | | | | | | | | | | | | | Polal |
| | | | | | | | | | | | | | Pcytlb |
| | | | | | | | | | | | | | Pdk3 |
| | | | | | | | | | | | | | AU015 836 |
| | | | | | | | | | | | | | Gml47 98 |
| | | | | | | | | | | | | | Zfx |
| | | | | | | | | | | | | | Eif2s3x |
| | | | | | | | | | | | | | Klhll5 |
| | | | | | | | | | | | | | Fam90 alb |
| | | | | | | | | | | | | | Apoo |
| | | | | | | | | | | | | | Gml48 27 |
| | | | | | | | | | | | | | Maged 1 |
| | | | | | | | | | | | | | Gspt2 |
| | | | | | | | | | | | | | Zxdb |
| | | | | | | | | | | | | | RP23- 9K14.6 |
| | | | | | | | | | | | | | Gm266 17 |
| | | | | | | | | | | | | | Spin4 |
| | | | | | | | | | | | | | Arhgef 9 |
| | | | | | | | | | | | | | Amerl |
| | | | | | | | | | | | | | Asbl2 |
| | | | | | | | | | | | | | Zc4h2 |
| | | | | | | | | | | | | | Zc3hl2 b |
| | | | | | | | | | | | | | 170001 ODOIRi k |
| | | | | | | | | | | | | | Lasll |
| | | | | | | | | | | | | | Msn |
| | | | | | | | | | | | | | F63002 8O10Ri k |
| | | | | | | | | | | | | | Vsig4 |
| | | | | | | | | | | | | | Hsf3 |
| | | | | | | | | | | | | | Heph |
| | | | | | | | | | | | | | Gprl65 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Pgrl5l |
| | | | | | | | | | | | | Eda2r |
| | | | | | | | | | | | | Ar |
| | | | | | | | | | | | | Ophnl |
| | | | | | | | | | | | | Yipf6 |
| | | | | | | | | | | | | Stard8 |
| | | | | | | | | | | | | Efnbl |
| | | | | | | | | | | | | Gml4812 |
| | | | | | | | | | | | | Gml4809 |
| | | | | | | | | | | | | Gml4808 |
| | | | | | | | | | | | | Pjal |
| | | | | | | | | | | | | Tmem28 |
| | | | | | | | | | | | | Eda |
| | | | | | | | | | | | | Awat2 |
| | | | | | | | | | | | | Otud6a |
| | | | | | | | | | | | | Igbpl |
| | | | | | | | | | | | | Dgat2l6 |
| | | | | | | | | | | | | Awatl |
| | | | | | | | | | | | | P2ry4 |
| | | | | | | | | | | | | Arr3 |
| | | | | | | | | | | | | Pdzdll |
| | | | | | | | | | | | | Kif4 |
| | | | | | | | | | | | | Gdpd2 |
| | | | | | | | | | | | | Gml4902 |
| | | | | | | | | | | | | Dlg3 |
| | | | | | | | | | | | | Texll |
| | | | | | | | | | | | | Slc7a3 |
| | | | | | | | | | | | | Snxl2 |
| | | | | | | | | | | | | Foxo4 |
| | | | | | | | | | | | | Gm614 |
| | | | | | | | | | | | | Gm20489 |
| | | | | | | | | | | | | H2rg |
| | | | | | | | | | | | | Medl2 |
| | | | | | | | | | | | | Nlgn3 |
| | | | | | | | | | | | | Gjbl |
| | | | | | | | | | | | | Zmym3 |

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Nono |
| | | | | | | | | | | | | | Itgblb P2 |
| | | | | | | | | | | | | | Tafl |
| | | | | | | | | | | | | | Ogt |
| | | | | | | | | | | | | | Cxcr3 |
| | | | | | | | | | | | | | Gm477 9 |
| | | | | | | | | | | | | | 803047 4K03Ri k |
| | | | | | | | | | | | | | Nhsl2 |
| | | | | | | | | | | | | | Rgag4 |
| | | | | | | | | | | | | | Pin4 |
| | | | | | | | | | | | | | Ercc6l |
| | | | | | | | | | | | | | Rps4x |
| | | | | | | | | | | | | | Citedl |
| | | | | | | | | | | | | | Hdac8 |
| | | | | | | | | | | | | | Phkal |
| | | | | | | | | | | | | | Gm911 2 |
| | | | | | | | | | | | | | Dmrtcl b |
| | | | | | | | | | | | | | Dmrtcl cl |
| | | | | | | | | | | | | | Dmrtcl c2 |
| | | | | | | | | | | | | | 170003 lF05Ri k |
| | | | | | | | | | | | | | Dmrtcl a |
| | | | | | | | | | | | | | 170001 1M02R ik |
| | | | | | | | | | | | | | Napll2 |
| | | | | | | | | | | | | | Cdx4 |
| | | | | | | | | | | | | | Chicl |
| | | | | | | | | | | | | | Gm269 52 |
| | | | | | | | | | | | | | Tsx |
| | | | | | | | | | | | | | Gm269 92 |
| | | | | | | | | | | | | | Tsix |
| | | | | | | | | | | | | | Xist |
| | | | | | | | | | | | | | Jpx |

| | | | | | | | | | | | | | Ftx |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Zcchcl3 |
| | | | | | | | | | | | | | Slcl6a2 |
| | | | | | | | | | | | | | Rlim |
| | | | | | | | | | | | | | C77370 |
| | | | | | | | | | | | | | Abcb7 |
| | | | | | | | | | | | | | Uprt |
| | | | | | | | | | | | | | Zdhhcl5 |
| | | | | | | | | | | | | | 1700121L16Rik |
| | | | | | | | | | | | | | Magee2 |
| | | | | | | | | | | | | | Pbdcl |
| | | | | | | | | | | | | | Magee1 |
| | | | | | | | | | | | | | 5330434G04Rik |
| | | | | | | | | | | | | | Cypt2 |
| | | | | | | | | | | | | | Fgfl6 |
| | | | | | | | | | | | | | Atrx |
| | | | | | | | | | | | | | Magtl |
| | | | | | | | | | | | | | Cox7b |
| | | | | | | | | | | | | | Atp7a |
| | | | | | | | | | | | | | Tlrl3 |
| | | | | | | | | | | | | | Pgkl |
| | | | | | | | | | | | | | Taf9b |
| | | | | | | | | | | | | | Fnd3c2 |
| | | | | | | | | | | | | | Fndc3c1 |
| | | | | | | | | | | | | | Cysltrl |
| | | | | | | | | | | | | | Gm5127 |
| | | | | | | | | | | | | | Zcchc5 |
| | | | | | | | | | | | | | Lpar4 |
| | | | | | | | | | | | | | P2rylO |
| | | | | | | | | | | | | | A630033H20Rik |
| | | | | | | | | | | | | | Gprl74 |

| | | | | | | | | | | | | Itm2a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Tbx22 |
| | | | | | | | | | | | | 261000 2M06R ik |
| | | | | | | | | | | | | Fam46 d |
| | | | | | | | | | | | | Gm732 |
| | | | | | | | | | | | | Gm379 |
| | | | | | | | | | | | | Brwd3 |
| | | | | | | | | | | | | Hmgn5 |
| | | | | | | | | | | | | Sh3bgr 1 |
| | | | | | | | | | | | | Gm637 7 |
| | | | | | | | | | | | | RP23- 240M8 .2 |
| | | | | | | | | | | | | Pou3f4 |
| | | | | | | | | | | | | Cylcl |
| | | | | | | | | | | | | GmlOl 12 |
| | | | | | | | | | | | | Rps6ka 6 |
| | | | | | | | | | | | | Hdx |
| | | | | | | | | | | | | RP23- 466J17 .3 |
| | | | | | | | | | | | | Texl6 |
| | | | | | | | | | | | | 493340 3O08Ri k |
| | | | | | | | | | | | | Apool |
| | | | | | | | | | | | | Satll |
| | | | | | | | | | | | | 201010 6E10Ri k |
| | | | | | | | | | | | | Zfp711 |
| | | | | | | | | | | | | Poflb |
| | | | | | | | | | | | | Gml49 36 |
| | | | | | | | | | | | | Chm |
| | | | | | | | | | | | | Dach2 |
| | | | | | | | | | | | | Klhl4 |
| | | | | | | | | | | | | Ube2d nil |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Ube2dnl2 |
| | | | | | | | | | | | | 4930555B12Rik |
| | | | | | | | | | | | | Cpxcrl |
| | | | | | | | | | | | | H2afb2 |
| | | | | | | | | | | | | Gml4920 |
| | | | | | | | | | | | | Gm28579 |
| | | | | | | | | | | | | Tgif2lx2 |
| | | | | | | | | | | | | Tgif2lx1 |
| | | | | | | | | | | | | Gml4929 |
| | | | | | | | | | | | | Pabpc5 |
| | | | | | | | | | | | | Pcdhllx |
| | | | | | | | | | | | | H2afb3 |
| | | | | | | | | | | | | Napll3 |
| | | | | | | | | | | | | Gml7521 |
| | | | | | | | | | | | | Cldn34cl |
| | | | | | | | | | | | | Astx6 |
| | | | | | | | | | | | | Srsx |
| | | | | | | | | | | | | Gml7577 |
| | | | | | | | | | | | | Gml4951 |
| | | | | | | | | | | | | Astx2 |
| | | | | | | | | | | | | Gml7412 |
| | | | | | | | | | | | | Cldn34c2 |
| | | | | | | | | | | | | Gml4950 |
| | | | | | | | | | | | | Gml7467 |
| | | | | | | | | | | | | Cldn34c3 |
| | | | | | | | | | | | | Astx5 |
| | | | | | | | | | | | | Vmn2r121 |
| | | | | | | | | | | | | Astxla |
| | | | | | | | | | | | | Gml75 |

| | | | | | | | | | | | | 84 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Astx4a |
| | | | | | | | | | | | | Gml7469 |
| | | | | | | | | | | | | Astx4b |
| | | | | | | | | | | | | Astxlb |
| | | | | | | | | | | | | Gml7361 |
| | | | | | | | | | | | | Gm21616 |
| | | | | | | | | | | | | Astx4c |
| | | | | | | | | | | | | Gml7693 |
| | | | | | | | | | | | | Astxlc |
| | | | | | | | | | | | | Gml7522 |
| | | | | | | | | | | | | Astx4d |
| | | | | | | | | | | | | Gml7267 |
| | | | | | | | | | | | | Astx3 |
| | | | | | | | | | | | | 4932411N23Rik |
| | | | | | | | | | | | | Gm382 |
| | | | | | | | | | | | | 4921511C20Rik |
| | | | | | | | | | | | | Cldn34c4 |
| | | | | | | | | | | | | 4930558G05Rik |
| | | | | | | | | | | | | Diaph2 |
| | | | | | | | | | | | | Pcdhl9 |
| | | | | | | | | | | | | Gm26851 |
| | | | | | | | | | | | | Tnmd |
| | | | | | | | | | | | | Tspan6 |
| | | | | | | | | | | | | Srpx2 |
| | | | | | | | | | | | | Sytl4 |
| | | | | | | | | | | | | Cstf2 |
| | | | | | | | | | | | | Noxl |
| | | | | | | | | | | | | Xkrx |
| | | | | | | | | | | | | Arll3a |
| | | | | | | | | | | | | Trmt2b |

| | | | | | | | | | | | | | Tmem 35 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Cenpi |
| | | | | | | | | | | | | | Drp2 |
| | | | | | | | | | | | | | Taf7l |
| | | | | | | | | | | | | | Timm8 al |
| | | | | | | | | | | | | | Btk |
| | | | | | | | | | | | | | Rpl36a |
| | | | | | | | | | | | | | Gla |
| | | | | | | | | | | | | | Hnrnp h2 |
| | | | | | | | | | | | | | Armcx 4 |
| | | | | | | | | | | | | | Armcx 1 |
| | | | | | | | | | | | | | Armcx 6 |
| | | | | | | | | | | | | | Armcx 3 |
| | | | | | | | | | | | | | Armcx 2 |
| | | | | | | | | | | | | | Nxf2 |
| | | | | | | | | | | | | | Zmatl |
| | | | | | | | | | | | | | Gml50 23 |
| | | | | | | | | | | | | | Tceal6 |
| | | | | | | | | | | | | | Pramel 3 |
| | | | | | | | | | | | | | Gm512 8 |
| | | | | | | | | | | | | | Gm790 3 |
| | | | | | | | | | | | | | AV320 801 |
| | | | | | | | | | | | | | Nxf7 |
| | | | | | | | | | | | | | Prame |
| | | | | | | | | | | | | | Tcpllx 2 |
| | | | | | | | | | | | | | Tmsbl 5a |
| | | | | | | | | | | | | | Armcx 5 |
| | | | | | | | | | | | | | Gprasp 1 |
| | | | | | | | | | | | | | Bhlhb9 |
| | | | | | | | | | | | | | Gprasp |

| | | | | | | | | | | | | 2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Arxes2 |
| | | | | | | | | | | | | Arxesl |
| | | | | | | | | | | | | Bex2 |
| | | | | | | | | | | | | Nxf3 |
| | | | | | | | | | | | | Bex4 |
| | | | | | | | | | | | | Tceal8 |
| | | | | | | | | | | | | Tceal5 |
| | | | | | | | | | | | | Bexl |
| | | | | | | | | | | | | Tceal7 |
| | | | | | | | | | | | | Wbp5 |
| | | | | | | | | | | | | Ngfrap1 |
| | | | | | | | | | | | | Kir3dl2 |
| | | | | | | | | | | | | Kir3dll |
| | | | | | | | | | | | | Tceal3 |
| | | | | | | | | | | | | Tceall |
| | | | | | | | | | | | | Morf4l2 |
| | | | | | | | | | | | | Glra4 |
| | | | | | | | | | | | | Plpl |
| | | | | | | | | | | | | Rab9b |
| | | | | | | | | | | | | H2bfm |
| | | | | | | | | | | | | Tmsbl5 1 |
| | | | | | | | | | | | | Tmsbl5b2 |
| | | | | | | | | | | | | Tmsbl5bl |
| | | | | | | | | | | | | Slc25a53 |
| | | | | | | | | | | | | Zcchcl8 |
| | | | | | | | | | | | | Faml99x |
| | | | | | | | | | | | | Esxl |
| | | | | | | | | | | | | Illrapl2 |
| | | | | | | | | | | | | Texl3a |
| | | | | | | | | | | | | Nrk |
| | | | | | | | | | | | | Serpina7 |
| | | | | | | | | | | | | 4930513O06Rik |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | 4933428M09Rik |
| | | | | | | | | | | | | Mumll1 |
| | | | | | | | | | | | | Trapla |
| | | | | | | | | | | | | D330045A20Rik |
| | | | | | | | | | | | | Rnfl28 |
| | | | | | | | | | | | | Tbcld8b |
| | | | | | | | | | | | | Gml5013 |
| | | | | | | | | | | | | Ripplyl |
| | | | | | | | | | | | | Cldn2 |
| | | | | | | | | | | | | Morc4 |
| | | | | | | | | | | | | Rbm41 |
| | | | | | | | | | | | | Nup62cl |
| | | | | | | | | | | | | Pihlh3b |
| | | | | | | | | | | | | Gml5046 |
| | | | | | | | | | | | | Frmpd3 |
| | | | | | | | | | | | | Prpsl |
| | | | | | | | | | | | | Tsc22d3 |
| | | | | | | | | | | | | Mid2 |
| | | | | | | | | | | | | Eif2c5 |
| | | | | | | | | | | | | Texl3 |
| | | | | | | | | | | | | Vsigl |
| | | | | | | | | | | | | Psmdl0 |
| | | | | | | | | | | | | Atg4a |
| | | | | | | | | | | | | Col4a6 |
| | | | | | | | | | | | | Col4a5 |
| | | | | | | | | | | | | Irs4 |
| | | | | | | | | | | | | Gml5295 |
| | | | | | | | | | | | | Gml5294 |
| | | | | | | | | | | | | Gml5298 |
| | | | | | | | | | | | | Gucy2f |

|  |  |  |  |  |  |  |  |  |  |  |  |  | Nxt2 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  |  |  |  |  |  |  | Kcnell |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Acsl4 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Tmem164 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Ammecrl |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Rgagl |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Chrdll |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Pak3 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Capn6 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Dcx |
|  |  |  |  |  |  |  |  |  |  |  |  |  | A730046J19Rik |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Algl3 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Trpc5 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Trpc5os |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Zcchcl6 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Lhfpll |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Amot |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Htr2c |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Ill3ra2 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Lrch2 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5128 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5080 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5107 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5114 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gm8334 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5127 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Luzp4 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5099 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Ott |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5092 |
|  |  |  |  |  |  |  |  |  |  |  |  |  | Gml5093 |

| | | | | | | | | | | | | Gml51 00 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Gml50 85 |
| | | | | | | | | | | | | Gml50 86 |
| | | | | | | | | | | | | Gml04 39 |
| | | | | | | | | | | | | Gml50 97 |
| | | | | | | | | | | | | Gml50 91 |
| | | | | | | | | | | | | Gml51 04 |
| | | | | | | | | | | | | Tmem 29 |
| | | | | | | | | | | | | Apex2 |
| | | | | | | | | | | | | Alas2 |
| | | | | | | | | | | | | Pfkfbl |
| | | | | | | | | | | | | Tro |
| | | | | | | | | | | | | Maged 2 |
| | | | | | | | | | | | | Gm271 91 |
| | | | | | | | | | | | | Gnl3l |
| | | | | | | | | | | | | Fgdl |
| | | | | | | | | | | | | Tsr2 |
| | | | | | | | | | | | | Gml51 38 |
| | | | | | | | | | | | | Wnk3 |
| | | | | | | | | | | | | A2300 72E10 Rik |
| | | | | | | | | | | | | Faml2 Oc |
| | | | | | | | | | | | | Phf8 |
| | | | | | | | | | | | | Huwel |
| | | | | | | | | | | | | Hsdl7 blO |
| | | | | | | | | | | | | Ribcl |
| | | | | | | | | | | | | Smcla |
| | | | | | | | | | | | | Iqsec2 |
| | | | | | | | | | | | | Kdm5c |
| | | | | | | | | | | | | Kantr |
| | | | | | | | | | | | | Tspyl2 |
| | | | | | | | | | | | | Gprl73 |

| | | | | | | | | | | | | | Cldn34a |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Shroom2 |
| | | | | | | | | | | | | | Gprl43 |
| | | | | | | | | | | | | | Usp51 |
| | | | | | | | | | | | | | Mageh1 |
| | | | | | | | | | | | | | Foxr2 |
| | | | | | | | | | | | | | Rragb |
| | | | | | | | | | | | | | Klf8 |
| | | | | | | | | | | | | | Ubqln2 |
| | | | | | | | | | | | | | Cypt3 |
| | | | | | | | | | | | | | Kctdl2b |
| | | | | | | | | | | | | | RP23-106P7.5 |
| | | | | | | | | | | | | | 2210013021Rik |
| | | | | | | | | | | | | | Spin2c |
| | | | | | | | | | | | | | Samtl |
| | | | | | | | | | | | | | 4921511M17Rik |
| | | | | | | | | | | | | | GmlOO57 |
| | | | | | | | | | | | | | Gml5140 |
| | | | | | | | | | | | | | 4930524N10Rik |
| | | | | | | | | | | | | | Samt4 |
| | | | | | | | | | | | | | Samt2 |
| | | | | | | | | | | | | | Cldn34bl |
| | | | | | | | | | | | | | Magea6 |
| | | | | | | | | | | | | | Magea3 |
| | | | | | | | | | | | | | Magea8 |
| | | | | | | | | | | | | | Magea2 |
| | | | | | | | | | | | | | Magea5 |
| | | | | | | | | | | | | | Magea |

| | | | | | | | | | | | | | 1 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | Cldn34b2 |
| | | | | | | | | | | | | | Satl |
| | | | | | | | | | | | | | Acot9 |
| | | | | | | | | | | | | | Prdx4 |
| | | | | | | | | | | | | | Ptchdl |
| | | | | | | | | | | | | | Gml5156 |
| | | | | | | | | | | | | | Gml5155 |
| | | | | | | | | | | | | | Phex |
| | | | | | | | | | | | | | Sms |
| | | | | | | | | | | | | | Mbtps2 |
| | | | | | | | | | | | | | Yy2 |
| | | | | | | | | | | | | | Smpx |
| | | | | | | | | | | | | | Gml5169 |
| | | | | | | | | | | | | | Klhl34 |
| | | | | | | | | | | | | | Cnksr2 |
| | | | | | | | | | | | | | Rps6ka3 |
| | | | | | | | | | | | | | Eiflax |
| | | | | | | | | | | | | | Map7d2 |
| | | | | | | | | | | | | | A830080D01Rik |
| | | | | | | | | | | | | | Sh3kbp1 |
| | | | | | | | | | | | | | Map3k15 |
| | | | | | | | | | | | | | Pdhal |
| | | | | | | | | | | | | | Adgrg2 |
| | | | | | | | | | | | | | Gml5241 |
| | | | | | | | | | | | | | Phka2 |
| | | | | | | | | | | | | | Gml5243 |
| | | | | | | | | | | | | | Ppefl |
| | | | | | | | | | | | | | Rsl |
| | | | | | | | | | | | | | Cdkl5 |
| | | | | | | | | | | | | | Gja6 |
| | | | | | | | | | | | | | Scml2 |

| | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Gml5262 |
| | | | | | | | | | | | | Rai2 |
| | | | | | | | | | | | | Scmll |
| | | | | | | | | | | | | Gml5205 |
| | | | | | | | | | | | | Nhs |
| | | | | | | | | | | | | Gml5202 |
| | | | | | | | | | | | | Reps2 |
| | | | | | | | | | | | | Rbbp7 |
| | | | | | | | | | | | | Txlng |
| | | | | | | | | | | | | Syapl |
| | | | | | | | | | | | | Ctps2 |
| | | | | | | | | | | | | SlOOg |
| | | | | | | | | | | | | Grpr |
| | | | | | | | | | | | | Rnfl38rtl |
| | | | | | | | | | | | | Apls2 |
| | | | | | | | | | | | | Zrsr2 |
| | | | | | | | | | | | | Car5b |
| | | | | | | | | | | | | Siahlb |
| | | | | | | | | | | | | Tmem27 |
| | | | | | | | | | | | | Ace2 |
| | | | | | | | | | | | | Bmx |
| | | | | | | | | | | | | Pir |
| | | | | | | | | | | | | Figf |
| | | | | | | | | | | | | Piga |
| | | | | | | | | | | | | Asbll |
| | | | | | | | | | | | | Asb9 |
| | | | | | | | | | | | | Mospd2 |
| | | | | | | | | | | | | Fancb |
| | | | | | | | | | | | | Gml7604 |
| | | | | | | | | | | | | Glra2 |
| | | | | | | | | | | | | Gemin8 |
| | | | | | | | | | | | | Gpm6b |
| | | | | | | | | | | | | Ofdl |
| | | | | | | | | | | | | Trappc2 |
| | | | | | | | | | | | | Rab9 |

| | | | | | | | | | | | | Tceanc |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | Egfl6 |
| | | | | | | | | | | | | Gml52 26 |
| | | | | | | | | | | | | Gml72 0 |
| | | | | | | | | | | | | Gml52 30 |
| | | | | | | | | | | | | Gm881 7 |
| | | | | | | | | | | | | Gml52 32 |
| | | | | | | | | | | | | Gml52 28 |
| | | | | | | | | | | | | Tmsb4 x |
| | | | | | | | | | | | | Tlr8 |
| | | | | | | | | | | | | Tlr7 |
| | | | | | | | | | | | | Prps2 |
| | | | | | | | | | | | | Gml52 39 |
| | | | | | | | | | | | | Frmpd 4 |
| | | | | | | | | | | | | Msl3 |
| | | | | | | | | | | | | Arhgap 6 |
| | | | | | | | | | | | | Gml52 61 |
| | | | | | | | | | | | | Amelx |
| | | | | | | | | | | | | Hccs |
| | | | | | | | | | | | | Gml52 45 |
| | | | | | | | | | | | | Midi |
| | | | | | | | | | | | | 493340 OAllRi k |
| | | | | | | | | | | | | Gml57 26 |
| | | | | | | | | | | | | Gml52 47 |
| | | | | | | | | | | | | Gm218 87 |
| | | | | | | | | | | | | Asmt |

[0148]    As an additional validation, we modified an existing trajectory finding technique, Wishbone(SlO) — based on shortest paths in k-NN graphs — to include information about time and proliferation. This gives trajectories whose overall shape agrees with the transports displayed in **FIG.** 8A.

**Learning gene regulatory networks**

[0149]    How to set up an optimization problem to solve for a regulatory function that fits the transport maps is described above.

[0150]    In order to make this concrete, a function class F was specified over which to optimize. Consider a rectified-linear function class defined in terms of a specific generalized logistic function

$$\ell(x; k, b, y_0, x_0) = \frac{k\,y\,\circ}{y_0 + (k - yo)e^{-b(x - x_0)}},$$

where k, b, yO, xO $\in$ R are parameters of the generalized logistic function l(x). A function class F is defined consisting of functions f : RG $\to$ RG of the form

$$f(x) = U\ell\{WTx),$$

where **1** is applied entry-wise to the vector WZx $\in$ R$^M$ to obtain a vector that we multiply against U $\in$ RGxM . Here T $\in$ RGTF ×G denotes a projection operator that selects only the coordinates of x that are transcription factors, and GTF is the number of transcription factors.

[0151]    The following optimization over matrices U $\equiv$ RGxM and W $\equiv$ RMxGTF

$$\min_{U,W} \quad \mathbb{E}_r \left\| \frac{X_{t_i} - X_{t_{i+1}}}{\Delta_t} - U\ell(WTX_{t_i}) \right\|^2 + \eta_1 \|U\|_1 + \eta_2 \|W\|_1, + \eta_3 \|W\|_2^2$$
$$\text{s.t.} \quad U \geq 0.$$

where (Xti , Xti+i ) is a pair of random variables distributed according to the normalized transport map $r$ and // U // i denotes the sparsity-promoting $l_1$ norm of U, viewed as a vector (that is, the sum of the absolute value of the entries of U ). Each rank one component (row of U or column of W ) gives us a group of genes controlled by a set of transcription factors. The regularization parameters $\eta\iota$ and $\eta_2$ control the sparsity level (i.e. number of genes in these groups).

[0152]    **Implementation:** A stochastic gradient descent algorithm was designed to solve [10]. Over a sequence of epochs, the algorithm samples batches of points (Xti , Xti+i ) from the transport maps, computes the gradient of the loss, and updates the optimization variables U and

W. The batch sizes are determined by the Shannon diversity of the transport maps: for each pair of consecutive time points, the Shannon diversity S was computed of the transport map, then randomly sample max(S x 10-5, 10) pairs of points to add to the batch. We run for a total of 10, 000 epochs.

[0153]    This algorithm was implemented in Python.

## 7. Clustering cells

[0154]    Cells were clustered using the Louvain-Jaccard community detection algorithm (SI 9-S21) in 20 dimensional diffusion component space. This algorithm maximizes the Louvain modularity — a value between - 1 and 1 that measures the density of links inside communities compared to links between communities.

[0155]    As a first step, the 20-nearest neighbor graph in 20 dimensional diffusion component space (computed on cells from both 2i and serum) were computed. The edges are weighted in this graph by the Jaccard similarity coefficient. The resulting graph was partitioned into clusters using the Louvain community detection algorithm (SI 9) implemented in the function multilevel. community from the R pack- age IGRAPH (1.0.1) (S22). The default parameters for automatically selecting the number of clusters gave us 33 clusters, displayed in **FIG. 7D.**

## 8. Gene correlation modules reveal biological signatures

[0156]    In this section technique for identifying modules of correlated genes are described, with the goal of revealing coherent biological processes.

[0157]    The procedure consists of two steps. In the first step, the Graphical Lasso (S23) was used to compute a regularized estimate of the covariance matrix for the 66,000 expression profiles. The Graphical Lasso fits a covariance matrix to the data, regularized so that the inverse of the covariance matrix is *sparse* (i.e. has only a few non-zeros). The motivation for selecting a sparse inverse covariance is based on the fact that if a collection of observations have a multivariate Gaussian distribution with mean $\mu$ and covariance $\Sigma$, then the zero pattern of $\Sigma^{-1}$ completely specifies the conditional independence structure of the observations:

$$\Sigma_{ij}^{-1} = 0 \quad \Leftrightarrow \quad \text{variables } i \text{ and } j \text{ are conditionally independent given the other variables.}$$

Let $\Theta = \Sigma^{-1}$ and let $S$ denote the empirical covariance for our expression profiles

[0158]    The Graphical Lasso maximizes the Gaussian log likelihood:

$$\underset{\Theta}{\text{maximize}} \quad \log \det \Theta - \mathbf{tr}(S\Theta) - \rho \|\Theta\|_1.$$

Here $\|\Theta\|_1$ is a regularization term that promotes sparse solutions. The optimal $\Theta$ is a (regularized) maximum-likelihood estimate of the inverse covariance matrix $\Sigma^{-1}$ for a Gaussian ensemble.

**[0159]** Gene modules were identifed as tightly knit communities in the network specified by $\Theta$ (see below). Based on these gene modules, we then identified gene signatures related to specific pathways, cell types, and conditions. We did this by functional enrichment analysis (see below). The gene modules are displayed in **FIG. 13.**

**[0160]** Computing gene modules: The glasso package was used (S23) to solve the graphical lasso optimization problem. The regularization parameter p was tuned to achieve a desirable sparsity level for $\Theta$. In particular, we select a value of p that gave around 10,000 total genes (i.e. 10,000 non-zero rows and columns of $\Theta$).

**[0161]** Viewing $\Theta$ as an adjacency matrix defining a network of genes, we partitioned the network using with the Infomap community detection algorithm (S24) from the R package IGRAPH (vl.1.0) (S22), retaining modules that contain more than 10 genes. This yields 44 gene modules, each consisting of a set of genes. The modules are visualized in **FIG. 13.**

**[0162]** **Functional enrichments:** Functional enrichment analysis was performed on the gene sets defined by the modules using the findGO.pl program from the HOMER suite (Hypergeometric Optimization of Motif Enrichment, version: 4.9.1) (S12) with Benjamini and Hochberg correction for multiple hypothesis testing (retaining terms at adjusted p-value $< 0.05$). All genes that passed quality-control filters were used as a background set.

**[0163]** This yielded a set of biological signatures related to each module.

**[0164]** Computing scores from gene sets Given a set of genes (coming from a gene module or biological signature), cells were scored based on their gene expression. In particular, for a given cell the z-score for each gene in the set was determined. The z-scores were then truncated at 5 or -5, and define the signature of the cell to be the mean z-score over all genes in the gene set. The scores for the gene modules are visualized in **FIG. 13** and the scores for the biological signatures are visualized in **FIGs. 7A-7F.**

**Example 2 Reprogramming to iPSCs as a test case for analysis of developmental landscapes.**

103

[0165]    WADDINGTON-OT was used to analyze the reprogramming of fibroblasts to iPSCs (39-42).

[0166]    Studies have applied scRNA-Seq, but they have involved only several dozen cells or several dozen genes (13, 43). Studies have proposed that reprogramming involves two "transcriptional waves," with gain of proliferation and loss of fibroblast identity followed by transient activation of developmental regulators and gradual activation of embryonic stem cell (ESC) genes (12). Some studies (16, 44, 45), have noted strong upregulation of lineage-specific genes from unrelated lineages (e.g., related to neurons), but it has been unclear whether this largely reflects disorganized gene activation by TFs or coherent differentiation of specific (off-target) cell types (45).

[0167]    scRNA-seq profiles of 65,781 cells were collected across a 16-day time course of iPSC induction, under two conditions (FIGs. 6A,6B). An efficient "secondary" reprogramming system was used (46), as described hereinbelow.

[0168]    Mouse embryonic fibroblasts (MEFs) were obtained from a single female embryo homozygous for ROSA26-M2rtTA, which constitutively expresses a reverse transactivator controlled by doxycycline (Dox), a Dox-inducible polycistronic cassette carrying Pou5fl (Oct4), Klf4, *Sox2,* and *Myc (OKSM),* and an EGFP reporter incorporated into the endogenous *Oct4* locus (Oct4-IRES-EGFP). MEFs were plated in serum-containing induction medium, with Dox added on day 0 to induce the OKSM cassette (Phase-l(Dox)). Following Dox withdrawal at day 8, cells were transferred to either serum-free N2B27 2i medium (Phase-2(2i)) or maintained in serum (Phase-2(serum)). Oct4 EGFP+ cells emerged on day 10 as a reporter for "successful" reprogramming to endogenous Oct4 expression (FIG. 6C). Single or duplicate samples were collected at the various time points (FIG. 6A), single cell suspensions were generated and scRNA-Seq (Table 8, FIGs. 11A-11D) was performed. Samples were also collected from established iPSC lines reprogrammed from the same MEFs, maintained in either 2i or serum conditions. Overall, 68,339 cells were programed to an average depth of 38,462 reads per cell (Table 8). After discarding cells with less than 1,000 genes detected, a total of 65,781 cells were retained, with a median of 2,398 genes and 7,387 unique transcripts per cell.

**Table 8**

| Sample (Day) | Phase | Number of Cells | Number of cells (filtered) | Number of reads | Mean Reads per Cells | Median Genes per Cell | Median UMI Counts per Cell | cDNA PCR Duplication % |
|---|---|---|---|---|---|---|---|---|
| D O | Dox | 4241 | 4060 | 111,286,101 | 26240 | 2446 | 6495 | 50.5 |
| D2-1 | Dox | 2909 | 2890 | 143,713,479 | 49403 | 2867 | 8401 | 55.6 |
| D2-2 | Dox | 2758 | 2729 | 109,907,870 | 39850 | 2521 | 6271 | 70.2 |
| D4-1 | Dox | 2889 | 2882 | 126,824,856 | 43899 | 2447 | 7349 | 57.3 |
| D4-2 | Dox | 3976 | 3962 | 99,109,221 | 24926 | 2386 | 7446 | 34.1 |
| D6-1 | Dox | 3676 | 3198 | 132,565,146 | 36062 | 1453 | 3147 | 84 |
| D6-2 | Dox | 3534 | 3168 | 99,748,307 | 28225 | 1533 | 3567 | 76.5 |
| D8-1 | Dox | 2177 | 2142 | 98,462,446 | 45228 | 2332 | 8216 | 65.7 |
| D8-2 | Dox | 3677 | 2625 | 95,807,550 | 26055 | 1486 | 3862 | 62.6 |
| D9-1 | 2i | 2445 | 2441 | 122,451,561 | 50082 | 2843 | 11799 | 51.8 |
| D9-2 | 2i | 2183 | 2174 | 125,014,976 | 57267 | 2734 | 11183 | 57 |
| DlO-1 | 2i | 2878 | 2878 | 129,837,247 | 45113 | 2625 | 9570 | 58.1 |
| D10-2 | 2i | 2620 | 2619 | 126,364,110 | 48230 | 2647 | 9930 | 59.5 |
| Dll | 2i | 1532 | 1529 | 119,736,956 | 78157 | 2892 | 10744 | 65.9 |
| D12-1 | 2i | 5144 | 5139 | 158,679,538 | 30847 | 2269 | 6299 | 41 |
| D12-2 | 2i | 2156 | 2155 | 112,512,277 | 52185 | 2651 | 8633 | 54.8 |
| D16 | 2i | 4621 | 4500 | 117,242,910 | 25371 | 2203 | 7761 | 39.5 |
| iPSCs | 2i | 2917 | 2916 | 139,441,360 | 47803 | 3172 | 12775 | 38.2 |
| D10 | serum | 2094 | 2088 | 115,832,953 | 55316 | 2717 | 9733 | 58.4 |
| D12 | serum | 2913 | 2895 | 96,402,567 | 33093 | 2711 | 8819 | 44.2 |
| D16 | serum | 3875 | 3703 | 119,329,130 | 30794 | 1953 | 4984 | 53.6 |
| iPSCs | serum | 3124 | 3088 | 128,207,617 | 41039 | 2637 | 9689 | 46.1 |
| | | | | | | | | |
| | Total | 68339 | 65781 | | | | | |
| | | | Average depth per cell: | 38,462 | | | | |

**Example 3 The reprogramming landscape reveals relationships among biological features.**

[0169]        WADDINGTON-OT was used to generate a transport map across the cells in the time course described in the previous example. Based on similarity of expression profiles, the 16,339 detected genes were partitioned into 44 gene modules and the 65,781 cells into 33 cell clusters. Some of the clusters contained cells from more than one time point, reflecting asynchrony in the reprogramming process. The landscape of reprogramming was explored by identifying cell

subsets of interest (e.g., successfully reprogrammed cells at day 16, or each of the cell clusters), studying the trajectories to and from these subsets (e.g., characterizing the pattern of gene expression in ancestors at day 8 of successfully reprogrammed target cells at day 16), and considering contemporaneous interactions between them. The analyses were visualized in a two-dimensional embedding using FLE (Fig. 7A), annotated in various ways. FLE reflects better global structures in the data presented herein than other modes of visualization (Figs. 12A-12C). These annotations include time points and growth conditions (Figs. 7B,7C), gene modules (Figs. 13, 14A-14B, Table 1), cell clusters (Fig. 7D, Fig. 14A-14D, Table 9), expression of gene signatures (curated gene sets associated with specific cell types, pathways, and responses, such as MEF identity, proliferation, pluripotency, and apoptosis; Fig. 7E, Table 7), expression of individual genes (Fig. 7F, Fig. 15), and ancestor and descendant distributions (Figs. 8A-8F). Extensive sensitivity analysis showed that key biological results for the reprogramming data were largely robust to the details of the formulation. Finally, the WADDINGTON-OT landscape was compared to the landscapes produced by various graph-based methods. The results show the following. Cell trajectories start at the lower right corner at day 0, proceed leftward to day 2 and then upward towards two regions identified as the Valley of Stress and the Horn of Transformation (Fig. 7B, Fig. 8A). The Valley is characterized by signatures of cellular stress, senescence, and, in some regions, apoptosis (Fig. 7E); it appears to be a terminal destination. By contrast, the Horn is characterized by increased proliferation, loss of fibroblast identity, a mesenchymal-to-epithelial transition (Fig. 7E), and early appearance of certain pluripotency markers (e.g., Nanog and Zfp42, Fig. 7F), which are predictive features of successful reprogramming (47). Some of the cells in the Horn proceed toward pre-iPSCs by day 12 and iPSCs by day 16, while others encounter alternative fates of placental-like development and neurogenesis (in serum, but not 2i condition; Figs. 7B, 7C). A more detailed account of the landscape is in the following examples.

**Table 9**

| Cluster | Phase-l(Dox) | | | | | Phase-2 (2i) | | | | | | Phase-2 (serum) | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | D0 | D2 | D4 | D6 | D8 | D9 | D10 | Dll | D12 | D16 | iPSCs | D10 | D12 | D16 | iPSCs |
| 1 | 97.4 | 0.1 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.4 | 0.1 | 0.9 |

| | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| **2** | 2.0 | 0.3 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.1 | 0.1 |
| **3** | 0.1 | 22.0 | 0.9 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **4** | 0.0 | 31.7 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **5** | 0.2 | 33.5 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 0.1 | 0.1 | 0.0 | 0.0 |
| **6** | 0.0 | 12.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **7** | 0.0 | 0.1 | 60.7 | 5.8 | 0.1 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **8** | 0.0 | 0.0 | 23.9 | 8.3 | 2.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **9** | 0.0 | 0.0 | 0.9 | 16.5 | 16.8 | 1.2 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.0 | 0.0 | 0.0 | 0.0 |
| **10** | 0.0 | 0.0 | 0.0 | 2.4 | 15.1 | 19.3 | 0.5 | 0.3 | 0.0 | 0.0 | 0.0 | 21.8 | 0.0 | 0.1 | 0.0 |
| **11** | 0.0 | 0.0 | 0.0 | 0.2 | 1.3 | 22.6 | 14.1 | 7.1 | 1.5 | 0.1 | 0.0 | 14.4 | 2.9 | 0.7 | 0.1 |
| **12** | 0.2 | 0.0 | 0.0 | 0.0 | 0.0 | 3.2 | 16.0 | 11.4 | 9.7 | 1.1 | 0.6 | 3.0 | 13.9 | 2.6 | 0.2 |
| **13** | 0.1 | 0.0 | 0.0 | 0.0 | 0.4 | 9.1 | 11.5 | 8.6 | 3.4 | 0.2 | 0.0 | 18.1 | 16.8 | 1.8 | 0.1 |
| **14** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 2.9 | 4.8 | 12.3 | 1.4 | 1.5 | 0.0 | 2.5 | 0.6 | 0.0 |
| **15** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 1.2 | 5.6 | 11.6 | 6.2 | 5.3 | 0.0 | 0.2 | 0.6 | 0.0 |
| **16** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.7 | 5.9 | 14.2 | 16.0 | 2.5 | 0.0 | 0.3 | 1.0 | 1.5 | 0.0 |
| **17** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.6 | 10.5 | 11.9 | 6.7 | 0.2 | 0.0 | 0.0 | 0.9 | 0.2 | 0.0 |
| **18** | 0.0 | 0.1 | 12.5 | 15.9 | 1.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 |
| **19** | 0.0 | 0.0 | 0.0 | 10.6 | 27.5 | 11.6 | 0.0 | 0.1 | 0.0 | 0.0 | 0.0 | 5.6 | 0.0 | 0.0 | 0.0 |
| **20** | 0.0 | 0.0 | 0.6 | 31.7 | 20.0 | 4.3 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.0 | 0.0 |
| **21** | 0.0 | 0.0 | 0.0 | 8.5 | 15.5 | 24.9 | 0.1 | 0.1 | 0.1 | 0.0 | 0.0 | 32.5 | 0.2 | 0.6 | 0.1 |
| **22** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.6 | 25.8 | 10.1 | 0.5 | 0.1 | 0.0 | 1.2 | 1.0 | 0.3 | 0.1 |
| **23** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.3 | 0.1 | 0.5 | 0.1 | 0.0 | 0.7 | 29.2 | 16.5 | 1.7 |
| **24** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.3 | 8.6 | 11.6 | 6.3 | 1.6 | 0.1 | 0.2 | 16.8 | 7.7 | 0.1 |
| **25** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.3 | 7.3 | 0.4 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 |
| **26** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.6 | 1.0 | 0.3 | 0.1 | 0.0 | 0.0 | 0.8 | 30.7 | 0.0 |
| **27** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.6 | 0.1 | 0.0 | 0.0 | 0.0 | 3.0 | 0.0 |
| **28** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 1.8 | 12.7 | 23.0 | 2.3 | 0.7 | 0.6 | 12.7 | 0.6 | 0.0 |
| **29** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 0.0 | 31.6 | 0.0 | 0.0 | 0.0 | 1.1 | 0.0 |
| **30** | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 33.4 | 0.1 | 0.0 | 0.1 | 0.4 | 0.0 |

| 31 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 15.4 | 1.6 | 0.0 | 0.1 | 23.3 | 1.1 |
| 32 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 6.6 | 95.5 |
| 33 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.1 | 3.1 | 90.2 | 0.0 | 0.0 | 0.8 | 0.1 |

**Example 4**

**[0170]** Predictive markers of reprogramming success are detectable by day 2.

**[0171]** The vast majority (>98%) of cells at day 0 fall into a single cluster characterized by a strong signature of MEF identity, with clear bimodality in the proliferation signature (Fig. 16A). By day 2 after Dox treatment, cells show high levels of expression of the OKSM cassette and have begun to diverge in their responses (clusters 3, 4, 5, 6, Fig. 7D). Overall, they score highly for expression signatures of proliferation, MEF identity, and endoplasmic reticulum (ER) stress (reflecting high secretion in mesenchymal cells) (Fig. 7E).

**[0172]** However, the cells exhibit considerable heterogeneity, seen most clearly by comparing the cells in clusters 4 and 6, which vary in their expression signatures and in their fates (Figs. 8A, 8B and Figs. 17A-17C). While cells in both clusters are highly proliferative, cells in cluster 4 have begun to lose MEF identity, show lower ER stress, and have higher OKSM-cassette expression, while cells in cluster 6 have the opposite properties (FIGs. 7D, 7E and Fig. 16B). The cells in the two clusters show clear differences in their enrichment in the ancestral distribution of iPSCs (Fig. 8D). The majority (54%) of the day 2 ancestors of iPSCs lie in cluster 4, while only a small fraction (3%) lie in cluster 6. Clusters 4 and 6 also show clear differences in their descendants (Figs. 8A, 8C and Fig. 17A): the descendants of cells in cluster 6 are strongly biased toward the Valley of Stress (e.g., 81% of Cluster 6 cell descendants are in clusters 8-11 by day 8 vs. 18% for cluster 4), while cluster 4 is strongly biased toward the Horn of Transformation (e.g., 81% in clusters 19-21 vs. 12% for cluster 6).

**[0173]** The strongest difference in gene expression between clusters 4 and 6 was seen for Shisa8 (detected in 67% vs. 3% of cells in clusters 4 and 6, respectively) (Fig. 7F, fig. 16B) and Shisa8+ cells are enriched among the day 2 ancestors of iPSCs (Fig. 16B). Notably, Shisa8 is strongly associated with the entire trajectory toward successful reprogramming (Fig. 7F): it is expressed in the Horn, pre-iPSCs, and iPSCs, but not in the Valley or in the alternative fates of neurogenesis and placental development. The expression pattern of Shisa8 is similar to, but

stronger than, that of Fut9 (Fig. 15), a known early marker of successful reprogramming that synthesizes the surface glyco-antigen SSEA-1 (12). Shisa8 is a little-studied mammalian specific member of the Shisa gene family in vertebrates, which encodes single-transmembrane proteins that play roles in development and are thought to serve as adaptor proteins (48). The analysis suggests that Shisa8 may serve as a useful early predictive marker of eventual reprogramming success and may play a functional role in the process.

**Example 5 Cells in the valley of stress induce a Senescence Associated Secretion Phenotype (SASP).**

[0174]     By day 4, cells display a bimodal distribution of properties that is strongly correlated with their eventual descendants: cells in cluster 8 (low proliferation, high MEF identity, Fig. 7D, E and Fig. 16C) have 95% of their descendants in the Valley (Figs. 8A, 8B and Fig. 17A), while cells in cluster 18 (high proliferation, low MEF identity, Figs. 7D, 7E and Fig. 16C) have 94% of their descendants in the Horn (Figs. 8A, 8B and Fig. 17A and Table 10). Cells in cluster 7 show intermediate properties and have roughly equal probabilities of each fate (Fig. 8A, 8B and Fig. 17A).

**Table 10**

| Cluster | To1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 | 11 | 12 | 13 | 14 | 15 | 16 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| From 1 | | 0.001 | 0.920 | 0.980 | 0.978 | 0.987 | 0.001 | 0.001 | 0.000 | 0.000 | 0.000 | 0.001 | 0.008 | 0.001 | 0.002 | 0.003 |
| 2 | | 0.790 | 0.000 | 0.003 | 0.003 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | | 0.000 | 0.012 | 0.005 | 0.000 | 0.000 | 0.206 | 0.166 | 0.012 | 0.002 | 0.002 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | | 0.007 | 0.058 | 0.002 | 0.000 | 0.000 | 0.265 | 0.044 | 0.004 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 5 | | 0.106 | 0.008 | 0.003 | 0.006 | 0.003 | 0.293 | 0.298 | 0.004 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 |
| 6 | | 0.000 | 0.000 | 0.000 | 0.007 | 0.010 | 0.100 | 0.074 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 |
| 7 | | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.131 | 0.169 | 0.383 | 0.143 | 0.040 | 0.000 | 0.005 | 0.000 | 0.000 | 0.000 |
| 8 | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.003 | 0.240 | 0.171 | 0.126 | 0.018 | 0.000 | 0.005 | 0.000 | 0.000 | 0.000 |
| 9 | | 0.002 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.006 | 0.163 | 0.197 | 0.062 | 0.031 | 0.168 | 0.021 | 0.001 | 0.046 |
| 10 | | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.0 | 0.2 | 0.0 | 0.3 | 0.0 | 0.0 |

| | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | 05 | 00 | 00 | 00 | 00 | 00 | 00 | 11 | 63 | 88 | 83 | 93 | 77 | 25 | 37 |
| **11** | | 0.004 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.002 | 0.001 | 0.031 | 0.216 | 0.081 | 0.211 | 0.085 | 0.065 |
| **12** | | 0.012 | 0.000 | 0.004 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.020 | 0.127 | 0.032 | 0.166 | 0.269 | 0.152 |
| **13** | | 0.012 | 0.001 | 0.003 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.013 | 0.112 | 0.236 | 0.085 | 0.514 | 0.578 |
| **14** | | 0.002 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.003 | 0.017 | 0.002 | 0.028 | 0.037 | 0.017 |
| **15** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.001 | 0.006 | 0.005 |
| **16** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.003 | 0.005 | 0.003 | 0.025 | 0.026 |
| **17** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.003 | 0.003 | 0.003 | 0.026 | 0.027 |
| **18** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.002 | 0.003 | 0.201 | 0.079 | 0.013 | 0.003 | 0.001 | 0.000 | 0.000 | 0.000 |
| **19** | | 0.007 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.029 | 0.120 | 0.357 | 0.123 | 0.272 | 0.036 | 0.001 | 0.032 |
| **20** | | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.018 | 0.172 | 0.270 | 0.047 | 0.052 | 0.001 | 0.000 | 0.002 |
| **21** | | 0.010 | 0.000 | 0.000 | 0.004 | 0.000 | 0.000 | 0.000 | 0.001 | 0.094 | 0.075 | 0.021 | 0.036 | 0.035 | 0.001 | 0.005 |
| **22** | | 0.002 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.004 | 0.001 | 0.006 | 0.003 | 0.002 |
| **23** | | 0.027 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.005 | 0.004 | 0.001 | 0.021 | 0.004 | 0.003 |
| **24** | | 0.010 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.002 | 0.001 | 0.005 | 0.003 | 0.002 |
| **25** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **26** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **27** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **28** | | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **29** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **30** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **31** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **32** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| **33** | | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |

**Table 10 (Cont'd)**

| Cluster To | 17 | 18 | 19 | 20 | 21 | 22 | 23 | 24 | 25 | 26 | 27 | 28 | 29 | 30 | 31 | 32 | 33 |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| From 1 | 0.003 | 0.003 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.004 | 0.006 | 0.000 | 0.006 | 0.002 | 0.001 | 0.006 | 0.001 |
| 2 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 3 | 0.000 | 0.051 | 0.001 | 0.004 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 4 | 0.000 | 0.276 | 0.000 | 0.005 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 5 | 0.000 | 0.009 | 0.000 | 0.001 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 6 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 7 | 0.000 | 0.578 | 0.183 | 0.340 | 0.044 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 8 | 0.000 | 0.008 | 0.008 | 0.001 | 0.005 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 9 | 0.026 | 0.004 | 0.047 | 0.003 | 0.073 | 0.011 | 0.001 | 0.005 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 |
| 10 | 0.058 | 0.000 | 0.033 | 0.001 | 0.069 | 0.080 | 0.065 | 0.026 | 0.015 | 0.001 | 0.001 | 0.009 | 0.001 | 0.003 | 0.000 | 0.001 | 0.000 |
| 11 | 0.111 | 0.000 | 0.003 | 0.001 | 0.006 | 0.005 | 0.000 | 0.000 | 0.000 | 0.007 | 0.012 | 0.001 | 0.012 | 0.004 | 0.003 | 0.012 | 0.001 |
| 12 | 0.084 | 0.000 | 0.000 | 0.000 | 0.000 | 0.014 | 0.000 | 0.000 | 0.000 | 0.025 | 0.046 | 0.002 | 0.043 | 0.015 | 0.009 | 0.041 | 0.004 |
| 13 | 0.650 | 0.000 | 0.001 | 0.000 | 0.001 | 0.015 | 0.000 | 0.000 | 0.000 | 0.037 | 0.066 | 0.003 | 0.057 | 0.020 | 0.011 | 0.055 | 0.005 |
| 14 | 0.006 | 0.000 | 0.000 | 0.000 | 0.000 | 0.003 | 0.000 | 0.000 | 0.000 | 0.006 | 0.010 | 0.000 | 0.010 | 0.004 | 0.002 | 0.010 | 0.001 |
| 15 | 0.002 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 16 | 0.020 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.001 | 0.002 | 0.000 | 0.002 | 0.001 | 0.000 | 0.002 | 0.000 |
| 17 | 0.015 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.001 | 0.002 | 0.000 | 0.001 | 0.000 | 0.000 | 0.001 | 0.000 |
| 18 | 0.000 | 0.064 | 0.264 | 0.227 | 0.116 | 0.007 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 19 | 0.014 | 0.003 | 0.143 | 0.057 | 0.107 | 0.104 | 0.050 | 0.073 | 0.017 | 0.001 | 0.000 | 0.045 | 0.003 | 0.013 | 0.000 | 0.002 | 0.000 |
| 20 | 0.001 | 0.006 | 0.304 | 0.309 | 0.336 | 0.276 | 0.011 | 0.005 | 0.000 | 0.001 | 0.000 | 0.002 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 |
| 21 | 0.006 | 0.000 | 0.014 | 0.052 | 0.235 | 0.387 | 0.339 | 0.260 | 0.083 | 0.032 | 0.013 | 0.744 | 0.021 | 0.082 | 0.006 | 0.017 | 0.003 |
| 22 | 0.001 | 0.000 | 0.000 | 0.000 | 0.000 | 0.008 | 0.014 | 0.001 | 0.001 | 0.008 | 0.007 | 0.000 | 0.009 | 0.003 | 0.002 | 0.008 | 0.001 |
| 23 | 0.001 | 0.000 | 0.000 | 0.000 | 0.005 | 0.076 | 0.498 | 0.008 | 0.089 | 0.663 | 0.396 | 0.005 | 0.243 | 0.076 | 0.047 | 0.223 | 0.021 |
| 24 | 0.001 | 0.000 | 0.000 | 0.000 | 0.001 | 0.010 | 0.020 | 0.622 | 0.793 | 0.145 | 0.201 | 0.011 | 0.197 | 0.111 | 0.095 | 0.183 | 0.067 |

| | | | | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 25 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 26 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.061 | 0.228 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 27 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.005 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |
| 28 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.000 | 0.000 | 0.000 | 0.006 | 0.004 | 0.174 | 0.364 | 0.640 | 0.804 | 0.406 | 0.885 |
| 29 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.002 | 0.002 | 0.002 | 0.002 | 0.001 |
| 30 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.004 | 0.003 | 0.003 | 0.004 | 0.002 |
| 31 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.009 | 0.008 | 0.007 | 0.010 | 0.004 |
| 32 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.001 | 0.001 | 0.000 | 0.015 | 0.010 | 0.008 | 0.016 | 0.005 |
| 33 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 | 0.000 |

[0175]    Along the trajectory from cluster 8 to the Valley (days 10-16; Fisg. 8A, 8B and 8E,F), cells show a strong decrease in cell proliferation (Fig. 7E), accompanied by increased expression of various cell-cycle inhibitors, such as Cdkn2a, which encodes pi6, an inhibitor of the Cdk4/6 kinase and halts Gl/S transition (Fig. 7F), Cdknla (p21), and Cdkn2b (pl5) (Fig. 16D), which peaks in the Valley. The cells show increased expression of D-type cyclin gene Ccnd2 (Figs. 15, 16D) associated with growth arrest (49). A subset of the cells in the Valley (29%; clusters 12 and 14) showed high activity for a gene module that is correlated with a p53 pro-apoptotic signature, compared to all other cells inside the Valley (p-value< 10-16, average difference 0.17, Mest) and outside the Valley (p-value< 10-16, average difference 0.32, Mest) (Fig. 7E, fig. 16E).

[0176]    Cells in the Valley also show activation of signatures of extracellular-matrix (ECM) rearrangement and secretory functions (Fig. 7E, Fig. 16E). Because these properties are consistent with a senescence associated secretory phenotype (SASP), a SASP signature involving 60 genes (50) was used. Cells with this signature appear on day 10 and continue through day 16, consistent with previous reports concerning the timing of onset of stress-induced senescence (50) (Fig. 7E, Fig. 16E).

[0177]    SASP, which has key roles in wound healing and development that are relevant for reprogramming biology, includes the expression of various soluble factors (including 116), chemokines (including 118), inflammatory factors (including Ifng), and growth factors (including Vegf) that can promote proliferation and inhibit differentiation of epithelial cells (50). Recent

reports have suggested that secretion of 116 and other soluble factors by senescent cells can enhance reprogramming (51). Although detectable levels of 116 mRNA were present in only a small fraction of cells both in 2i and serum (0.2%) at days 12 and 16 (0.34% in all cells), the overall SASP signature was evident in 72% of cells in the Valley (vs. 11% elsewhere, primarily in day 0 MEFs). This suggests that the senescent cells in the Valley are likely to have paracrine effects on cells that successfully emerge from the Horn.

**Example 6 Other cells at day 4 are strongly biased toward the Horn of Transformation.**

**[0178]**    For the remaining cells at day 4, the forward trajectory is characterized by high proliferation and loss of MEF identity (Figs. 7B, 7E), and the descendants are strongly biased toward the Horn at day 8 (Figs. 8A, 8B and Fig. 17A and Table 10). The Horn is distinguished as a point of transformation, where cells that have lost their mesenchymal identity are beginning their transitions to an epithelial fate. As discussed below, a minority of cells in the Horn have begun to express activators of a pluripotency expression program.

**[0179]**    Following Dox withdrawal and media replacement on day 8, the cells in the Horn adopt one of four alternative outcomes by day 12 (senescence, neuronal program, placental program, and pre-iPSCs). Roughly half appear to become senescent, migrating through clusters 19 and 10 to the Valley (Fig. 8A). The fate of the remaining cells is strongly influenced by the culture medium. In serum conditions, the proportion of these cells that transition to neuronal, placental and pre-iPSC states is 62%, 13% and 26%, respectively. By contrast, the proportions in 2i condition are 3%, 37% and 59% (Table 10). These results are consistent with the presence in the 2i medium of two small-molecule inhibitors to inhibit differentiation, including one reported to inhibit neuronal differentiation (52).

**Example 7**

**[0180]**    Neuronal-like and placental-like cells arise during reprogramming.

**[0181]**    Two unusual cell populations were analyzed: placental-like cells (clusters 24 and 25, Figs. 7B, 7D and Figs. 8A, 8B, 8E, 8F) at day 12 and neural-like cells (clusters 26 and 27, Figs. 7B, 7D and Figs. 8A, 8B, 8E, 8F) at day 16. The first group was characterized by high activity of two gene modules enriched in signatures for "epithelial cell differentiation," "placenta development," and "reproductive structure development," while the second group showed high

113

activity of signature for "neuron differentiation," "axon development," and "regulation of nervous system development" (Table 1, and Figs. 7B, 8C, 8E).

**[0182]**     Both populations showed a substantial decrease in proliferation (Fig. 7E, fig. 16E). To explore if a common mechanism was responsible for this change, 98 cell-cycle related genes (53) were examined to identify those that were differentially upregulated in the placenta and neural clusters compared to all other clusters. The most distinctive characteristic was the high expression of Cdknlc, which encodes a cell-cycle inhibitor (p57) that promotes G1 arrest (Fig. 7F) and is required for maintenance of some adult stem cells (54). Other features are also shared between these two alternative lineages and adult stem cells-including the expression of Lgr5, a marker of adult epithelial stem cells in certain tissues (55) (Fig. 15).

**[0183]**     The neural-like cells reside in a large "spike" observed at day 16 in serum but not 2i conditions (16% vs. 0.1% of cells), presumably due to differentiation inhibitors in the latter conditions. Cells near the base of the spike (cluster 26, Fig. 7D and Figs. 8E, 8F) expressed neural stem-cell markers (including Pax6 and Sox2, Fig. 7E, fig. 15), while cells further out along the spike (cluster 27, Fig. 7D) expressed markers of neuronal differentiation (including Neurog2 and Map2, fig. 15). The cells thus appear to span multiple stages of neurogenesis along the length of the spike (Fig. 7E).

**[0184]**     Analysis of the developmental landscape suggests a potential mechanism for triggering neural differentiation. The ancestors of neural-like cells are largely found in cluster 23 on day 12 (Figs. 8A, 8F and fig. 17C and Table 10). At least 19% of cells in cluster 23 express Cntfr, an 116-family receptor that plays a critical role in neuronal differentiation and survival (56) (Fig. 7F); the true proportion is likely to be higher because the gene has low expression. Contemporaneously, senescent cells in the Valley at day 12 express activating ligands (Crlfl and Clcfl) of Cntfr (fig. 15). Thus, neural differentiation may be triggered by paracrine signals from senescent cells to Cntfr-expressing cells.

**[0185]**     The placental-like cells express high levels of certain imprinted genes on chromosome 7 (Cdknlc, Igf2, Peg3, H19 and Ascl2; Fig. 7F, Fig. 15), as well as TFs (Cdx2 and Soxl7) associated with placental development (57, 58) (Fig. 15). They also show elevated levels of an ER stress signature (Fig. 3E), consistent with the secretory nature of placental cells and observations of placental cells in vivo (59). Analysis was performed to address whether the

placental-like cells resembled recently described extraembryonic endodermal (XEN) cells from an iPSC reprogramming study (44). It was found that they do not share the distinctive XEN signature of the cells disclosed in that analysis. The proportion of cells in the placental-like population decreased substantially from day 12 to day 16 in 2i conditions, although the optimal-transport analysis could not confidently infer whether the decrease is due to cells dying, being overtaken by faster-growing cells, or transitioning to other fates (fig. 14A).

[0186]     The following two tables provide a list of candidate reprogramming factors.

**Example 8**

**Trajectory to successful reprogramming reveals a continuous program of gene activation.**

[0187]     We next studied the trajectory leading to reprogramming (Figs. 8D, 8E), which passes through pre-iPSCs (cluster 28; Figs. 8A, 8B) at day 12 en route to iPSC-like cells at day 16. The iPSC-like cells in serum conditions (which reside in cluster 31) closely resemble fully reprogrammed cells grown in serum (cluster 32). By contrast, the iPSC-like cells under 2i conditions are spread across three clusters (cluster 29-31). While the cells in cluster 31 resemble fully reprogrammed cells grown in 2i (cluster 33), those in cluster 29 show distinct properties suggestive of partial differentiation. In particular, cluster 29 shows lower proliferation, lower Nanog expression, and increased expression of genes related to differentiation (Figs. 7D, 7F).

[0188]     In contrast to initial descriptions of reprogramming as involving two "waves" of gene expression, the trajectory of successful reprogramming reveals a more complex regulatory program of gene activity (Fig. 9A). By grouping genes according to their temporal patterns of activation in cells on the OT-defined trajectory to successful reprogramming, a rich collection of markers for particular stages can be obtained **(Fig. 9A).** In particular, 47 genes that appear late in successfully reprogrammed cells (for example, *Obox6, Spic, Dppa4)* were identified. These genes may provide useful markers to enrich fully reprogrammed iPSCs **(Table 2).**

**Example 9**

**Paracrine signaling from the Valley may influence late stages of reprogramming.**

[0189]     The simultaneous presence of multiple cell types raises the possibility of paracrine signaling, with secreted factors from one cell type binding to receptors on another cell type. One

such potential interaction above, is SASP+ cells in the Valley secreting Crlfl, Clcfl and neural-like cells on days 12 and 16 expressing the cognate receptor Cntfr.

[0190] To systematically identify potential opportunities for paracrine signaling, we defined an interaction score, $I_{A,B,X,Y,t}$, as the product of (1) the fraction of cells in cluster A expressing ligand X and (2) the fraction of cells in cluster B expressing the cognate receptor Y, at time t. Using a curated list of 149 expressed ligands and their associated receptors, we studied potential interactions between all pairs of clusters for each ligand-receptor pair, as well as the aggregate signal across all pairs and across those pairs related to the SASP signature. The potential for paracrine signaling varied sharply across the time course, as well as across cell types. Potential interactions are initially high, as cells with MEF identity retain their secretory functions; drop dramatically by day 6 (Fig. 18A), after cells have lost their MEF identity (Fig. 7B, 7C, 7E); rise steadily from day 8 to day 11, as secretory cells in the Valley emerge; and then drop again from days 12 to 16, as the abundance of cells in the Valley decreases (Fig. 18A). The same pattern is seen when considering only the 20 ligands in the SASP signature (Fig. 18B).

[0191] Notably, potential interactions are observed between cells in the Valley and each of iPSC, neural-like and placental-like cells. At day 16, cells in the Valley (clusters 15 and 16) express SASP ligands, while iPSCs (clusters 29-33) express receptors for these ligands (Fig. 18C), with the highest frequency seen for the chemokine Cxcll2 and receptor Dpp4 (Fig. 18D). As noted above, at days 12 and 16, the ligands Crlfl and Clcfl cells are expressed in the Valley while their receptor Cntfr is expressed in the neural spike (Fig. 7E, Fig. 18E). The interaction between Cntfr and Crlfl is ranked as the top interaction among all ligand-receptor pairs (Fig. 18E).

[0192] At day 12, many placental-like cells express the ligand Igf2 while cells in the Valley express receptors *Igflr and Igf2r* (**Fig. 18F).**

**Example 10**

**X-chromosome reactivation follows activation of early and late pluripotency genes.**

[0193] The reversal of X-chromosome inactivation in female cells is known to occur in the late stages of reprogramming and is an example of chromosome-wide chromatin remodeling. A recent study (60) reported that X-reactivation follows the activation of various pluripotency genes, based on immunofluorescence and RNA FISH in single cells. To assess X-reactivation,

116

from scRNA-Seq data, each cell was characterized with respect to signatures of X-inactivation (Xist expression), X-reactivation (proportion of transcripts derived from X-linked genes, normalized to cells at day 0), and early and late pluripotency genes. Along the trajectory to successful reprogramming (but not elsewhere, **Fig. 7E),** cells at day 12 show strong downregulation of *Xist* but do not yet display X-reactivation. X-reactivation is complete at day 16, with the signature having risen from 1.0 to -1.6, consistent with the expected increase in X-chromosome expression (61). Analysis of the trajectory confirms that activation of both early and late pluripotency genes precedes *Xist* downregulation and X-reactivation.

**Example 11**

**Some cell populations are enriched for aberrant genomic events.**

**[0194]**    Anaylsis was done to identify other coherent increases or decreases in gene expression across large genomic regions, which might indicate the presence of copy-number variations (CNVs) in specific cells. Particularly, analysis done to identify whole chromosome aberrations, demonstrated that 0.9% of cells showed significant up- or down-regulation across an entire chromosome; the expression-level changes were largely consistent with gain or loss of a single chromosome.

**[0195]**    Next, evidence of large subchromosomal events was identified by analyzing regions spanning 25 consecutive housekeeping genes (median size -25 Mb). Significant events were found in -0.8% of cells. The frequency was highest (2.8%) in cluster 14, consisting of cells in the Valley of Stress enriched for a DNA damage-induced apoptosis signature. The frequency was 2-to-3-fold lower in other cells in the Valley (enriched for senescence but not apoptosis), in cells en route to the Valley (clusters 8 and 11), and in fibroblast-like cells at days 0 and 2. Notably, it was much lower (6-fold) in cells on the trajectory to successful reprogramming (Figs. 22B, 22C). Direct experimental evidence would be needed to confirm these events, and to clarify if the aberrations were preexisting in the MEF population, or if they accumulated during the course of reprogramming. 9

**Example 12**

**Inferred trajectories agree with experimental results from cell sorting.**

**[0196]**    To test the accuracy of the probabilistic trajectories calculated for each cell based on optimal transport, results based on the trajectories were compared to experimental data from a

recent study of reprogramming of secondary MEFs (16). In that study, cells were flow-sorted at day 10, based on the cell-surface markers CD44 and ICAM1 and a Nanog-EGFP reporter gene, and each sorted population was grown for several days thereafter to monitor reprogramming success. Gene expression profiles were obtained from each population at day 10 and CD44-ICAMl+Nanog+ population at day 15, together with mature iPSCs and ESCs. Reprogramming efficiency was lowest for CD44+ICAM-Nanog- cells, intermediate for CD44-ICAMl+Nanog- and CD44-ICAMl-Nanog+ cells, and highest for CD44-ICAMl+Nanog+ cells.

[0197]    The flow-sorting-and-growth protocol was emulated in silico, by partitioning cells based on transcript levels of the same three genes at day 10 and predicting the fates of each population at day 16 based on the inferred trajectory of each cell in the optimal transport model. The computational predictions showed good agreement with these earlier experimental results (Fig. 5B), with respect to both reprogramming efficiency and changes in gene-expression profiles. In particular, the *in silico* results showed 93% correlation with results from the earlier study concerning relative reprogramming efficiencies for six categories of sorted cells (p value= 0.0023) **(Fig. 9B***).* Notably, the computationally inferred trajectory of double positive cells rapidly transitioned toward iPSCs and continued in this direction through the end of the time course **(Fig. 9B).** Only one category (CD44-ICAM+Nanog-) differed significantly.

[00138]    Differences may reflect the fact that experimental protocols were not identical *(e.g.,* the earlier study (16) maintains continuous expression of OSKM and supplements the medium with an ALK-inhibitor and vitamin C).

**Example 13**

**Inferring transcriptional regulators that control the reprogramming landscape.**

[0198]    The optimal transport map provides an opportunity to infer regulatory models, based on association between TF expression in ancestors and gene expression patterns in descendants. TFs were identified by two approaches **(Fig. 9C):** (i) a global regulatory model, to identify modules of TFs and target genes and (ii) enrichment analysis, to identify TFs in cells having many vs. few descendants in a target cell population of interest. Gene regulation along the trajectories to placental-like and neural-like cells was examined **(Fig. 19).** For placental-like cells, the analysis pointed to 22 TFs **(Figs. 19A, 19B** and **Table 3).** Of the four most enriched *(Pparg, Cebpa, Gcml,* and *Gata2),* all have been reported to play roles in placenta development

(62). For example, *Gcml* was detected in 42% of cells at day 10 with a high proportion (>80%) of descendants in the placental-like fate but only 0.7% of those cells with a low proportion (<20%) (57-fold enrichment). For neural-like cells, the analysis pointed to 10 TFs *(Pax3, Msxl, Msx3, Sox3, Soxll, Tall, Enl, Foxa2, Gbx2,* and *Foxbl)*. All have been implicated in various aspects of neural development **(fig. *19C) (62-70).***

**[0199]**    Additional analysis focused on identifying TFs that play roles along the trajectory to successful reprogramming (Fig. 9D and fig. 19D, 19E). The global regulatory model generated two regulatory modules, A and B, with 61 TFs in module A, 16 in module B, and 11 in both (Figs. 19D, 19E).

**[0200]**    Module A involves target genes active across clusters 29-31, while Module B involves target genes that are more active in cluster 31, which contains more fully reprogrammed cells. The TFs in these modules are progressively activated across the trajectory of successful reprogramming. For Module B, the TFs are active in 13% of cells in the Horn on day 8, while target-gene activity is evident (at >80% of the levels observed in iPSCs) in 1.3%, 10%, and 21% of their descendant cells in days 10, 11, and 12 in 2i conditions; the pattern in serum conditions is similar, although with lower overall frequency (11%) of cells by day 12). The onset of TFs and target genes in Module A lags by 1-2 days (Fig. 9D).

**[0201]**    To identify TFs likely to play a key role in the final stages of reprogramming, we used enrichment analysis to identify TFs enriched in cells at day 12 with a high vs. low proportion (>80% vs. <20%) of successfully reprogrammed descendants and then focused on the intersection of this set with the 66 TFs from the global regulatory analysis above. The analysis pointed to 9 TFs associated with a high probability of success in the late stages of reprogramming (Fig. 19F). Of these, five (Sox2, Nanog, Hesxl, Esrrb, Zfp42) have establishedroles in regulation of pluripotency *(71-73),* while the remaining four *(Obox6, Spic, Mybl2,* and *Msc)* have not previously been implicated. Among these novel factors, *Obox6* stands out as having the greatest enrichment in high- vs. low-probability cells (68-fold, 9.3% vs -0.14%) **(fig. 19F).**

**Example 14**

**Forced expression of *Obox6* enhances reprogramming.**

[0202]     Obox6 was identified by the regulatory analysis described herein as strongly correlating to reprogramming success. Obox6 (oocyte-specific homeobox 6) is a homeobox gene of unknown function that is preferentially expressed in the oocyte, zygote, early embryos and embryonic stem cells (74).

[0203]     To test whether Obox6 also plays an active role in the process of reprogramming, experiments were performed to address whether expressing Obox6 along with OKSM during days 0-8 can boost reprogramming efficiency. Secondary MEFs were infected with a Dox-inducible lentivirus carrying either Obox6, the known pluripotency factor Zfp42 (73), or no insert as a negative control. Both Obox6 and Zpf42 increased reprogramming efficiency of secondary MEFs by ~2-fold in 2i and even more so in serum. The results were confirmed in multiple independent experiments (Figs. 10A and 10B, and fig. 20). Assays in primary MEFs showed similar increases in reprogramming efficiency (fig. 20). These results demonstrate the importance of Obox6 in the context of cellular reprogramming.

[0204]     Figs. 1OA-IOC demonstrate the effect of overexpression of Obox6 and Zpf42 on reprogramming efficiency in secondary MEFs. Figs. 10 A and 10B show bright field and fluorescence images of iPSC colonies generated by lentiviral overexpression of *Oct4, Kl/4, Sox2,* and *Myc* (OKSM) with either an empty control, *2fp42* or *Obox6* expression cassette, in either Phase-l(Dox)/Phase-2(2i)(A)   and Phase-l(Dox)/Phase-2(serum)  (B) conditions (indicated). Cells were imaged at day 16 to measure Oct4-EGFP+ cells. Bar plots representing average percentage of Oct4-EGFP+ colonies in each condition on day 16 are included below the images. Shown are data from one of five independent experiments, with three biological replicates each. Error bars represent standard deviation for the three biological replicates. Figure 6C is a schematic of the overall reprogramming landscape highlighting: the progression of the successful reprogramming trajectory, alternative cell lineages, and specific transition states (Horn of Transformation). Also highlighted are transcription factors (orange) predicted to play a role in the induction and maintenance of indicated cellular states, and putative cell-cell interactions between contemporaneous cells in the reprogramming system.

**Example 15**

**Definition of gene signatures**

[0205]    From gene set enrichment analysis of 44 gene modules (Table 1, Figs. 12A-12C), significant enrichments for terms that shed light on the reprogramming landscape were found. Analysis was done to investigate whether similar expression patterns from well-defined gene signatures could be identified. To investigate this, a list of gene sets from various databases of gene signatures was curated (see Table 11, a list of genes for each gene signature is shown in Table 2). A pluripotency gene signature was determined.

[0206]    Differential gene expression analysis was performed between two groups of cells: mature iPSCs and cells along the time course D0 to D16, and the top 100 genes with increased expression in mature iPSCs were identified. A proliferation gene signature was obtained by combining genes expressed at Gl/S and G2/M phases. For epithelial and neural gene signatures, canonical markers of epithelial and neuronal cell lineage markers, respectively were collected.

Table 11.: List of gene signatures used in this work. List of genes for each gene signature are shown in Table 2.

| Gene Signature | Source |
|---|---|
| MEF identity | Mouse Gene Atlas (S29, S30) |
| Pluripotency | this work, iPSCs vs. D0 to D16 cells |
| Proliferation | G1/S and G2/M genes, (S31) |
| ER stress | GO:0034976, Biological Process Ontology |
| Epithelial identity | (S32–S35) |
| ECM rearrangement | GO:0030198, Biological Process Ontology |
| Apoptosis | Hallmark P53 Pathway, MSigDB |
| Senescence | Table 1 in (S36) |
| Neural identity | (S37-S43) |
| Placental identify | Mouse Gene Atlas, (S29, S30) |
| X reactivation | chromosome X |

**Computing descendant distributions for clusters of cells**

[0207]    The descendant distributions for the 33 clusters of cells, some of which span multiple days were computed. To put each cluster on equal footing, 100 cells in each cluster were initialized. These 100 cells were distributed proportionally over the days represented in the cluster.

For each day $d$ and cluster $i$, let $n^i_d$ denote the number of day $d$ cells in cluster $i$. We denote the total number of cells in cluster %by $N^i = \Sigma_d n^i_d$. With this notation, we initialize $100 \times \frac{n^i_d}{N^i}$ cells in cluster $i$ on day $d$ and compute the descendant distribution of these cells at the next time point. We denote this descendant distribution by $D_d$. We then compute the mass of this descendant distribution residing in each cluster $j$ by summing up the mass $D_d$ assigns to each cell in cluster $j$. Finally, to obtain the $i,j$ entry of the cluster - cluster transition table, wc sum over $d$.

This give the total mass transferred from from cluster $i$ to cluster $j$, per 100 cells initialized in cluster $i$. We compute this separately for 2i and serum.

## Extraembryonic gene signatures

**[0208]**    Previous reports have shown that extraembryonic endoderm stem cells (XEN) were induced in the reprogramming process in parallel of reprograming to iPSCs (S48). To determine if XEN cells were induced in the reprogramming system described herein, the XEN gene signature from in vivo XEN cells, trophoblast and placental gene signatures was analyzed (**Table 12**). While a small fraction of cells (180 cells) displays a high XEN score at day 16 (under serum condition), a larger fraction of cells in clusters 24 and 25 displays high trophoblast and placental signature scores . This indicates that the alternative placental-like cell lineage does not share the distinctive XEN signature as previously reported.

| Gene .Signature | Genes | Reference |
|---|---|---|
| XEN | Dab2 Fst Pdgfra Pthlr Gata6 Foxql Fxyd3 Tet3 Sox 17 Foxa2 Lama l Lamb l Gata4 Krt8 | (S49) |
| Trophoblast | Ascl2 Bmp4 Bmp8b Cdx2 Elf5 Eomes Esrrb Els2 Fgfr2 Grn Tgf2 .Tade l Lipg Pcsfc6 Ptpra Smad3 Snai l Tead4 Tfap2c Vavl Yap l Gata3 Krt7 Krtl 8 | (S50) |
| Placenta | Table A 1 | |

**Table 12.: List of XEN, trophoblast and placenta gene signatures**

## Example 16

## Identifying markers for reprogramming success

**[0209]**    To gain further insights into the mechanisms of reprogramming success, categories of genes that changed their expression in characteristic patterns (Figs. 5A-5G) along the successful trajectory determined by optimal transport were characterized. Genes that exhibited significant changes along the trajectory (2,872 genes) were clustered using k-means clustering and the

number of clusters was determined by the gap statistic (S44). 14 distinct expression patterns among cells that would end up succesfully reprogrammed (Table 10) were identified. Genes were divided into two obvious patterns, upregulated (Al to A10) and downregulated (Al l to A14). After dox induction, a large number of genes that were mainly involved with MEF identify were downregulated. Instead of "two waves" indicated by a previous report (S45), continuous activation patterns after dox induction were observed. In early stage of reprogramming, they were involved with metabolic changes and were targets of Myc (Al to A3). In late stage (A6 and A7) they were associated with activation of pluripotency networks. Two categories of pluripotency-associated genes were identifed. Genes in category A6 gradually upregulated after dox withdrawal, such as Nanog, Sox2, Dppa3 (early pluripotency-associated genes). Genes in category A7 upregulated after genes in A6, such as Obox6, Dppa4 (late pluripotency-associated genes).

[0210]    Genes that were upregulated preferentially in cells that were successfully reprogrammed from A6 and A7 were identifed. The fraction of cells in clusters 28 to 33 vs. all other clusters were calculated. By setting a threshold of 1%, genes that were expressed in less than 1% of cells in all other clusters were ranked. 47 genes that were preferentially expressed in the late stage of reprogramming on successful trajectory and were mostly absent from other cells (Table 10) were identified.

**Example 17**

**Cell-Cell Interactions**

[0211]    To characterize potential cell-cell interactions between contemporaneous cells during reprogramming, a list of ligands and receptors found in the GO database were collected. The set of ligands (415 genes) is a union of three gene sets from the following GO terms: 1) cytokine activity (GO:0005125), 2) growth factor activity (GO:0008083), and 3) hormone activity (GO:0005179). The set of receptors (2335 genes) is defined by the GO term receptor activity (GO:0004872). Next, a curated database of mouse protein-protein interactions (S46) was used to identify 580 potential ligand-receptor pairs. Two aspects of potential cell-cell interactions in the data were the focus of the analysis: 1) determining global trends in the expression of all potential contemporaneous ligand-receptor pairs across the reprogramming time course and 2) ranking individual ligand-receptor pairs at a specific day and condition. First, an interaction score

I$\lambda$,B$\chi\gamma$,i as the product of (1) the fraction of cells $(F_{A,X,t})$ in cluster A expressing ligand X at time

$t$ and (2) the fraction of cells $(F_{B,Y,t})$ in cluster B expressing the cognate receptor Y at time t was

defined. Aggregate interaction score $I_{A,B,t}$ was defined as a sum of the individual interaction

scores across all pairs:

$$I_{A,B,t} = \sum_{All\ X\text{-}Y\ pairs} I_{A,B,X,Y,t} = \sum_{All\text{-}X\text{-}Y\ pairs} F_{A,X,t} F_{B,Y,t}$$

[0212]    The aggregate interaction scores for all combinations of cell clusters in **figs. 18A-B**

were depicted. Second, individual ligand-receptor pairs at a given day and condition between cell

subsets of interest were examined. Values of the interaction scores $I_{A,B,X,Y,t}$ are high for

ubiquitously expressed ligands and receptors at a given day and may be nonspecific to a pair of

cell subsets of interest. Thus, permutations were used to generate an empirical null distribution of

interaction scores between two random groups of cells. In each of the 10,000 permutations, two

groups R1 and R2 of 100 cells each from time t were selected and the interaction score between

the ligand in group R1 and the receptor in group R2 was calculated. Each ligand-receptor

interaction score was standardized by taking the distance between the interaction score $I_{A,B,x,Yt}$

and the mean interaction score in units of standard deviations from the permuted data $((I_{A,B,X,Y,t}$

$— mean(I_{Ri,R2},x,Y,t))/sd(I_{Ri,R2,xj,t}))$. Examples of standardized interaction scores ranked by their

values are depicted in **Figs. 18D-F.**

**Example 18**

**X-chromosome reactivation**

[0213]    Analysis was performed to identify X-chromosome reactivation from our scRNA-seq

dataset. The set of all detected genes (16,339) was split to X-chromosomal and autosomal genes.

Then the mean X/autosome expression ratio for each cell (normalized by the average

X/autosome expression ratio at day 0 cells) as a measurement of X-chromosome reactivation was

calculated.

[0214]    The mean X/Autosome expression ratio reached mean value of 1.6 in late stage of

reprogramming indicating X-chromosome reactivation . Interestingly, cells in cluster 32 (mature

iPSCs in serum) had their X-chromosome inactivated but no Xist expression, which might be

due to partial differentiation of iPSCs in serum condition or that the established female iPSCs lost one of their X chromosomes, which happens frequently in serum cultured female ESCs or iPSCs but less often in 2i cultured female ESCs/iPSCs (S47). This was specific to mature iPSCs in serum as day-16 cells in serum exhibited similar X-chromosome reactivation to day 16 cells in 2i

[0215] Downregulation of Xist expression (cluster 28, day 12 cells) preceded X-chromosome reactivation (clusters 29,30,31, and 33; day 16, mature iPSCs) (Figs. 21A-21C). The upregulation of early and late pluripotency genes (activation pattern A6 and A7, respectively) preceded X-chromosome reactivation (Figs. 21D-21F).

[0216] The fraction of cells that activated late pluripotency genes A7 and reactivated the X-chromosome were analyzed. The X/Autosome expression ratio and A7 gene signature score show bimodal distribution across all cells (fig. 21G and fig. 21H, respectively). We classified cells to those that had reactivated their X-chromosome if the X/Autosome expression ratio > 1.4 and those that induced A7 genes if the A7 average z-score > 0.25 **(figs. 21G, 21H).** Using the above thresholds the fraction of cells in clusters 28-33 that reactivated their X-chromosome and activated the A7 program **(Table 13) were calculated.** Around a 10-fold difference is observed in the percentage of cells that upregulated A7 genes and reactivated X chromosome in clusters 28 and 32.

| **Cluster** | 28 | 29 | 30 | 3 1 | 32 | 33 |
|---|---|---|---|---|---|---|
| X/A | 7.6 | 79.3 | 84.2 | 89.1 | 7.2 | 81.9 |
| A7 | 72,9 | 98.9 | 99.7 | 99.1 | 93.3 | 99. 1 |

Table 13. Percentage of cells in clusters 28-33 that exhibited X-chromosome reactivation and induction of A7 genes.

**Example 19**

**Identifying large chromosomal aberrations**

[0217] Methodology. Two types of analysis were performed to detect aberrant expression in large chromosomal regions. First, analysis was performed to identify cells with significant up- or down-regulation at the level of entire chromosomes. Second, analysis was performed to identify cells with significant subchromosomal aberrations spanning windows of 25 consecutive broadly-

expressed genes. Empirical p-values and false discovery rates (FDRs) for both analyses were computed by randomly permuting the arrangement of genes in the genome, as described below.

[0218]    Permutations for both types of analysis are done as follows. In each of 100,000 permutations the labels of genes in the entire dataset were randomly shuffled, while preserving the genomic positions of genes (with each position having a new label each time) and the expression levels in each cell (so that each cell has the same expression values, but with new labels). Either whole chromosome or subchromosomal aberration scores for each cell were calculated. To identify whole-chromosome aberrations scores in each cell, the sum of expression levels in 25Mbp sliding windows along each chromosome, with each window sliding IMbp so that it overlaps the previous window by 24Mbp was calculated. For each window in each cell, the Z-score of the net expression, relative to the same window in all other cells was calculated. The fraction of windows on each chromosome with an absolute value Z-score > 2 was counted. This fraction serves as the whole-chromosome aberration score for each chromosome in each cell. To assign a p-value to the whole-chromosome score for cellj chromosomej, the empirical probability that the score for cellj chromosomej in the randomly permuted data was at least as large as the score in the original data was calculated.

[0219]    Subchromosomal aberration scores were computed as follows. The 20% of genes with the most uniform expression across the entire dataset were identified. This is done by calculating the Shannon Diversity (eentropy(gene)) for each gene, and taking the 20% of genes with the largest values. Using these genes, the sum of expression in sliding windows of 25 consecutive genes, with each window sliding by one gene and overlapping the previous window (on the same chromosome) by 24 genes was calculated. In each window, the Z-score relative to all cells at day 0 was calculated. The net subchromosomal aberration score for a cell is calculated as the 12-norm of the Z-scores across all windows. To assign a p-value to the subchromosomal aberration score for celli, the empirical probability that the score for celli in the randomly permuted data was at least as large as the score in the original data was calculated.

[0220]    For subchromosomal aberration scores chromosomal aberrations (vs. locally coordinated programs of gene expression) were enriched for by excluding recurrent events. Recurrent events were identified by clustering cells based on their aberration profiles (net expression levels across all windows). Clustering was completed by calculating the SVD of all

aberration profiles, and performing KMeans clustering on the the top 10 singular vectors (with k=100). For each cluster, we quantified cluster compactness and separation using the silhouette score. Cells that were in compact, well-separated clusters (with a silhouette score > 0.08) were removed from consideration for subchromosomal aberrations.

[0221]    For both types of scores, p-values were used to calculate false discovery rates (FDRs). To identify cells with aberrations at an FDR of q, the largest p-value, $\hat{p}$ was identified, such that pN/sum(p< p), where N represents the total number of p-values for a score and sum(p< $\hat{p}$) represents the number of p-values less than $\hat{p}$.

[0222]    Since recurrent aberrations are expected in this setting (due to clonal expansion) cells based on clustering recurrent patterns were not removed. Applied to these data, this method detected aberrations in 35% of malignant cells (classified in the original study as containing significant copy number variation) and 0% of non-malignant cells (FDR 5%). This demonstrates the specificity and conservative nature of the approach.

[0223]    Results. The results of this analysis are displayed in Figs. 22A-22C. In analysis designed to look for whole chromosome aberrations, it was found that 0.9% of cells showed significant up- or downregulation across an entire chromosome; the expression-level changes were largely consistent with gain or loss of a single chromosome (AHA). Next, analysis performed to look for evidence of large subchromosomal events, found significant events in 0.8% of cells. The frequency was highest (2.8%>) in cluster 14, consisting of cells in the Valley of Stress enriched for a DNA damage-induced apoptosis signature. The frequency was 2-to-3-fold lower in other cells in the Valley (enriched for senescence but not apoptosis), in cells en route to the Valley (clusters 8 and 11), and in fibroblast-like cells at days 0 and 2. Notably, it was much lower (6-fold) in cells on the trajectory to successful reprogramming (Figs. 22B, 22C). Direct experimental evidence would be needed to confirm these events, and to clarify if the aberrations were preexisting in the MEF population, or if they accumulated during the course of reprogramming.

Example 20

[0224]    Forced expression of transcriptional regulators enhances reprogramming.

[0225]    To test whether any of the transcriptional regulators provided in Tables 2, 3 and 4, for example, Obox6, Spic, Zfp42, Sox2, Mybl2, Msc, Nanog, Hesxl and Esrrb, play an active role

in the process of reprogramming, experiments are performed to address whether expressing these transcription regulators along with OKSM during days 0-8 can boost reprogramming efficiency. Secondary MEFs or primary MEFS are infected with a Dox-inducible lentivirus carrying any one of the transcription regulators provided in Tables 2, 3 and 4, the known pluripotency factor Zfp42 (73), or no insert as a negative control. Reprogramming efficiency is assessed in 2i or in serum. Multiple independent experiments are performed. An increase in reprogramming efficiency by a transcriptional regulator identifies the regulator as important in the context of cellular reprogramming.

[0226]    Reprogramming efficiency is assessed by analyzing bright field and fluorescence images of iPSC colonies generated by lentiviral overexpression of Oct4, Klf4, Sox2, and Myc (OKSM) with either an empty control, Zfp42 or an expression cassette for any one of the transcription regulators provided in Tables 2, 3 and 4, in either Phase-l(Dox)/Phase-2(2i)(A) and Phase-l(Dox)/Phase-2(serum). Cells are imaged at day 16 to measure Oct4-EGFP+ cells. Bar plots representing average percentage of Oct4-EGFP+ colonies in each condition on day 16 are generated. Error bars represent standard deviation for biological replicates.

**Example 20**

[0227]    Reconstruction of developmental landscapes by optimal-transport analysis of single-cell gene expression across time sheds light on reprogramming

[0228]    Here, we introduced Waddington-OT, a new approach for studying developmental time courses to infer ancestor-descendant fates and model the regulatory programs that underlie them. We applied Waddington-OT to reconstruct the landscape of reprogramming from 315,000 scRNA-seq profiles, collected mostly at half-day intervals across 18 days. We revealed a wider range of developmental programs than previously recognized. Cells gradually adopted either a terminal stromal state or a mesenchymal-to-epithelial transition state. The latter gave rise to populations related to pluripotent, extra-embryonic, and neural cells, with each harboring multiple finer subpopulations. We predicted transcription factors controlling various fates, of which we showed that Obox6 enhanced reprogramming efficiency. We also found rich potential for paracrine signaling. Our approach shedded new light on the process and outcome of reprogramming and provided a framework applicable to diverse temporal processes in biology.

[0229]    In the mid-20th century, Waddington introduced two metaphors that shaped biological thinking about cellular differentiation during development: first, trains moving along branching railroad tracks and, later, marbles following probabilistic trajectories as they roll through a developmental landscape of ridges and valleys (Waddington, 1936, 1957). Empirically reconstructing and studying the actual landscapes, fates and trajectories associated with cellular differentiation and de-differentiation — such as in organismal development, long-term physiological responses, and induced reprogramming — requires general approaches to answer questions such as: What classes of cells are present at each stage? What was their origin at earlier stages? What are their likely fates at later stages? What genetic regulatory programs control their dynamics? To what extent are events synchronous vs. asynchronous? To what extent are they stochastic vs. deterministic? Is there only a single path to a given fate, or are there multiple developmental paths?

[0230]    Traditional approaches based on bulk analysis of cell populations were not well suited to addressing these questions, because they did not provide general solutions to two challenges: discovering the cell classes in a population and tracing the development of each class. Progress had historically relied on ad hoc approaches for each question asked (e.g., sorting and following the development of a particular cell class by using an antibody to a class-specific cell-surface protein or a reporter construct).

[0231]    The first challenge has recently been largely solved by the advent of single-cell RNA-Seq (scRNA-Seq) (Klein et al., 2015; Kumar et al., 2014; Macosko et al., 2015; Ramskold et al., 2012; Shalek et al., 2013; Tanay and Regev, 2017; Tang et al., 2009; Wagner et al., 2016), which allowed cell classes to be discovered based on their expression profiles. The second challenge remained a work-in-progress. ScRNA-seq now offered the prospect of empirically reconstructing developmental trajectories based on snapshots of expression profiles from heterogeneous cell populations undergoing dynamic transitions (Bendall et al., 2014; Marco et al., 2014; Setty et al., 2016; Tanay and Regev, 2017; Trapnell et al., 2014; Wagner et al., 2016). But, to trace the trajectories of cell classes, one may connect the discrete 'snapshots' produced by scRNA-Seq into continuous 'movies.' At least at present, one may not be able to follow expression profiles of the same cell and its direct descendants across time because current methods may destroy cells to profile their state. While various approaches have been developed to record information about

cell lineage, they currently provide only very limited information about a cell's state at all earlier time points (Daniel T. Montoro et al., 2018; Kester and van Oudenaarden, 2018; McKenna et al., 2016).

[0232]    Comprehensive studies of cell trajectories thus relied heavily on computational reconstruction of paths in gene-expression space. Pioneering work introduced various methods to infer trajectories (Bendall et al., 2014; Cannoodt et al., 2016; Haghverdi et al., 2015; Matsumoto and Kiryu, 2016; Qiu et al., 2017; Rashid et al., 2017; Rostom et al., 2017; Setty et al., 2016; Street et al., 2017; Trapnell et al., 2014; Weinreb et al., 2017; Welch et al., 2016; Zwiessele and Lawrence, 2016). Profiles of heterogeneous populations can provide information about the temporal order of asynchronous processes—enabling cells to be ordered in pseudotime along trajectories, based on their state of differentiation (Bendall et al., 2014). Some approaches used k-nearest neighbor graphs (Bendall et al., 2014) or binary trees (Trapnell et al., 2014) to connect cells into paths. More recently, diffusion maps have been used to order cell-state transitions, by assigning cells to densely populated paths in diffusion-component space (Haghverdi et al., 2015; Haghverdi et al., 2016). Each such path was interpreted as a transition between cellular fates, with trajectories determined by curve fitting and cells pseudotemporally ordered based on the diffusion distance to the endpoints of each path. Recent work has grappled with incorporating branching paths, which were critical for understanding developmental decisions, and have been applied to analyze whole-organism development in zebrafish, frog, and planaria (Briggs et al., 2018; Farrell et al., 2018; Fincher et al., 2018; Plass et al., 2018; Wagner et al., 2018).

[0233]    While these approaches have shed important light on various biological systems, many important challenges remain. First, most methods neither directly modeled nor explicitly leveraged the temporal information in a developmental time course (Weinreb et al., 2017) because they were designed to extract information about stationary processes (such as adult stem cell differentiation or the cell cycle) in which all stages existed simultaneously across a single population of cells. However, with the rapidly decreasing cost of scRNA-Seq, time-courses may soon be commonplace. Second, many methods model trajectoried in the language of graph theory which imposesed strong structural constraints on the model, such as one-dimensional trajectories ("edges") and zero-dimensional branch points ("nodes"). Yet, some biological systems may show a gradual divergence of fates that were not captured well by these models

(Briggs et al., 2018; Farrell et al., 2018; Wagner et al., 2018). Third, few methods were able to account for cellular growth and death during development. One method capable of modeling nonuniform cellular growth rates was Population Balance Analysis (Weinreb et al., 2017). However, this method assumed the population of cells is in equilibrium, and therefore it was not suited for analyzing dynamical systems where the distribution of cells changed over time.

[0234]     One case in point was the challenge of understanding cellular reprogramming—such as converting fibroblasts to induced pluripotent stem cells (iPSCs) or trans-differentiating one mature cell type into another. These non-natural processes involved the transient overexpression of a set of transcription factors (TFs) designed to push a cell out of its current state and toward a new fate, even in the absence of the usual developmental context. Reprogramming had great therapeutic potential, but it still tends to be slow, inefficient, and asynchronous (Takahashi and Yamanaka, 2016). Single-cell analysis of trajectories during reprogramming could shed light on questions such as: What is the full range of cell classes that arise during reprogramming? What are the developmental paths that lead to reprogramming and to any alternative fates? Which cell intrinsic factors and cell-cell interactions drive progress along these paths? To what extent do cells activate normal developmental programs vs. unnatural hybrid programs? Can the programs that are activated provide information about the normal developmental landscape? Can the information gleaned be used to improve the efficiency of reprogramming toward a desired destination?

[0235]     In particular, reprogramming of fibroblasts to induced pluripotent stem cells (iPSCs), as pioneered by Yamanaka (Hou et al., 2013; Shu et al., 2013; Takahashi and Yamanaka, 2006; Yu et al., 2007), has been largely characterized to date by a combination of fate-tracing of cells based on a handful of markers (e.g., Thyl and CD44 as markers of the fibroblast state, and ICAM1, Oct4, and Nanog as markers of successful reprogramming), together with RNA- and chromatin-profiling studies of bulk cell populations (Buganim et al., 2012; Hussein et al., 2014; O'Malley et al., 2013; Polo et al., 2012; Tonge et al., 2014). With limited cellular resolution, the profiling studies have provided only coarse-grained analyses, such as describing two "transcriptional waves," with gain of proliferation and loss of fibroblast identity followed by transient activation of developmental regulators and gradual activation of embryonic stem cell (ESC) genes (Polo et al., 2012). Some studies (Mikkelsen et al., 2008; O'Malley et al., 2013;

Parenti et al., 2016), including from our own group (Mikkelsen et al., 2008), have noted strong upregulation of several lineage-specific genes from unrelated lineages (e.g., neurons), but it has been unclear whether this reflects coherent differentiation of specific cell types or disorganized gene expression (Kim et al., 2015; Mikkelsen et al., 2008). Most studies that used single-cell methods to study genetic reprogramming have involved few genes or few cells (Buganim et al., 2012, Kim et al., 2015). Recently, a study (Zhao et al., 2018) profiled -36,000 cells during chemical reprogramming, but focused only on a single bifurcation separating successful and failed trajectories.

[0236]     Here, we described a framework, implemented in a method called Waddington-OT, that aimed to capture the notion that cells at any time were drawn from a probability distribution in gene-expression space and cells at any time and position within the landscape had a distribution of both probable origins and probable fates **(FIGs. 23A-23F).** It then used scRNA-seq data collected across a time-course to infer how these probability distributions evolved over time, by using the mathematical approach of Optimal Transport (OT). We applied and tested this framework in the context of scRNA-seq data we profiled from more than 315,000 cells, sampled across a dense time course over 18 days under two different reprogramming conditions. We found that reprogramming unleashed a much wider range of developmental programs and subprograms than previously recognized, resulting in multiple large distinct populations of cells related to pluripotent, extraembryonic, neural, and stromal cells, with evidence for large-scale genomic amplifications and deletions in trophoblast-like and stromal-like cells. Within each population, there were subsets with distinct programs associated with specific cell types in vivo, including programs associated with 2-, 4-, 8-, 16-, and 32-cell stage embryos; with several distinct types of trophoblasts and primitive endoderm; with astrocytes, oligodendrocytes, and neurons; and with a wider range of stromal cells than MEFs. Trajectory analysis with Waddington-OT showed that differentiation among these classes occurred gradually, including an early gradual transition to either stroma-like cells or a mesenchymal-to-epithelial transition state, with the latter state serving as the ancestor population of both eventual iPSC-like cells and extraembryonic and neural. These differentiation fates were predicted by various sets of TFs, including well studied factors and others not previously implicated. We tested one TF found by our analysis to be associated with pluripotency and showed that it enhanced reprogramming

efficiency. Finally, we also found evidence for potential paracrine interactions between the stromal cells and other cell types, which may be important cell extrinsic forces in reprogramming, and for genomic aberrations in certain cells types, with different features in stromal cells and trophoblasts.

**[0237]**   **Results**

**[0238]**   <u>Reconstruction of probabilistic trajectories by Optimal Transport</u>

**[0239]**   **A** goal of the study was to learn the relationship between ancestor cells at one time point and descendant cells at another time point: given that a cell has a specific expression profile at one time point, where will its descendants likely be at a later time point and where are its likely ancestors at an earlier time point? To this end, we modeled a differentiating population of cells as a time-varying probability distribution (i.e., stochastic process) on a high-dimensional gene expression space. By sampling this probability distribution Pt at various time points t, we aimed to infer how the differentiation process it modeled evolves over time **(FIG. 23A).** By sampling a large number of cells at a given time point, we approximated the distribution at that time point. However, this alone did not tell us the ancestor or descendant relationships between cells at different time points: Because different cells were sampled at different time points, we lost this temporal coupling of the stochastic process Pt that specified the joint distribution of expression between pairs of time points. In the absence of any constraint on cellular transitions (e.g., if cells may "jump" about gene-expression space arbitrarily rapidly), we could not infer the temporal coupling. But if we assumed that, over sufficiently short time periods, cells could only move relatively short distance, we could infer the temporal coupling by using the classical mathematical technique of optimal transport **(FIG. 23A, Methods).**

**[0240]**   Optimal transport was originally developed by Monge in 1781 to redistribute earth for the purpose of building fortifications with minimal work (Villani, 2008). In the 1940s, Kantorovich generalized it to identify an optimal coupling of probability distributions via linear programming (Kantorovitch, 1958). This classical linear program minimized the total squared distance that earth travels, subject to conservation of mass constraints. Recent work, which added entropic regularization, dramatically accelerated the numerical computation of large-scale optimal transport problems (Chizat et al., 2017; Cuturi, 2013).

[0241]    However, matching cells to their descendants differed in one important aspect: unlike earth or particles, cells can proliferate. We therefore modified the classical conservation of mass constraints to accommodate cell growth and death. In particular, we allowed the mass of cells to grow as cells proliferate and shrink as cells die (STAR Methods). By leveraging techniques from unbalanced transport (Chizat et al., 2017), we automatically learned cellular growth and death rates, initializing with prior estimates from signatures of cellular proliferation and apoptosis (STAR Methods).

[0242]    Using optimal transport, we calculated couplings between consecutive time points and then inferred couplings over longer time-intervals by composing the transport maps between every pair of consecutive intermediate time points. We noted that the optimal-transport calculation (i) implicitly assumed that a cell's fate depended on its current position but not on its previous history (i.e., the stochastic process is Markov) and (ii) captured only the time-varying components of the distribution, rather than processes at dynamic equilibrium. We returned to these points in the Discussion.

[0243]    We defined trajectories in terms of "descendant distributions" and "ancestor distributions" as follows. For any set C of cells at time $t_i$ , its "descendant distribution" at a later time $t_{i+1}$ referred to the mass distribution over all cells at time $t_{i+1}$ obtained by transporting C according to the transport maps (FIG. 23C). Branching events, for example, were revealed by the (potentially gradual) emergence of bimodality in the descendant distribution (FIG. 23C). Conversely, its "ancestor distribution" at an earlier time $t_{i-1}$ was defined as a mass distribution over all cells at time $t_{i-1}$, obtained by transporting C in the opposite direction (that is, as though one "rewinds" time) (FIG. 23D). Shared ancestry between two cell sets at $t_i$ was revealed by convergence of the ancestor distributions (FIG. 23E). The "trajectory from C" referred to the sequence of descendant distributions at each subsequent time point, and the trajectory to C similarly referred to the sequence of ancestor distributions (FIGs. 23C, 23D). For convenience below, we sometimes referred simply to the 'ancestors, 'descendants',  and 'trajectories' of cells. These terms referred to probability distributions over a set of observed cells that served as proxies for the actual ancestors or descendants. In summary, we used the inferred coupling to calculate a distribution over representative ancestors and descendants at any other time. We then

determined the expression of any gene or gene signature along a trajectory by computing the mean expression level weighted by the distribution over cells at each time point.

[0244]    To identify TFs that regulated the trajectory, we inferred regulatory models by sampling cells from the joint distribution given by the couplings. We developed two approaches: one used 'local' enrichment analysis, identifying TFs that were enriched in cells having many vs. few descendants in the target cell population; a second built a global regulatory model, composed of modules of TFs and modules of target genes, to predict expression levels of target gene signatures (FIG. 23F, left) at later time points from expression levels of TFs at earlier time points (FIG. 23F, middle, right).

[0245]    We implemented our approach in a method, Waddington-OT, for exploratory analysis of developmental landscapes and trajectories, including a public software package (STAR Methods). The method included: (1) Performing optimal-transport analyses on scRNA-seq data from a time course, by calculating optimal-transport maps and using them to find ancestors, descendants and trajectories; (2) Inferring regulatory models that drive the temporal dynamics by sampling pairs of cells from the joint distribution specified by the OT couplings; (3) Visualizing the developmental landscape in two dimensions, by using Force-Directed Layout Embedding (FLE) to visualize the graph of nearest neighbor relationships in diffusion component space (Jacomy et al., 2014; Weinreb et al., 2016; Zunder et al., 2015), and (4) annotating the landscape by cell types, ancestors, descendants, trajectories, gene expression patterns, and other features.

[0246]    <u>A dense experimental scRNA-Seq time course of iPS reprogramming</u>

[0247]    To study the trajectories of reprogramming, we generated iPSCs via a secondary reprogramming system (FIG. 24A), which is more efficient than derivation of iPSCs by primary infection (Stadtfeld et al., 2010). We obtained mouse embryonic fibroblasts (MEFs) from a single female embryo homozygous for ROSA26-M2rtTA, which constitutively expresses a reverse transactivator controlled by doxycycline (Dox), a Dox-inducible polycistronic cassette carrying Pou5fl (Oct4), Klf4, Sox2, and Myc (OKSM), and an EGFP reporter incorporated into the endogenous Oct4 locus (Oct4-IRES-EGFP). We plated MEFs in serum-containing induction medium, with Dox added on day 0 to induce the OKSM cassette (Phase-1(Dox)). Following Dox withdrawal at day 8, we transferred cells to either serum-free N2B27 2i medium (Phase-2(2i)) or

maintained the cells in serum (Phase-2(serum)). Oct4-EGFP+ cells emerged on day 10 as a reporter for successful reprogramming to endogenous Oct4 expression (FIGs. 24A, 30G).

[0248]　　　We performed two dense time-course experiments. In the first we collected -65,000 scRNA-seq profiles at 10 time points across 16 days, with samples taken every 48 hours. In the second we profiled -250,00 0 cells collected at 39 time points across 18 days, with samples taken every 12 hours (and every 6 h o urs between days 8 and 9) (FIG. 24A, Methods, Table 14). The density allows us to ensure that the model is fit on a smoothly progressing process, as well as to use some time points as test data for predictions (below). We also collected samples from established iPSC lines reprogrammed from the same MEFs, maintained in either 2i or serum conditions. The two experiments were consistent (STAR Methods). We focused on the second experiment, where we profiled 259,155 cells to an average depth of 46,523 reads per cell (Table 14). After discarding cells with less than 2,000 transcripts detected, we retained a total of 251,203 cells, with a median of 2,565 genes and 9,132 unique transcripts detected per cell.

Table 14 - Summary of single cell sequencing statistics and sample information.

| Sample Name | Estimated Number of Cells | Mean Reads per Cell | Median Genes per Cell | Number of Reads | Valid Barcodes | Reads Mapped Confidently to Transcriptome | Reads Mapped Confidently to Exonic Regions | Reads Mapped Confidently to Intronic Regions | Reads Mapped Confidently to Intergenic Regions | Reads Mapped Antisense to Gene | Sequencing Saturation | Q30 Bases in Barcode | Q30 Bases in RNA Read | Q30 Bases in Sample Index | Q30 Bases in UMI | Fraction Reads in Cells | Total Genes Detected | Median UMI Counts per Cell |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| D0_Dox_C1 | 3495 | 17263 | 2308 | 60336236 | 98 | 62.7 | 66.1 | 10.8 | 5.4 | 4.4 | 17.4 | 97.9 | 90.9 | 95.8 | 97.7 | 92.2 | 16467 | 7421 |
| D0_Dox_C2 | 1125 | 41979 | 3559 | 47227004 | 98 | 64.2 | 67.6 | 10.5 | 4.9 | 4.3 | 30.8 | 97.9 | 90.6 | 96.3 | 97.7 | 92.4 | 15884 | 15756 |
| D0.5_Dox_C1 | 1220 | 65642 | 4258 | 80083266 | 97.9 | 63.4 | 66.9 | 11.3 | 5 | 4.4 | 38.7 | 97.9 | 90.6 | 95.8 | 97.7 | 95.5 | 16658 | 22429 |

| D0.5_Dox_C2 | D1_Dox_C1 | D1_Dox_C2 | D1.5_Dox_C1 | D1.5_Dox_C2 | D2_Dox_C1 | D2_Dox_C2 | D2.5_Dox_C1 | D2.5_Dox_C2 | D3_Dox_C1 | D3_Dox_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 12851 | 6263 | 8318 | 27357 | 48498 | 11247 | 2275 | 5041 | 7728 | 8215 | 18216 |
| 16911 | 15028 | 16161 | 17182 | 15562 | 17003 | 14980 | 15423 | 16143 | 16144 | 17099 |
| 90.3 | 89 | 94 | 91.8 | 78.5 | 92.5 | 87.1 | 92.7 | 95.6 | 94.4 | 93.5 |
| 97.5 | 97.4 | 97.5 | 97.7 | 97.6 | 97.6 | 97.6 | 97.8 | 97.5 | 97.5 | 97.5 |
| 96.2 | 95.8 | 96 | 95.5 | 96.1 | 96.1 | 95.9 | 96.3 | 95.6 | 95.9 | 96.1 |
| 97.8 | 97.7 | 97.8 | 97.9 | 97.9 | 97.9 | 97.8 | 98 | 97.8 | 97.8 | 97.7 |
| 87.8 | 85.3 | 88.2 | 91.3 | 89 | 90.5 | 88.8 | 90.6 | 87.4 | 87.6 | 87.1 |
| 97.8 | 97.7 | 97.8 | 97.9 | 97.9 | 97.9 | 97.8 | 98 | 97.8 | 97.8 | 97.7 |
| 22.5 | 12.8 | 13.5 | 33.3 | 64.6 | 18.9 | 10.2 | 13 | 14.7 | 15.8 | 26.1 |
| 4.6 | 6.6 | 5.2 | 4 | 4.7 | 3.5 | 4.4 | 3.9 | 4.2 | 4.4 | 4.2 |
| 5.2 | 2.9 | 7.4 | 9.2 | 3.1 | 9.8 | 3.3 | 5.8 | 4.4 | 3.3 | 3.8 |
| 10.2 | 9.7 | 11.4 | 12.6 | 8.9 | 12.4 | 7.9 | 10.7 | 9.4 | 9.5 | 9.1 |
| 65.7 | 73.6 | 55.8 | 50.2 | 74.9 | 47.6 | 75.6 | 60.4 | 69 | 73.7 | 71.9 |
| 61.9 | 67.8 | 51.8 | 47.4 | 71.1 | 45.3 | 71.9 | 57.5 | 65.4 | 69.8 | 68.1 |
| 98.3 | 98.1 | 98.1 | 97.9 | 98.3 | 97.9 | 98.2 | 98.4 | 98.3 | 98.2 | 98.2 |
| 72036482 | 17538332 | 49231019 | 1.7E+08 | 80424447 | 1.64E+08 | 235931131 | 37988832 | 59391343 | 58972209 | 1.3E+08 |
| 3230 | 2366 | 2776 | 4926 | 6159 | 3154 | 1007 | 1838 | 2296 | 2314 | 3630 |
| 32317 | 12500 | 21111 | 103491 | 253704 | 37710 | 4443 | 11931 | 15914 | 16055 | 41424 |
| 2229 | 1403 | 2332 | 1639 | 317 | 4360 | 5310 | 3184 | 3732 | 3673 | 3148 |

| D3.5_Dox_C1 | D3.5_Dox_C2 | D4_Dox_C1 | D4_Dox_C2 | D4.5_Dox_C1 | D4.5_Dox_C2 | D5_Dox_C1 | D5_Dox_C2 | D5.5_Dox_C1 | D5.5_Dox_C2 | D6_Dox_C1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 6138 | 3562 | 11428 | 16183 | 20437 | 20725 | 20293 | 28005 | 16917 | 12974 | 19034 |
| 15929 | 14788 | 16574 | 17265 | 17466 | 17681 | 17882 | 17837 | 17425 | 16996 | 18190 |
| 96.3 | 96.6 | 97 | 97.6 | 95.9 | 96.2 | 96.3 | 94.9 | 96 | 96 | 95.1 |
| 97.6 | 97.6 | 97.6 | 97.6 | 97.6 | 97.3 | 97.3 | 97 | 97.3 | 97.3 | 97.8 |
| 95.7 | 96.3 | 96.1 | 95.9 | 96 | 95.3 | 95.7 | 95.2 | 95.3 | 95.9 | 96 |
| 89.3 | 89.7 | 89.6 | 89.7 | 87.9 | 83.1 | 82.9 | 84.1 | 83.4 | 84.3 | 92 |
| 97.9 | 97.9 | 97.9 | 97.9 | 97.8 | 97.6 | 97.6 | 97.5 | 97.6 | 97.6 | 98 |
| 15.3 | 12.1 | 22.5 | 28.9 | 38.2 | 31.5 | 34.4 | 42.1 | 37.5 | 27.4 | 56.6 |
| 4.6 | 4.6 | 4.5 | 4.5 | 4.7 | 5.5 | 5.5 | 5.1 | 5.4 | 5 | 3.7 |
| 3.3 | 3 | 3 | 2.6 | 2.8 | 2.7 | 2.7 | 2.5 | 2.8 | 2.7 | 3.1 |
| 9 | 9 | 9 | 8.4 | 8.9 | 7.6 | 7.3 | 7.4 | 8 | 7.5 | 10 |
| 74.5 | 76.3 | 76.1 | 77.8 | 75.8 | 74.6 | 74.7 | 76.5 | 74.5 | 75.8 | 75.5 |
| 70.7 | 72.4 | 72.3 | 74 | 71.8 | 69.6 | 69.7 | 71.7 | 69.7 | 71.4 | 73 |
| 98.3 | 98.3 | 98.4 | 98.1 | 98.3 | 98.3 | 98.4 | 98.3 | 98.4 | 98.4 | 98.4 |
| 55079302 | 21741409 | 94013331 | 1.69E+08 | 1.88E+08 | 1.78E+08 | 2.01E+08 | 2.5E+08 | 1.48E+08 | 92501384 | 4.19E+08 |
| 1782 | 1284 | 2532 | 3078 | 3490 | 3460 | 3308 | 3986 | 3032 | 2586 | 3223 |
| 11906 | 6320 | 23014 | 34713 | 52881 | 49701 | 49996 | 77855 | 44353 | 28798 | 75461 |
| 4626 | 3440 | 4085 | 4877 | 3551 | 3576 | 4018 | 3209 | 3338 | 3212 | 5554 |

| D6_Dox_C2 | D6.5_Dox_C1 | D6.5_Dox_C2 | D7_Dox_C1 | D7_Dox_C2 | D7.5_Dox_C1 | D7.5_Dox_C2 | D8_Dox_C1 | D8_Dox_C2 | D8.25_2i_C1 | D8.25_2i_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 39404 | 32776 | 25293 | 27686 | 25478 | 19859 | 11274 | 6435 | 4995 | 6758 | 5702 |
| 18938 | 16277 | 17548 | 18209 | 18024 | 17416 | 16519 | 15616 | 15285 | 15657 | 15714 |
| 95.6 | 96.7 | 96.2 | 94.8 | 95.5 | 94.3 | 92.7 | 90.6 | 91.7 | 93.1 | 92.6 |
| 97.9 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.6 | 97.6 | 97.6 | 97.6 |
| 96.4 | 96.4 | 96 | 96.2 | 96 | 96 | 95.7 | 95.8 | 96.1 | 96 | 96 |
| 93.2 | 92.6 | 92.1 | 92.1 | 92.2 | 92 | 92.3 | 90.9 | 90.4 | 90.3 | 90.3 |
| 98.1 | 98 | 98 | 98 | 98 | 98 | 98 | 97.9 | 97.9 | 97.9 | 97.9 |
| 85.2 | 81.8 | 54.1 | 65.5 | 47.9 | 51.1 | 26.3 | 23.2 | 20.7 | 21.2 | 19.1 |
| 4 | 4.5 | 3.9 | 4.1 | 4 | 3.9 | 3.8 | 3.9 | 3.9 | 4.4 | 4.5 |
| 3.5 | 2.8 | 2.5 | 3.2 | 3 | 3.1 | 3.7 | 5.7 | 4.2 | 3.3 | 4.1 |
| 9.7 | 11.6 | 9.1 | 11.2 | 10.7 | 11.1 | 10 | 10.4 | 9.1 | 9.2 | 9.3 |
| 73.7 | 73.3 | 77.1 | 73.1 | 73.9 | 73.7 | 72.3 | 64.3 | 71.4 | 75.2 | 71.4 |
| 71.2 | 70.2 | 74.4 | 70.2 | 71.1 | 70.9 | 69.8 | 61.3 | 68.2 | 71.5 | 67.8 |
| 98.5 | 98.4 | 98.4 | 98.3 | 98.3 | 98.4 | 98.4 | 98.2 | 98.4 | 98.3 | 98.3 |
| 1.35E+09 | 1.55E+08 | 2.21E+08 | 4.31E+08 | 2.72E+08 | 1.78E+08 | 65541812 | 33456383 | 24003361 | 28066499 | 27516277 |
| 4897 | 4717 | 4114 | 4327 | 4154 | 3667 | 2494 | 1644 | 1374 | 1692 | 1587 |
| 471033 | 290563 | 85899 | 137190 | 80817 | 68735 | 26535 | 17805 | 11221 | 15122 | 12979 |
| 2868 | 535 | 2576 | 3138 | 3369 | 2591 | 2470 | 1879 | 2139 | 1856 | 2120 |

| D8.25_serum_C1 | D8.25_serum_C2 | D8.5_2i_C1 | D8.5_2i_C2 | D8.5_serum_C1 | D8.5_serum_C2 | D8.75_2i_C1 | D8.75_2i_C2 | D8.75_serum_C1 | D8.75_serum_C2 | D9_2i_C1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 7892 | 6359 | 19378 | 14092 | 14336 | 12381 | 4785 | 5962 | 5629 | 10133 | 15871 |
| 15808 | 15972 | 16274 | 16219 | 16335 | 16274 | 15033 | 15231 | 15445 | 16266 | 16091 |
| 90.7 | 88.9 | 92.6 | 92.8 | 90.7 | 91.9 | 92.2 | 92.2 | 89.6 | 87.1 | 85.3 |
| 97.7 | 97.7 | 97.8 | 97.8 | 97.8 | 97.8 | 97.7 | 97.7 | 97.6 | 97.7 | 97.8 |
| 95.6 | 96.1 | 96.3 | 96.2 | 95.7 | 96 | 96.1 | 95.7 | 95.8 | 96.1 | 96.4 |
| 91.4 | 90.7 | 93.5 | 93.5 | 93.4 | 93.6 | 91.7 | 90.5 | 90.1 | 90.6 | 93.7 |
| 97.9 | 97.9 | 98 | 98 | 98 | 98 | 98 | 97.9 | 97.9 | 97.9 | 98 |
| 25.9 | 25.2 | 50.1 | 36.2 | 39.6 | 35.8 | 17.6 | 19.1 | 18.8 | 26.3 | 52.1 |
| 3.8 | 3.6 | 3.8 | 3.9 | 4 | 3.9 | 3.7 | 3.9 | 3.9 | 3.7 | 3.9 |
| 5.4 | 4.5 | 2.4 | 2.4 | 2.3 | 2.4 | 3.2 | 2.9 | 2.9 | 3.1 | 2.3 |
| 10.7 | 8.9 | 7.2 | 7 | 7.8 | 7.7 | 9 | 8.8 | 9.4 | 9.5 | 7.3 |
| 65 | 70.7 | 79.6 | 79.8 | 78.9 | 78.6 | 75.3 | 76.8 | 76 | 75 | 79.5 |
| 62.2 | 67.9 | 76.5 | 76.6 | 75.6 | 75.6 | 72.5 | 73.5 | 72.7 | 71.9 | 76.4 |
| 98.2 | 98.4 | 98.2 | 98 | 98 | 97.9 | 98.4 | 98.4 | 98.3 | 98.3 | 97.8 |
| 34670761 | 38854100 | 71646422 | 57753221 | 66514572 | 60937426 | 17654865 | 20225030 | 20630020 | 40237550 | 64328422 |
| 1901 | 1601 | 3119 | 2534 | 2653 | 2451 | 1333 | 1552 | 1529 | 2275 | 2817 |
| 22382 | 16332 | 60410 | 35193 | 40214 | 31754 | 9830 | 12257 | 12766 | 26367 | 59016 |
| 1549 | 2379 | 1186 | 1641 | 1654 | 1919 | 1796 | 1650 | 1616 | 1526 | 1090 |

| D9_2i_C2 | D9_serum_C1 | D9_serum_C2 | D9.5_2i_C1 | D9.5_2i_C2 | D9.5_serum_C1 | D9.5_serum_C2 | D10_2i_C1 | D10_2i_C2 | D10_serum_C1 | D10_serum_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 13794 | 6160 | 8071 | 9665 | 13737 | 8356 | 8383 | 5660 | 9422 | 7906 | 3321 |
| 15694 | 15502 | 15526 | 15662 | 15572 | 15936 | 15754 | 15323 | 15798 | 16178 | 14888 |
| 94.5 | 95 | 95.2 | 90.5 | 89.9 | 87.2 | 86.6 | 91.3 | 92.5 | 83.5 | 85.1 |
| 97.8 | 97.8 | 97.9 | 97.6 | 97.7 | 97.7 | 97.6 | 97.8 | 97.8 | 97.8 | 97.8 |
| 96.2 | 96.2 | 96 | 95.9 | 96.3 | 96.1 | 96.2 | 95.9 | 95.9 | 96 | 95.6 |
| 93.7 | 93.5 | 93.6 | 90.4 | 90.7 | 90.8 | 90.3 | 92.5 | 92.3 | 92.2 | 91.9 |
| 98 | 98 | 98 | 97.9 | 97.9 | 97.9 | 97.9 | 98 | 98 | 98 | 98 |
| 42.9 | 52.1 | 64.2 | 40.4 | 49.8 | 39.1 | 41.1 | 24.7 | 33.7 | 31.1 | 15.8 |
| 3.6 | 3 | 3.1 | 3.3 | 3.5 | 3.5 | 3.2 | 3.5 | 3.5 | 3.6 | 3.4 |
| 2.2 | 1.8 | 2 | 3.3 | 4 | 3.9 | 2.9 | 5.9 | 5 | 4.1 | 3.3 |
| 7 | 4.4 | 5.2 | 9.7 | 9.6 | 10.9 | 8.7 | 12 | 11.8 | 12.7 | 11.9 |
| 80.3 | 85.3 | 83.8 | 75.9 | 72.9 | 71.4 | 78.1 | 63.8 | 67.1 | 69.3 | 73.6 |
| 77.5 | 83.2 | 81.7 | 73.1 | 70 | 68.6 | 75.3 | 61.3 | 64.7 | 66.7 | 71.1 |
| 98.1 | 98.5 | 98.5 | 98.3 | 98.2 | 98.2 | 98.3 | 98.1 | 98.2 | 98.1 | 98.3 |
| 34630027 | 33750278 | 40057020 | 29703571 | 31593148 | 31931324 | 29811637 | 17333643 | 27704152 | 33583765 | 895917 |
| 2753 | 1977 | 2317 | 2185 | 2732 | 2056 | 1892 | 1645 | 2358 | 2068 | 1210 |
| 36684 | 18322 | 32382 | 29973 | 52831 | 27622 | 26127 | 16523 | 30277 | 26013 | 7939 |
| 944 | 1842 | 1237 | 991 | 598 | 1156 | 1141 | 1049 | 915 | 1291 | 1128 |

| 11465 | 9225 | 8158 | 6896 | 8173 | 8421 | 4054 | 4176 | 11511 | 14816 | 15611 |
|---|---|---|---|---|---|---|---|---|---|---|
| 16115 | 15697 | 15951 | 15650 | 15758 | 15560 | 15335 | 15379 | 16398 | 16538 | 17172 |
| 92.4 | 91.8 | 72.5 | 78.8 | 79.2 | 89.8 | 86 | 80.8 | 88.4 | 90.7 | 85.8 |
| 97.7 | 97.7 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.7 | 97.8 | 97.7 | 97.8 |
| 95.5 | 95.7 | 96 | 96.1 | 96.2 | 95.7 | 96.1 | 95.7 | 95.5 | 96.3 | 96.2 |
| 91.8 | 91.9 | 91.7 | 92.2 | 92 | 92.6 | 91.5 | 90.4 | 92 | 91.9 | 91.6 |
| 98 | 98 | 98 | 98 | 98 | 98 | 98 | 97.9 | 98 | 97.9 | 98 |
| 30.1 | 25.8 | 29 | 30.8 | 29.4 | 27.2 | 19.4 | 25.6 | 40.9 | 49 | 60.1 |
| 3.7 | 3.7 | 3.8 | 3.5 | 3.7 | 3.8 | 3.5 | 3.6 | 3.7 | 3.6 | 3.5 |
| 3.6 | 3.6 | 2.8 | 2.7 | 2.4 | 3 | 2.3 | 2 | 2.7 | 2.7 | 2.5 |
| 13 | 11.9 | 11.8 | 11 | 9.2 | 10.5 | 10.7 | 8.5 | 10.7 | 10.5 | 11.6 |
| 71.4 | 71.5 | 74.7 | 76 | 78.3 | 74.4 | 78.3 | 81.5 | 76.6 | 76.7 | 76.7 |
| 68.5 | 68.8 | 72 | 73.6 | 75.6 | 71.9 | 75.7 | 78.8 | 73.9 | 74.1 | 74.1 |
| 98.1 | 98.1 | 98.2 | 98.2 | 98.2 | 98.2 | 98.3 | 98.4 | 98.3 | 98.3 | 98.2 |
| 24523951 | 17574924 | 26189701 | 22243909 | 18033999 | 13379426 | 12888357 | 12788655 | 27834347 | 35823619 | 90774725 |
| 2717 | 2369 | 2313 | 2171 | 2171 | 2308 | 1585 | 1692 | 2783 | 3298 | 3586 |
| 31973 | 25324 | 27167 | 21765 | 23981 | 22188 | 9160 | 10612 | 38658 | 54360 | 77058 |
| 767 | 694 | 964 | 1022 | 752 | 603 | 1407 | 1205 | 720 | 659 | 1178 |
| D10.5_2i_C1 | D10.5_2i_C2 | D10.5_serum_C1 | D10.5_serum_C2 | D11_2i_C1 | D11_2i_C2 | D11_serum_C1 | D11_serum_C2 | D11.5_2i_C1 | D11.5_2i_C2 | D11.5_serum_C1 |

| D11.5_serum_C2 | D12_2i_C1 | D12_2i_C2 | D12_serum_C1 | D12_serum_C2 | D12.5_2i_C1 | D12.5_2i_C2 | D12.5_serum_C1 | D12.5_serum_C2 | D13_2i_C1 | D13_2i_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 5562 | 10044 | 12519 | 8119 | 7210 | 10070 | 15004 | 10108 | 21756 | 12776 | 11522 |
| 15665 | 16604 | 16529 | 16471 | 16513 | 16343 | 16879 | 16850 | 18479 | 16853 | 16820 |
| 86.2 | 86.2 | 85 | 84.8 | 85.4 | 84.3 | 86 | 84.7 | 81.5 | 66.3 | 49.1 |
| 97.8 | 97.8 | 97.8 | 97.7 | 97.7 | 97.7 | 97.7 | 97.7 | 97.7 | 97.7 | 97.7 |
| 95.6 | 96.2 | 96 | 96 | 96.2 | 96.1 | 96.1 | 96 | 96.1 | 96.1 | 95.8 |
| 91.9 | 92 | 91.4 | 91 | 90.6 | 91 | 91.2 | 90.8 | 90.8 | 90.8 | 90.8 |
| 98 | 98 | 98 | 98 | 97.9 | 97.9 | 97.9 | 97.9 | 97.9 | 98 | 98 |
| 23.6 | 51.4 | 55.3 | 35.4 | 29.9 | 37.9 | 47.7 | 35 | 67.1 | 56.4 | 72.9 |
| 3.5 | 4.1 | 3.8 | 3.6 | 3.6 | 4.1 | 4 | 3.7 | 3.8 | 4.3 | 4.3 |
| 2.4 | 2.8 | 2.7 | 2.4 | 2.3 | 2.9 | 2.9 | 2.3 | 2.4 | 3.1 | 2.8 |
| 10.9 | 8.6 | 7.8 | 9.4 | 9.3 | 8.5 | 8.4 | 8.5 | 8.8 | 8.8 | 8.3 |
| 77.4 | 77.1 | 78.7 | 78.7 | 79.2 | 76.8 | 76.8 | 79.7 | 79.2 | 75.5 | 76.8 |
| 74.9 | 74.3 | 76 | 76.1 | 76.4 | 73.7 | 73.8 | 76.8 | 76.3 | 72.1 | 73.4 |
| 98.2 | 98.5 | 98.5 | 98.4 | 98.4 | 98.4 | 98.4 | 98.3 | 98.3 | 98.3 | 98.3 |
| 15149367 | 34932625 | 36075300 | 27804384 | 27170840 | 22372820 | 36585438 | 30987716 | 1.66E+08 | 49269432 | 1.01E+08 |
| 1903 | 2523 | 2880 | 2468 | 2358 | 2560 | 3214 | 2816 | 4369 | 2938 | 2866 |
| 14238 | 42704 | 58092 | 25116 | 20552 | 32471 | 54768 | 29456 | 138451 | 75220 | 156892 |
| 1064 | 818 | 621 | 1107 | 1322 | 689 | 668 | 1052 | 1201 | 655 | 643 |

| D13_serum_C1 | D13_serum_C2 | D13.5_2i_C1 | D13.5_2i_C2 | D13.5_serum_C1 | D13.5_serum_C2 | D14_2i_C1 | D14_2i_C2 | D14_serum_C1 | D14_serum_C2 | D14.5_2i_C1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 12190 | 15494 | 5599 | 5146 | 5287 | 5360 | 15207 | 20543 | 10816 | 14705 | 12798 |
| 17377 | 18070 | 16769 | 15987 | 16853 | 16725 | 18525 | 18764 | 18461 | 18884 | 18532 |
| 77.6 | 85.4 | 74.6 | 75.3 | 77.2 | 71.1 | 91.6 | 93.6 | 87.9 | 87.1 | 81.5 |
| 97.8 | 97.8 | 97.8 | 97.7 | 97.8 | 97.8 | 97.8 | 97.7 | 97.7 | 97.7 | 97.8 |
| 96.3 | 96.1 | 96.3 | 95.7 | 95.9 | 95.9 | 96.3 | 96.2 | 96 | 96.1 | 96 |
| 92.1 | 92.2 | 92.5 | 90.8 | 90.9 | 91 | 91.8 | 91.7 | 91.4 | 91.5 | 92 |
| 98 | 98 | 98 | 97.9 | 98 | 97.9 | 98 | 97.9 | 97.9 | 98 | 98 |
| 73.7 | 67.1 | 69.4 | 52.4 | 70.2 | 68.1 | 37 | 42.1 | 39.5 | 50.7 | 36.7 |
| 4 | 4 | 5.7 | 5.3 | 5.6 | 5.5 | 4.9 | 4.8 | 4.1 | 3.9 | 5.6 |
| 2.4 | 2.5 | 4.9 | 5.4 | 3.4 | 3.5 | 3.1 | 3.1 | 2.5 | 2.6 | 3.4 |
| 9.6 | 10.3 | 16 | 15.7 | 14.9 | 14.6 | 7.6 | 7.4 | 11.2 | 10.4 | 9.3 |
| 78.1 | 77 | 65.4 | 63.3 | 70.3 | 70.6 | 77.5 | 77.6 | 77.2 | 77.6 | 74.4 |
| 75 | 73.8 | 60.7 | 59 | 65.9 | 66.3 | 73.3 | 73.5 | 73.7 | 74.3 | 69.7 |
| 98.3 | 98.3 | 97.5 | 97.5 | 97.3 | 97.3 | 98.3 | 98.4 | 98.3 | 98.3 | 98.3 |
| 97956936 | 1.09E+08 | 49186630 | 22110011 | 54618691 | 46544722 | 74206890 | 1.05E+08 | 70077873 | 97873582 | 79779089 |
| 3179 | 3646 | 1996 | 1853 | 2056 | 2126 | 3022 | 3577 | 2897 | 3539 | 2744 |
| 99956 | 93789 | 46666 | 26735 | 43074 | 42121 | 39097 | 54136 | 34487 | 56705 | 39164 |
| 980 | 1166 | 1054 | 827 | 1268 | 1105 | 1898 | 1938 | 2032 | 1726 | 2037 |

| D14.5_2i_C2 | D14.5_serum_C1 | D14.5_serum_C2 | D15_2i_C1 | D15_2i_C2 | D15_serum_C1 | D15_serum_C2 | D15.5_2i_C1 | D15.5_2i_C2 | D15.5_serum_C1 | D15.5_serum_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 15068 | 8409 | 14650 | 5664 | 7023 | 11915 | 5252 | 8467 | 9841 | 15905 | 13986 |
| 18770 | 18018 | 18580 | 18159 | 17960 | 18739 | 18103 | 18490 | 18358 | 19807 | 19970 |
| 89.7 | 78.9 | 79.2 | 85.3 | 92.1 | 66.9 | 63.9 | 94.4 | 94.3 | 76.9 | 82.2 |
| 97.8 | 97.8 | 97.8 | 97.8 | 97.7 | 97.8 | 97.8 | 97.7 | 97.7 | 97.8 | 97.8 |
| 95.6 | 96.1 | 96.4 | 96.2 | 95.7 | 95.7 | 96 | 96 | 96.3 | 95.9 | 96 |
| 92 | 91.6 | 91.9 | 91.6 | 91.5 | 91.5 | 91.6 | 91.6 | 92.1 | 92 | 91.9 |
| 98 | 98 | 98 | 98 | 97.9 | 98 | 98 | 97.9 | 97.9 | 98 | 98 |
| 33.7 | 42 | 59.7 | 61.6 | 38.4 | 39.9 | 46 | 21.3 | 23 | 66.5 | 54.1 |
| 5.3 | 4.9 | 4.1 | 7.9 | 6 | 3.9 | 5.1 | 4.5 | 4.3 | 4.3 | 4.4 |
| 3.3 | 2.7 | 2.4 | 5 | 4.8 | 2.3 | 2.9 | 3.4 | 3.4 | 2.9 | 3 |
| 8.7 | 12 | 10 | 18 | 14.1 | 10 | 13.5 | 7.7 | 7.6 | 12.6 | 12.5 |
| 75.4 | 75.8 | 78.9 | 63.1 | 67.5 | 79 | 74 | 76.8 | 76.8 | 74.5 | 74.2 |
| 71 | 71.6 | 75.6 | 56.2 | 62.2 | 75.7 | 69.5 | 72.7 | 73 | 70.8 | 70.4 |
| 98.3 | 98.2 | 98.4 | 97.4 | 97.9 | 98.3 | 97.8 | 98.2 | 98.3 | 98.2 | 98.2 |
| 78954514 | 45618882 | 1.05E+08 | 82113379 | 46137688 | 86460491 | 55835189 | 70722479 | 66435427 | 2.48E+08 | 1.65E+08 |
| 3074 | 2505 | 3705 | 1935 | 2111 | 3162 | 2007 | 1964 | 2143 | 3685 | 3347 |
| 37795 | 33892 | 76526 | 32100 | 20244 | 48958 | 25885 | 16535 | 19528 | 107956 | 64367 |
| 2089 | 1346 | 1377 | 2558 | 2279 | 1766 | 2157 | 4277 | 3402 | 2295 | 2556 |

| D16_2i_C1 | D16_2i_C2 | D16_serum_C1 | D16_serum_C2 | D16.5_2i_C1 | D16.5_2i_C2 | D16.5_serum_C1 | D16.5_serum_C2 | D17_2i_C1 | D17_2i_C2 | D17_serum_C1 |
|---|---|---|---|---|---|---|---|---|---|---|
| 5076 | 9135 | 6791 | 8342 | 8471 | 5373 | 13361 | 6278 | 12668 | 10936 | 15523 |
| 17665 | 17761 | 18278 | 18336 | 18679 | 18374 | 19896 | 18796 | 18877 | 18501 | 19538 |
| 92.2 | 94.5 | 57 | 78.1 | 89.2 | 88.7 | 76.4 | 65.7 | 89.9 | 90.5 | 88.1 |
| 97.8 | 97.7 | 97.8 | 97.7 | 97.8 | 97.7 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 |
| 96.3 | 96.2 | 95.6 | 96.1 | 96.3 | 96.2 | 96.1 | 96.3 | 96.2 | 96.2 | 96.1 |
| 91.9 | 91.8 | 91.5 | 91.3 | 91.8 | 91.6 | 91.5 | 91.7 | 91.9 | 91.7 | 91.8 |
| 98 | 97.9 | 97.9 | 97.9 | 97.9 | 97.9 | 98 | 98 | 98 | 98 | 98 |
| 38.5 | 25.7 | 30.4 | 36.6 | 22.6 | 15.9 | 47.3 | 28.2 | 29.8 | 23.6 | 49.4 |
| 3.7 | 3.7 | 4 | 4.1 | 4.2 | 4.2 | 3.9 | 3.9 | 3.9 | 3.8 | 3.9 |
| 3.6 | 3.4 | 2.5 | 2.6 | 3.3 | 3.3 | 2.8 | 2.7 | 3.4 | 3.3 | 3 |
| 8.1 | 7.8 | 8.7 | 10.4 | 7.5 | 7.6 | 11.7 | 10.1 | 7.8 | 7.5 | 11.5 |
| 76.2 | 76.8 | 78.3 | 77.3 | 77.6 | 77.4 | 75.9 | 77.1 | 77.2 | 77.9 | 75.1 |
| 72.9 | 73.4 | 75 | 73.7 | 74 | 73.9 | 72.6 | 73.9 | 73.9 | 74.7 | 71.9 |
| 98.4 | 98.4 | 98.1 | 98.2 | 98.3 | 98.3 | 98.2 | 98.2 | 98.3 | 98.4 | 98.3 |
| 52290532 | 53190608 | 48858555 | 48904299 | 55829324 | 41911584 | 1.35E+08 | 52474229 | 67119554 | 46535861 | 96863752 |
| 1343 | 1921 | 2182 | 2467 | 2124 | 1618 | 3393 | 2119 | 2807 | 2539 | 3583 |
| 13315 | 18996 | 27763 | 28886 | 17424 | 10237 | 57651 | 22716 | 28918 | 22044 | 62052 |
| 3927 | 2800 | 1749 | 1693 | 3204 | 4094 | 2350 | 2310 | 2321 | 2111 | 1561 |

| D17_serum_C2 | D17.5_2i_C1 | D17.5_2i_C2 | D17.5_serum_C1 | D17.5_serum_C2 | D18_2i_C1 | D18_2i_C2 | D18_serum_C1 | D18_serum_C2 | DiPSC_2i_C1 | DiPSC_2i_C2 |
|---|---|---|---|---|---|---|---|---|---|---|
| 12979 | 14477 | 10753 | 12806 | 9998 | 18060 | 17916 | 9840 | 9029 | 10626 | 20527 |
| 19729 | 18309 | 18452 | 19556 | 19155 | 18821 | 18566 | 19294 | 19023 | 17918 | 18049 |
| 86.3 | 92.1 | 92.2 | 85.1 | 87.9 | 90.9 | 90.6 | 80 | 77.3 | 96.4 | 96.2 |
| 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.8 | 97.7 | 97.7 |
| 96.2 | 96.3 | 95.9 | 96.3 | 96 | 96.2 | 96.3 | 96 | 96.4 | 96.1 | 96.1 |
| 91.5 | 92.1 | 91.8 | 91.4 | 91.8 | 92.6 | 92.5 | 92.3 | 92 | 91.3 | 90.9 |
| 98 | 98 | 98 | 97.9 | 98 | 98 | 98 | 98 | 98 | 98 | 97.9 |
| 42 | 40.2 | 28.2 | 44.1 | 36.5 | 58.2 | 54.8 | 62.7 | 48.1 | 20.2 | 28.8 |
| 4 | 3.8 | 4 | 4 | 3.8 | 3.9 | 3.7 | 4.1 | 3.9 | 5.1 | 5.3 |
| 3 | 3.2 | 3.1 | 2.9 | 2.7 | 3.5 | 3.4 | 3 | 2.8 | 3.8 | 3.8 |
| 11.5 | 6.9 | 6.9 | 10.3 | 10.1 | 6.3 | 6 | 10.4 | 9.3 | 9.7 | 9.5 |
| 75 | 78.6 | 78.7 | 76.4 | 77.7 | 77.5 | 78.2 | 75.5 | 76.8 | 71.6 | 71.7 |
| 71.6 | 75.4 | 75.4 | 73.1 | 74.6 | 74.3 | 75 | 72.1 | 73.6 | 67.7 | 67.6 |
| 98.3 | 98.5 | 98.4 | 98.4 | 98.4 | 98.4 | 98.4 | 98.3 | 98.3 | 98.2 | 98.3 |
| 96965300 | 59918421 | 54120470 | 86540688 | 62361742 | 1.39E+08 | 1.04E+08 | 1.18E+08 | 77222647 | 74406713 | 87759016 |
| 3300 | 2900 | 2474 | 3292 | 2849 | 2774 | 2761 | 2472 | 2322 | 2524 | 3649 |
| 45803 | 36580 | 22428 | 44221 | 29527 | 69937 | 63038 | 62257 | 40600 | 21467 | 46879 |
| 2117 | 1638 | 2413 | 1957 | 2112 | 1989 | 1648 | 1898 | 1902 | 3466 | 1872 |

| | DiPSC_serum_C1 | DiPSC_serum_C2 |
|---|---|---|
| | 5247 | 4340 |
| | 18112 | 21502 |
| | 2241 | 2535 |
| | 95034273 | 93322919 |
| | 98.2 | 98.2 |
| | 65.9 | 67.5 |
| | 70.1 | 71.4 |
| | 10.3 | 9.8 |
| | 4.4 | 4 |
| | 5.1 | 4.9 |
| | 23.2 | 23.3 |
| | 97.9 | 97.9 |
| | 90.1 | 90.9 |
| | 95.9 | 96.1 |
| | 97.7 | 97.7 |
| | 93.2 | 90.8 |
| | 19202 | 19098 |
| | 7777 | 9449 |

[0249]     A model of the developmental landscape

[0250]     We visualized the developmental landscape of the 251,203 cells in a two-dimensional FLE (FIG. 24B) and annotated it according to sampling time (FIG. 24C), expression scores of gene signatures, and expression of individual genes (FIG. 24D, Table 15).

**Table 15 -** List of genes comprising gene signatures.

| MEF identity | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gm5571 | Il17rd | Gjd4 | Prss23 | Atp10a | Eif4g2 | Gulp1 | Sema3a |
| Rbfox2 | Ptk2 | Ccng1 | 9430030n17rik | Loxl1 | Vcl | Shank1 | Itgb1 |
| Btbd19 | Ehd2 | Gpr124 | Arntl2 | Loxl2 | Bcl2l2 | Bmp1 | Nxn |
| Actn1 | Lats2 | Fibin | Sh3rf1 | Fbln5 | Cd276 | Akt1s1 | Tmem41b |
| Gatad2a | Hspg2 | 8030476119rik | Mrc2 | Ctgf | Lrrc58 | Itga9 | Sec23a |
| Med6 | 4930456g14rik | Ddr2 | Mdh1 | Efnb2 | Wwc2 | Abcc1 | Gm22 |
| Mex3a | 4930429b21rik | Arf4 | Rictor | Rxra | Lpp | Eda | Itgb5 |
| Ccdc80 | | Ptprs | Map4k5 | Ccnd2 | Arl1 | B4galt2 | Dysf |
| Mex3c | Rps20 | Sprr2k | Plcl1 | Gpc2 | Ltbp1 | Nid1 | Thbs1 |
| Sdpr | Vgll3 | Adm | 11-Sep | Ntf3 | Ltbp2 | Ncam1 | Bc022687 |
| Pcdhb2 | Prr15 | A830029e22rik | Ryk | Kif5b | Wisp1 | Shc2 | Dnm3os |
| Trim16 | Fbxl7 | 9230114k14rik | Tgfb3 | Slit2 | Igf1r | Uba6 | Rnd3 |
| Obsl1 | Maged2 | Extl3 | Ube2i | Tpm1 | Rhobtb3 | Tradd | Pik3c2a |
| Epha1 | Galnt4 | Mecom | Tgfb2 | Gpc4 | Fam198b | Rtel1 | 2810008m24rik |
| Stx1b | Pdgfc | Qsox1 | Zfp319 | Flnb | Cnn2 | Bicd2 | Spred3 |
| Stau1 | Tmtc4 | | Gm10399 | 4930555b11rik | Glipr2 | Adamts12 | Senp5 |

Thbsl, Bc022687, Nidi, Ncaml, Ltbpl, Ltb 2, Gpc2, Ntf3, Plcll, Sprr2k, Rps20, Vgll3, Mex3c, Sd r

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Serpinel | Tmtc3 | Teadl | Fbxol7 | Fine | Sydel | Hs2stl | Aril 3b |
| Aa881470 | Lpar4 | Snx7 | Wnt5a | C76332 | Hhat | D10ertd610e | Polr2e |
| Coll2al | Pcdhl9 | Cdkl4 | Criml | Capn2 | Zmat3 | Cyr61 | Itgav |
| 2010300f17rik | Eda2r | Cdkn2a | Midi | Phlda3 | Caldl | Gtf3cl | Igf2bp3 |
| Cede 102a | Pcdhl8 | Cdkn2b | Displ | Map3k7 | Pmepal | Lbh | |
| Nradd | Gprl76 | Ccnyll | Ubox5 | MyhlO | E1301 121 23rik | Krt33b | |
| Pard6g | Loc 100503471 | Tubb2a-ps2 | St71 | D18ertd653e | Bag2 | Gm6607 | |
| Ntn4 | Mical2 | Aen | Col5a2 | Stox2 | Zfp583 | D3wsul67e | |
| 5730471h19rik | Dzipll | Farpl | Axl | Igf2r | Pibfl | Zc3h7b | |
| Sepnl | Hoxc6 | 4930402h24rik | Col5al | D15ertd621e | Pmaipl | 7630403g23rik | |
| Peg 12 | Hoxc5 | Sh3rf3 | Zyx | Arid5b | A130022j15rik | Tnpo2 | |
| Dpysl3 | Mettl4-psl | Adaml9 | Ror2 | TnfrsflOb | Bcl91 | Cepl70 | |
| 1110012d08rik | Sec63 | Ddbl | Wdfy3 | 261001 1e 03rik | Cpa6 | Pdlim5 | |
| Aktl | Ikbip | Cttn | Amotl2 | Ckap4 | D13ertd787e | Pdlim7 | |
| Zfp286 | Tsc22d2 | 92301 12e 08rik | Yapl | Efna2 | Pabpc41 | Cad | |
| Ubap21 | 23 10076g05rik | Dbnl | Phldb2 | Picalm | Zfhx3 | Unc5b | |
| Samd4 | Anxa6 | Fyttdl | 6330562c 20rik | CdhlO | Itga5 | 2410018113rik | |
| Phc2 | Nfatc4 | Lrrcl5 | Ctnndl | Ddahl | Txnrdl | Loc 100216343 | |
| Mcam | Fnl | FkbplO | Rock2 | Uba3 | Htrlb | Glrx3 | |
| Pla2g4c | Wnt9a | Trubl | Maspl | 0610038b21rik | Hmga2 | Kctd5 | |
| Fzd7 | Sorcs2 | Zdhhc20 | Pvtl | Gemin7 | 2-Sep | Loc269472 | |
| Pappa | Tmeffl | Stonl | Tnc | Ubal | Lambl | Myolc | |
| Ptk7 | C79491 | Hoxdl3 | Fbln2 | Fbnl | Zfp5 18b | 4930562c 15rik | |
| Nuakl | Crlfl | Nudt6 | Hdlbp | Lhx9 | Parva | Till | |
| | 2610034e 0lrik | Hoxdl2 | | | | | |

| Plunpotency | | | | | | | |
|---|---|---|---|---|---|---|---|
| Rhox5 | Mt2 | Asns | Taf7 | Folrl | Sox2 | Grhpr | Chmp4c |
| Tdgfl | Ube2a | Aldoa | Nudt4 | Gm7325 | Jam2 | Higdla | Hsf2bp |
| Utfl | Khdc3 | Tdh | Cox5a | Agtrap | Fkbp3 | Rpp25 | Polr2e |
| Mkrnl | Pycard | Gjb3 | Sod2 | Sppl | Cox7b | Rbpms | Blvrb |
| Dppa5a | Hsp90aal | Rbpms2 | S100al3 | Hells | Ash21 | Mmp3 | Ldhb |
| Uppl | Prrcl | Prpsl | Fkbp6 | Dppa4 | Dut | Apobec3 | Apocl |
| ChchdlO | Hatl | Fam25c | Rhox9 | Gabarapl2 | Dtymk | Spc24 | Syngrl |
| Klf2 | Calcoco2 | Eif2s2 | Gdf3 | Rhox6 | Gpx4 | Xlr3a | Bexl |
| Trap1a | Impa2 | Cenpm | 2700094K13Rik | Rhoxl | Eif4ebpl | Reel 14 | Nr2c2ap |
| Mylpf | Saa3 | Nanog | Fmrlnb | Cdc51 | Morel | Mtf2 | |
| 1700013H16Rik | Ooep | Ndufa412 | Hmgn2 | Texl9. 1 | Fabp3 | Snrpn | |
| AA467197 | Bnip3 | Syce2 | Ubald2 | Trim28 | Zfp428 | Gml3580 | |
| Dhxl6 | Mtl | Gml325 1 | Lactb2 | Atp5gl | Aqp3 | Gmnn | |

| Cell cycle | | | | | | | |
|---|---|---|---|---|---|---|---|
| Mcm4 | Lbr | Cdkl | Ndc80 | Cdca2 | Rrm2 | Hjurp | Rpa2 |
| Smc4 | Cenpf | Slbp | Mcm6 | Nasp | Tipin | Tacc3 | Gins2 |
| Gtsel | Birc5 | Aurkb | Rrm1 | Gmnn | Casp8ap2 | Mcm5 | E2f8 |
| Ttk | Dtl | Kif11 | Mlf1ip | Cdc6 | Tubb4b | Anp32e | Cdc25c |
| Rangap1 | Dscc1 | Cks1b | Top2a | Pold3 | Kif23 | Dlgap5 | Nek2 |
| Ccnb2 | Cbx5 | Blm | Hmgb2 | Ckap2l | Exo1 | Ect2 | Cdc20 |
| Cenpa | Usp1 | Msh2 | Ccne2 | Fam64a | Rfc2 | Nuf2 | Rad51ap1 |
| Cenpe | Hmmr | Gas2l3 | G2e3 | Ubr7 | Pola1 | Cdc45 | |
| Cdca8 | Wdr76 | Tyms | Tmpo | Fen1 | Mki67 | Ckap5 | |
| Ckap2 | Ung | Hjurp | Nusap1 | Bub1 | Tpx2 | Ctcf | |
| Rad51 | Hn1 | Hells | Ncapd2 | Brip1 | Aurka | Clspn | |

| Pcna | Cks2 | Priml | Mcm2 | Atad2 | Anln | Cdca7 | |
|------|------|-------|------|-------|------|-------|---|
| Ube2c | Kif20b | Uhrf1 | Kif2c | Psrc1 | Chaf1b | Cdca3 | |

| ER Stress | | | | | | | |
|-----------|---|---|---|---|---|---|---|
| Nck2 | Chac1 | Creb3 | Itpr1 | Os9 | Stt3b | Dnajb9 | Crebrf |
| Ankzf1 | Pdia3 | Sec61b | Edem1 | Ddit3 | Rnf185 | Tmx1 | Bak1 |
| Dnajb2 | Bcl2l11 | Erp44 | Bbc3 | Erlin2 | Xbp1 | Jkamp | Rnf5 |
| Rhbdd1 | Ddrgk1 | AI314180 | Psmc4 | Ppp2cb | Erlec1 | Sel11 | Atf6b |
| Bcl2 | Tmx4 | Jun | Bax | Ubxn8 | Stc2 | Psmc1 | Bag6 |
| Ubxn4 | Trib3 | Casp9 | Ppp1r15a | Casp3 | Trp53 | Atxn3 | Flot1 |
| Yod1 | H13 | Fbxo6 | Vimp | Pik3r2 | Alox15 | Derl1 | Eif2ak2 |
| Ppp1r15b | Edem2 | Fbxo2 | Rnf121 | Amfr | Derl2 | Rnf139 | Pmaip1 |
| Fam129a | Cebpb | Ube4b | Anks4b | Herpud1 | Trim25 | Foxred2 | Tmx3 |
| Edem3 | Ptpn1 | Ube2j2 | Ern2 | Aars | Cdk5rap3 | Pla2g6 | Syvn1 |
| Atf6 | Vapb | Psmc2 | Atp2a1 | Selk | Ccdc47 | Atf4 | Erlin1 |
| Ufc1 | Srpx | Tmub1 | Brsk2 | Ero1l | Psmc5 | Ep300 | |
| Atf3 | Aifm1 | Tmem129 | Ins2 | Psmc6 | Ern1 | Tmbim6 | |
| Man1b1 | Ubqln2 | Wfs1 | Ccnd1 | Trim13 | Nploc4 | Txndc11 | |
| Tor1a | Mbtps2 | Ube2k | Map3k5 | Dnajc3 | P4hb | Sdf2l1 | |
| Hspa5 | Usp13 | Tbl2 | Nrbf2 | Casp4 | Txndc5 | Ufd1l | |
| Dab2ip | Ufm1 | Get4 | Derl3 | Casp12 | Faf2 | Eif2b5 | |
| Nfe2l2 | Serp1 | Bhlha15 | Ube2g2 | Scamp5 | Ubqln1 | Nrros | |
| Dnajc10 | Creb3l4 | Creb3l2 | Tmem259 | Pml | Atg10 | Pdia5 | |
| Psmc3 | Tmem67 | Pdia4 | Creb3l3 | Parp16 | Thbs4 | Gsk3b | |
| Creb3l1 | Ufl1 | Eif2ak3 | Hsp90b1 | Nck1 | Col4a3bp | Park2 | |
| Thbs1 | Ube2j1 | Rnf103 | Apaf1 | Uba5 | Pik3r1 | Stub1 | |
| Eif2ak4 | Vcp | Aup1 | Ifng | Usp19 | Pdia6 | Pdia2 | |

| Epithelial Identity | | | | | | | |
|---|---|---|---|---|---|---|---|
| Cdhl | Cldn3 | Cldn7 | Ocln | Crb3 | Krtl9 | Dsp | Pkpl |
| Tgml | Cldn4 | Cldnl 1 | Epcam | Krt8 | Pkp3 | | |

| ECM Rearrangement | | | | | | | |
|---|---|---|---|---|---|---|---|
| Sulf1 | Creb3l1 | B4galt1 | Mia | Atxn11 | Adamts2 | Tnfrsf11b | Cyp1b1 |
| Col19a1 | Hsd17b12 | Reck | Spint2 | Crispld2 | Wnt3a | Col14a1 | Fshr |
| Col3a1 | Wt1 | Tgfbr1 | Aplp1 | Foxf1 | Mfap4 | Has2 | Mkx |
| Col5a2 | Grem1 | Col27a1 | Hpn | Foxc2 | Serpinf2 | Ptk2 | Lox |
| Fn1 | Spint1 | P3h1 | Klk4 | Agt | Vtn | Scx | Hpse2 |
| Ihh | Cst3 | Hspg2 | Acan | Exoc8 | Nf1 | Fbln1 | Kazald1 |
| Col4a4 | Fkbp1a | Vwa1 | Serpinh1 | Ero1l | Col1a1 | Adamts20 | Nfkb2 |
| Col4a3 | Mmp9 | Dnajb6 | Apbb1 | Lgals3 | Ramp2 | Col2a1 | |
| Serpinb5 | Sulf2 | Emilin1 | Ilk | Ripk3 | Gfap | Myh11 | |
| Fmod | Atp7a | Mpv17 | Ric8 | Loxl2 | Sox9 | Ccdc80 | |
| Elf3 | Nox1 | Apbb2 | Muc5ac | Lcp1 | Ero1lb | Abi3bp | |
| Lamc1 | Col4a6 | Pdgfra | Ctgf | Mmp13 | Nid1 | App | |
| Tnr | Prdx4 | Ambn | Nr2e1 | Mmp20 | Foxf2 | Serac1 | |
| Dpt | Gpm6b | Dmp1 | Nepn | Col5a3 | Foxc1 | Plg | |
| Ddr2 | Egfl6 | Ibsp | P4ha1 | Smarca4 | Ripk1 | Smoc2 | |
| Olfml2b | Postn | Tfip11 | Spock2 | Aplp2 | Tfap2a | Has1 | |
| Tgfb2 | Rxfp1 | Eln | Adamts14 | Mpzl3 | Ecm2 | Noxo1 | |
| Itga8 | Sfrp2 | Plod3 | Mmp11 | Thsd4 | B4galt7 | Col11a2 | |
| Adamtsl2 | Hapln2 | Col1a2 | Col18a1 | Anxa2 | Tgfbi | Tnxb | |
| Col5a1 | Ctss | Ndnf | Myf5 | Myo1e | Pxdn | Tnf | |
| Pomt1 | Adamtsl4 | Vhl | Col4a1 | Nphp3 | Smoc1 | 2300002M23Rik | |
| Eng | St7l | Mfap5 | Csgalnact1 | Dag1 | Ltbp2 | Flot1 | |
| Lmx1b | Col11a1 | Ercc2 | Comp | Lamb2 | Flrt2 | Hsp90ab1 | |

| Gsn | Npnt | Bcl3 | Gfod2 | Kif9 | Fbln5 | Washl | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| 01fml2a | Cyr61 | Tgfbl | Has3 | Sh3pxd2b | Egflam | Vit | |

| Apoptosis | | | | | | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| Ercc5 | Procr | Slc35d1 | Ldhb | Zfp365 | Zbtb16 | Sphk1 | Abcc5 |
| Serpinb5 | Blcap | Plk3 | Lrmp | Prmt2 | Rps27l | Rhbdf2 | Trp63 |
| Inhbb | Ada | Rnf19b | Tm7sf3 | Mknk2 | Mapkapk3 | Baiap2 | Fam162a |
| Steap3 | Fgf13 | Sfn | Tgfb1 | Dram1 | Ip6k2 | Dcxr | App |
| Btg2 | Irak1 | Fuca1 | Sertad3 | Apaf1 | Tcn2 | Hist1h1c | Rab40c |
| Phlda3 | Tspyl2 | Epha2 | Cebpa | Btg1 | Lif | Ninj1 | Bak1 |
| Tnni1 | Sat1 | Wrap73 | Klk8 | Mdm2 | Upp1 | Nol8 | Def6 |
| Rgs16 | Zmat3 | Mxd4 | Bax | Ddit3 | Ccng1 | F2r | Cdkn1a |
| Ier5 | Hspa4l | Rchy1 | Ppp1r15a | Gls2 | Cyfip2 | Ankra2 | Tap1 |
| Slc19a2 | Slc7a11 | Iscu | Rpl18 | Dgka | Gnb2ll | Plk2 | Ier3 |
| Adck3 | Tm4sf1 | Triap1 | Aen | Cdkn2aip | Hint1 | Sdc1 | Polh |
| Ephx1 | Rap2b | Prkab1 | Rrp8 | Hmox1 | Gm2a | Gpx2 | Ccnd3 |
| Ptpn14 | Fbxw7 | Trafd1 | Ccp110 | Rrad | Hist3h2a | Zfp36l1 | Hbegf |
| Atf3 | S100a4 | Pom121 | Nupr1 | Cdh13 | Alox8 | Fos | Hdac3 |
| Notch1 | S100a10 | Pdgfa | Ptpre | Osgin1 | Trp53 | Ccnk | Rad9a |
| Rxra | Txnip | Gadd45a | Hras | Cgrrf1 | Tax1bp3 | Jag2 | Ctsf |
| Ralgds | Nhlh2 | Vamp8 | Eps8l2 | Abhd4 | Traf4 | Ndrg1 | Slc3a2 |
| Ak1 | Dnttip2 | Retsat | Ctsd | Kif13b | Cdk5r1 | Pmm1 | Fas |
| Stom | Clca2 | Tprkb | Cd81 | Rb1 | Ppm1d | Plxnb2 | |
| Ddb2 | Wwp1 | Tgfa | Perp | Nudt15 | Rad51c | Vdr | |
| Cd82 | Klf4 | Mxd1 | Rps12 | Tsc22d1 | Tob1 | Csrnp2 | |
| Il1a | Ikbkap | Sec61a1 | Tpd52l1 | Casp1 | Krt17 | Acvr1b | |
| Pcna | Cdkn2a | Xpc | Sesn1 | St14 | Hexim1 | Sp1 | |
| Bmp2 | Cdkn2b | Ccnd2 | Foxo3 | Ei24 | Fdxr | Abat | |

| Trib3 | Jun | H2afj | Ddit4 | Vwa5a | Itgb4 | Socsl | |
|-------|-----|-------|-------|-------|-------|-------|--|

| SASP | | | | | | | |
|------|------|------|------|-------|-------|------|------|
| Il6 | Cxcl2 | Csf2 | Fgf7 | Igfbp4 | Mmp14 | Icam3 | Egfr |
| Il7 | Cxcl3 | Mif | Vegfa | Igfbp6 | Timp2 | Tnfrsf11b | Fn1 |
| Il1a | Ccl8 | Areg | Ang | Igfbp7 | Serpine1 | Tnfrsf1a | |
| Il1b | Ccl13 | Ereg | Kitl | Mmp1 | Serpinb2 | Tnfrsf1b | |
| Il13 | Ccl3 | Nrg1 | Cxcl12 | Mmp3 | Plat | Tnfrsf10b | |
| Il15 | Ccl20 | Egf | Pigf | Mmp10 | Plau | Fas | |
| Cxcl15 | Ccl16 | Fgf2 | Igfbp2 | Mmp12 | Ctsb | Plaur | |
| Cxcl1 | Ccl26 | Hgf | Igfbp3 | Mmp13 | Icam1 | Il6st | |

| Neural Identity | | | | | | | |
|------|------|--------|------|------|------|------|-------|
| Vtn | Zeb2 | Sox1 | Pax6 | Sox2 | Msx1 | Atoh1 | Tubb3 |
| Ednrb | Hes5 | Neurod1 | Cdh2 | Id2 | Msi1 | Rbfox3 | |
| Sox21 | Fabp7 | Pax3 | Sox9 | Hoxb1 | Msi2 | Map2 | |

| Placental Identity | | | | | | | |
|------|------|------|------|------|------|------|------|
| 4933433p14rik | Dusp9 | Pkp2 | Tnfrsf23 | Serpinb9d | Krt18 | 1600014k23rik | Hapln3 |
| Esx1 | H19 | 9630050e16rik | Sos1 | Plekhh1 | Nrn1l | Tbrg1 | Fam176a |
| Afap1 | Tmem37 | Pvrl2 | Dlx3 | 2210011c24rik | Sfi1 | Slit1 | Pdlim1 |
| Zfyve21 | Mmp15 | Zfp568 | Ippk | Cd320 | Tlr5 | A730090h04rik | Ube2q2 |
| Erv3 | Fam101b | Vtcn1 | Htr2b | Ccnjl | Rhou | 4931406p16rik | Au018091 |
| Atg12 | Phf16 | Il6ra | Dusp16 | Entpd2 | Arhgef6 | Opn3 | Bdkrb2 |
| Las1l | 4930422n03rik | Foxo4 | Cdc73 | Il1r2 | Tmem185b | Pdia4 | E130203b14rik |
| Rbp1 | Ada | Hsp90b1 | 1700025g04rik | Sfmbt2 | Tram2 | B930054o08 | S100g |
| Prl2b1 | Mmp1a | Prl7c1 | Prl4a1 | 1700011m02rik | Cited1 | 1700031f05rik | 4933402e13rik |
| Prl3d1 | Gpr126 | Prl6a1 | Zfp655 | Plekha7 | Cited2 | Inhba | Dapk2 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Rnf2 | Arf2 | Cdh5 | Slc13a4 | Sfrp5 | Zfand2a | Inhbb | Gml1985 |
| Set | Tinagll | Fgd6 | Ceacaml4 | Ppplr3f | Krt25 | Helz | Fndc3b |
| Mrgprg | Mfi2 | Cysltr2 | Ceacaml5 | Obsll | Klk4 | Sele | Twsgl |
| Aa763515 | Rpn2 | Rhox6 | Trap1a | Slc23a3 | Tnfrsfl1b | Pdia6 | Aldhla3 |
| Tfpi | Abhd2 | Cdh3 | Ceacaml2 | Tmem87b | 2010204kl3iHk | Pdia5 | Lnx2 |
| Etosl | Hrctl | Spp2 | Gml65 15 | Epasl | Torlaip2 | Creb3 | Taf7 |
| Slc5a6 | Adm | Ziml | Ceacaml3 | Ccdc68 | Fmrlnb | Efnal | Ai844869 |
| 1600025ml7Hk | Abhd6 | Flnb | 4930447f24rlk | Kdelr2 | Ctsr | Dlg5 | Clecl2b |
| Gm9 | Slc7al | Rbbp7 | Gzmd | Pramefl2 | Ctsq | Procr | Prkcsh |
| Creb312 | Tead4 | Map3k7 | Foxj2 | Lrp8 | Prl8a2 | Fgfrl | Lama5 |
| Bbx | Mbnl3 | Rhox9 | Fbxll9 | Pard6b | Ctsm | Gnb4 | Tchh |
| Prl3cl | Gprl | Whsclll | Gzmc | Peg 10 | Prl8al | 23 10030g06\|k | Lamal |
| Mta3 | 2900057e 15rik | Slc38al | Gzmf | N4bp2 | Ctsj | Gcml | Rps6ka6 |
| Prl2al | Ldocl | 1600012pl7\|k | Gzme | Pla2g4e | Mpzll | Psgl8 | Vhl |
| Gm9 112 | Adaml9 | Adra2b | Gzmg | Fam78b | Stra6 | Goltlb | Eps812 |
| Afap 112 | Rybp | Pgf | Patl2 | Arrdc3 | Bcap3 1 | Psgl9 | Polg |
| Erlin2 | Col4al | 1200009i06rjk | 3830417al3itk | Pla2g4d | Cregl | Psgl6 | |
| Pard3 | Fndc3cl | Mfsd7c | Tspanl4 | Rassf8 | Tcfap2c | Slc2al | |
| Aifll | Col4a2 | Esam | Handl | Au015836 | Prl7bl | Psgl7 | |
| Dmrtcla | 4930502e 18rik | Gprl07 | AtxnlO | Csnkle | Ghrh | Htra3 | |
| 4932442108rik | Pkn2 | AuO 15791 | Mgat4a | Stagl | 4930486124rlk | Klhll3 | |
| Gjb2 | Rlim | Arhgap8 | Unc50 | Vnnl | Neurog2 | Ets2 | |
| Gjb5 | 1600015il0r\|k | Ankrdl7 | I12rb | Tchhll | 5430425j l2rlk | Nppc | |
| Slco5al | Afp | Cul7 | Ceacaml 1 | Plala | Prl7al | Tgml | |
| Wdr61 | Tmeml40 | 23 10067p03\|k | Plekhgl | Slc45a4 | Prl7a2 | Tmeml08 | |
| Kitl | Fstl3 | Irs3 | Prl3bl | Tex264 | Mirl 199 | Usp53 | |
| 9430027b09i)ik | Ing4 | Prl5al | Folrl | Pcdhl2 | TbcldlOa | Mark3 | |
| Tfrc | Taf71 | Fntb | A830080d0 lHkCtr9 | | Ralbpl | Cbx8 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Slc6a2 | Sultlel | Tceanc | Blzfl | Ccrlll | Pdgfra | Hspa5 | |
| Wdr45 | Olrl | Lepr | Zfp667 | Htatsfl | Morc4 | Spats2 | |
| Zxda | 2610019f03rlk | Tnfrsf9 | Fltl | 9030409gl1rik | Rarres2 | Limk2 | |
| Prdx4 | F11 | Papola | Usp27x | Tspan9 | Arid3a | Mkl2 | |
| Faml22b | Fbxw8 | Srd5al | Hdac4 | Rassf6 | Lifr | Shroom4 | |
| Zxdb | Sema4c | Clqtnfl | Itgb3 | 463 1402f24rlk | Shisa3 | Shrooml | |
| Zxdc | Ctnnbipl | Slc38a4 | Sri | A2m | Uevld | Pou2f3 | |
| Pip5kla | Tfpi2 | Angpt4 | Sema3f | Rimklb | Scnnlb | Acvr2b | |
| Placl | ZbtblO | Ctla2a | Prl3al | Locl005045$9 | Dnajbl2 | Rbms2 | |
| Igf2as | Mitf | 9930012kl1rik | Bahdl | Apob | Brwd3 | Atg4b | |
| Usp9x | Gpr50 | Mical3 | Sin3b | Tmeml50a | Hhipll | Pappa2 | |
| Psg28 | Hic2 | Apoa4 | Gm2a | 9130404d08|k | Fbln7 | Rbm25 | |
| Bmp8b | Tpbpb | Cul4b | Serpinb9g | Prl8a6 | Maspl | Gm4793 | |
| Fnl | Slc9a6 | 3632454122rik | Bend4 | Cts6 | Nrk | Nidi | |
| Psg23 | Prl7dl | Psg-psl | Bend5 | Prl8a8 | Pvr | Uba6 | |
| Bmp8a | Tpbpa | Lcor | Serpinb9b | Prl8a9 | Atp2cl | Lamcl | |
| Psg21 | Slco2al | Tnfrsf22 | Serpinb9c | Cts3 | Amot | Slc40al | |

| X reactivation | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gm21950 | Slc9a7 | Rhox3h | Slitrk4 | Fam47c | Zdhhc15 | Bhlhb9 | Samt1 |
| Gm21364 | Rp2 | Rhox2h | Ctag2 | Gm7173 | 1700121L16Rik | Gprasp2 | 4921511 M17Rik |
| Gm14346 | Jade3 | Rhox5 | 4930447F04Rik | Mageb16 | Magee2 | Arxes2 | Gm10057 |
| Gm14345 | Rgn | Rhox6 | Slitrk2 | Gm26775 | Pbdc1 | Arxes1 | Gm15140 |
| Gm14351 | Ndufb11 | Rhox7a | 1700036O09Rik | Tmem47 | Magee1 | Bex2 | 4930524N10Rik |
| Gm3701 | Rbm10 | Rhox8 | | 4930595 M18Rik | 5330434G04Rik | Nxf3 | |
| Gm3706 | Uba1 | Rhox7b | Gm1140 | Dmd | Cypt2 | Bex4 | Samt4 |
| Gm14347 | Cdk16 | Rhox9 | Gm14692 | Tsga8 | Fgf16 | Tceal8 | Samt2 |
| Gm10921 | Usp11 | Btg1-ps1 | 4933436I | | | Tceal5 | Cldn34b1 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gml0922 | Araf | Btgl-ps2 | OlRik | Fthll7a | Atrx | Bexl | Magea6 |
| Gm3750 | Synl | RhoxlO | Fmrlos | Tab3 | Magtl | Tceal7 | Magea3 |
| Gm3763 | Timpl | Rhoxl 1 | Fmrl | Gk | Cox7b | Wbp5 | Magea8 |
| Mycs | Cfp | Rhoxl2 | Fmrlnb | Gml4764 | Atp7a | Ngfrapl | Magea2 |
| Gml4374 | Elkl | Rhoxl3 | Gml4698 | Gml4762 | Tlrl3 | Kir3dl2 | Magea5 |
| Nudtl 1 | Uxt | Zbtb33 | Gm6812 | 5430427019Rik | Pgkl | Kir3dll | Mageal |
| AU022751 | Zfpl82 | Tmem255a | Gml4705 | Samt3 | Taf9b | Tceal3 | Cldn34b2 |
| NudtlO | Spaca5 | Atplb4 | Aff2 | NrObl | Fnd3c2 | Tceall | Satl |
| Bmp15 | Zfp300 | Lamp2 | 17001 11N16Rik | Mageb4 | Fndc3cl | Morf412 | Acot9 |
| Shroom4 | Ssxal | Gm7598 | 1700020N15Rik | Illrapll | Cysltrl | Glra4 | Prdx4 |
| Dgkk | Gm21876 | Cul4b | Ids | Gm27000 | Gm5 127 | Plpl | Ptchdl |
| Ccnb3 | 4930453H23Rik | Mctsl | 1110012L19Rik | Pet2 | Zcchc5 | Rab9b | Gml5 156 |
| Akap4 | Gm6938 | Clgaltlcl | 4930567H17Rik | 4932429P05Rik | Lpar4 | H2bfm | Gml5 155 |
| Clcn5 | Gm26593 | Gml4565 | BC023829 | 4930415L06Rik | P2ryl0 | Tmsbl51 | Phex |
| Usp27x | Agtr2 | 6030498E09Rik | Gm44 | | A630033H20Rik | Tmsbl5b2 | Sms |
| Ppplr3f | Slc6al4 | Cyptl5 | Mamldl | Gml4773 | Gprl74 | Tmsbl5bl | Mbtps2 |
| Ppplr3fos | Gm28269 | Cyptl4 | Mtml | Mageb2 | Itm2a | Slc25a53 | Yy2 |
| Foxp3 | Gm28268 | Gria3 | Mtmrl | Gm5072 | Tbx22 | Zcchcl8 | Smpx |
| Ccdc22 | K1M13 | Thoc2 | Cd9912 | Gm8914 | 2610002M06Rik | Faml99x | Gml5 169 |
| Cacnalf | Wdr44 | Xiap | Gml6189 | 1700084M14Rik | Fam46d | Esxl | K1M34 |
| Syp | Gm4907 | Stag2 | Hmgb3 | Gml4781 | Gm732 | Illrapl2 | Cnksr2 |
| Gml4703 | Gm4985 | Gm43337 | Gpr50 | Mageb5 | Gm379 | Texl3a | Rps6ka3 |
| Prickle3 | Gm27192 | Sh2dla | Vma21 | Magebl | Brwd3 | Nrk | Eiflax |
| Plp2 | Gm5934 | Tenml | Gml l41 | Magebl8 | Hmgn5 | Serpina7 | Map7d2 |
| Magix | Gm4297 | Gm362 | Prrg3 | Gm5941 | Sh3bgrl | 4930513006Rik | A830080DOlRik |
| Gpkow | Gm5935 | Dcafl212 | Fatel | 1700003E24Rik | Gm6377 | 4933428M09Rik | Sh3kbpl |
| Wdr45 | Gm5 169 | Dcafl211 | Cnga2 | | RP23- | Mumlll | Map3kl5 |
| RP23- | | | | | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 109E24.10 | Gml993 | Prr32 | Magea4 | BC061 195 | 240M8.2 | Trap1a | Pdhal |
| Praf2 | E330010L02Rik | 4930515L19Rik | Gabre | Arx | Pou3f4 | D330045A20Rik | Adgrg2 |
| Cedc120 | Gm5168 | Actrtl | MagealO | Polal | Cylcl | Rnfl28 | Gml5241 |
| Tfe3 | Gm2012 | Gm29242 | Gabra3 | Pcytlb | Gml0112 | Tbcld8b | Phka2 |
| Girap1 | Gm2030 | Smarcal | Gabrq | Pdk3 | Rps6ka6 | Gml5013 | Gml5243 |
| Kcndl | Six | Ocrl | Cetn2 | AU015836 | Hdx | Ripply1 | Ppefl |
| Otud5 | Gml4525 | Apln | Nsdhl | Gml4798 | RP23-466J17.3 | Cldn2 | Rsl |
| Pim2 | Gm6121 | Xpnpep2 | Gml4684 | Zfx | Tex16 | Morc4 | Cdkl5 |
| Slc35a2 | Gml0230 | Sash3 | Zfpl85 | Eif2s3x | 493340300 08Rik | Rbm41 | Gja6 |
| Pqbpl | Gm2101 | Zdhhc9 | Pnma5 | K1M15 | Apool | Nup62cl | Scml2 |
| Timml7b | Gml0058 | Utpl4a | Pnma3 | Fam90alb | Satll | Pihlh3b | Gml5262 |
| Gml0491 | Gm2117 | 9530027J09Rik | Xlr4a | Apoo | 2010106El ORik | Gml5046 | Rai2 |
| Gml0490 | Gm4836 | Bcorll | Xlr3a | Gml4827 | Zfp711 | Frmpd3 | Scmll |
| Pcskln | Gml0147 | Elf4 | Xlr5a | Magedl | Poflb | Prpsl | Gml5205 |
| Eras | Gm2165 | Aifml | Gml4685 | Gspt2 | Gml4936 | Tsc22d3 | Nhs |
| Hdac6 | Gml0096 | Rab33a | DXBayl8 | Zxdb | Chm | Mid2 | Gml5202 |
| Gatal | Gm2200 | Zfp280c | Xlr5b | RP23-9K14.6 | Dach2 | Eif2c5 | Reps2 |
| Glod5 | Gm26818 | Slc25al4 | Spin2d | Gm26617 | K1M4 | Tex13 | Rbbp7 |
| Gml4820 | Gm3669 | Gprl19 | Xlr3b | Spin4 | Ube2dnll | Vsigl | Txlng |
| Suv39hl | Gml0488 | Rbmx2 | Xlr4b | Arhgef9 | Ube2dnl2 | PsmdlO | Syapl |
| Was | E330016L19Rik | Gm595 | F8a | Amerl | 4930555B 12Rik | Atg4a | Ctps2 |
| Wdrl3 | Gml4632 | Enox2 | Xlr4c | Asbl2 | Cpxcrl | Col4a6 | SlOOg |
| Rbm3 | Gm7437 | Gml4696 | Xlr3c | Zc4h2 | H2afb2 | Col4a5 | Grpr |
| Rbm3os | Gml4974 | Gml4697 | Xlr5c | Zc3hl2b | Gml4920 | Irs4 | Rnfl38rtl |
| Tbcld25 | Gml0487 | Arhgap36 | RP23-95K12.13 | 1700010D OlRik | Gm28579 | Gml5295 | Apls2 |
| Ebp | Gm21447 | Olfrl320 | Zfp275 | Las11 | Tgif21x2 | Gml5294 | Zrsr2 |
| Porcn | Spin2f | 01frl321 | Gml8336 | | | Gml5298 | Car5b |
| | | | | | | | Siahlb |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Ftsj l | Gm2784 | Igsfl | Gm26726 | Msn | Tgif21xl | Gucy2f | Tmem27 |
| Slc38a5 | Gm2777 | 01frl322 | Zfp92 | F630028O lORik | Gml4929 | Nxt2 | Ace2 |
| SsxblO | Gm21883 | 01frl323 | Trex2 | Vsig4 | Pabpc5 | Kcnell | Bmx |
| Ssxb9 | Spin2e | 01frl324 | Haus7 | Hsf3 | Pcdhl 1x | Acsl4 | Pir |
| Ssxbl | Gm21608 | Stk26 | Bgn | Heph | H2afb3 | Tmeml64 | Figf |
| Ssxb2 | Gm21637 | Frmd7 | Atp2b3 | Gprl65 | Napll3 | Ammecrl | Piga |
| Gml4459 | Gm21645 | Rap2c | Dusp9 | Pgrl51 | Gml7521 | Rgagl | Asbl 1 |
| Ssxb6 | Gm2799 | Mbnl3 | Pnck | Eda2r | Cldn34cl | Chrdll | Asb9 |
| Ssxb3 | Gmclll | Hs6st2 | Slc6a8 | Ar | Astx6 | Pak3 | Mospd2 |
| Ssxb8 | Gm5926 | Usp26 | Bcap3 1 | Ophnl | Srsx | Capn6 | Fancb |
| Ssx9 | Gm2195 1 | 17000800 16Rik | Abcdl | Yipf6 | Gml7577 | Dcx | Gml7604 |
| Ssxb5 | Gm21657 | Gpc4 | Plxnb3 | Stard8 | Gml495 1 | A730046J 19Rik | Glra2 |
| Gm6592 | Gm21789 | Gpc3 | Sφ k3 | Efnbl | Astx2 | Algl3 | Gemin8 |
| Gm575 1 | Gm2825 | Gml4582 | Idh3g | Gml4812 | Gml7412 | Tφς 5 | Gpm6b |
| B630019K06Rik | Spin2-ps6 | A630012P 03Rik | Ssr4 | Gml4809 | Cldn34c2 | Tφ c5os | Ofdl |
| Fthll7b | Gm2863 | Cede 160 | Pdzd4 | Gml4808 | Gml4950 | Zcchcl6 | Trappc2 |
| Fthll7c | Gm2854 | Phf6 | L Icam | Pjal | Gml7467 | Lhfpll | Rab9 |
| Fthll7d | Gm2913 | Hprt | Arhgap4 | Tmem28 | Cldn34c3 | Amot | Tceanc |
| Fthll7e | Gm2927 | Gm28730 | Avpr2 | Eda | Astx5 | Htr2c | Egfl6 |
| Fthll7f | Gm2933 | Placl | NaalO | Awat2 | Vmn2rl2 1 | I113ra2 | Gml5226 |
| 4930402K 13Rik | Gm2964 | Faml22b | Renbp | Otud6a | Astxla | Lrch2 | Gml720 |
| Lancl3 | Gm21870 | Faml22c | Hcfcl | Igbpl | Gml7584 | Gml5 128 | Gml5230 |
| Gml4862 | Gm21681 | Mospdl | Iraki | Dgat216 | Astx4a | Gml5080 | Gm8817 |
| Xk | Spin2g | Etd | Mecp2 | Awatl | Gml7469 | Gml5 107 | Gml5232 |
| 1700012L 04Rik | Gm21699 | Gml4597 | Opnlmw | P2ry4 | Astx4b | Gml5 114 | Gml5228 |
| Gml4501 | Gml4552 | Cxxlc | Tex28 | Arr3 | Astxlb | Gm8334 | Tmsb4x |
| | Gml0486 | Cxxla | Tktll | Pdzdl l | Gml7361 | Gml5 127 | Tlr8 |
| | Gm2309 | | Flna | | | | Tlr7 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Cybb | Gml4553 | Cxxlb | Emd | Kif4 | Gm21616 | Luzp4 | Prps2 |
| Gm5 132 | Gml4819 | 4930502E18Rik | RpllO | Gdpd2 | Astx4c | Gml5099 | Gml5239 |
| Dynlt3 | Dockl 1 | 1700013H16Rik | Dnaselll | Gml4902 | Gml7693 | Ott | Frmpd4 |
| Hypm | I113ral | Zfp3613 | Taz | Dlg3 | Astxlc | Gml5092 | Msl3 |
| 4930557A04Rik | Zcchcl2 | Xlr | Atp6apl | Tex 11 | Gml7522 | Gml5093 | Arhgap6 |
| Sytl5 | Lonrf3 | Gml6405 | Gdil | Slc7a3 | Astx4d | Gml5 100 | Gml5261 |
| Srpx | Gm6268 | Gml6430 | Fam50a | Snxl2 | Gml7267 | Gml5085 | Amelx |
| Rpgr | Gml4569 | Slxll | Plxna3 | Foxo4 | Astx3 | Gml5086 | Hccs |
| Otc | Pgrmcl | 3830403N18Rik | Lage3 | Gm614 | 493241 IN23Rik | Gml0439 | Gml5245 |
| Tspan7 | Akapl7b | Gm773 | Ubl4a | Gm20489 | Gm382 | Gml5097 | Midi |
| Gml0489 | Slc25a43 | 1600025 M17Rik | Slcl0a3 | I12rg | 492151 1C20Rik | Gml5091 | 4933400A11Rik |
| Midlipl | Slc25a5 | Zfp449 | Fam3a | Medl2 | Cldn34c4 | Gml5 104 | Gml5726 |
| Gml4493 | Gml4549 | Gm2155 | Ikbkg | Nlgn3 | 4930558G05Rik | Tmem29 | Gml5247 |
| Gml4483 | 23 10010 G23Rik | Smiml012a | G6pdx | Gjbl | Diaph2 | Apex2 | Gm21887 |
| Gml4474 | C330007P06Rik | Gm2174 | Gm6880 | Zmym3 | Pcdhl9 | Alas2 | Asmt |
| Gml4477 | Ube2a | Ddx26b | 01M326-psl | Nono | Gm2685 1 | Pfkfbl | |
| Gml4476 | Nkrf | Gml0477 | 01frl325 | Itgblbp2 | Tnmd | Tro | |
| Gml4484 | Gml5008 | Gm648 | Gm5640 | Tafl | Tspan6 | Maged2 | |
| Gml4479 | 43349 | Mmgtl | Gm6890 | Ogt | Srpx2 | Gm27191 | |
| Gml4482 | Sowahd | Slc9a6 | Gm5936 | Cxcr3 | Sytl4 | Gnl31 | |
| Gml4478 | Rpl39 | Fhll | Gab3 | Gm4779 | Cstf2 | Fgdl | |
| Gml4475 | Upf3b | Mtap7d3 | Dkcl | 8030474K03Rik | Noxl | Tsr2 | |
| Gm4906 | Nkap | Adgrg4 | Mppl | Nhsl2 | Xkrx | Gml5 138 | |
| Bcor | Akapl4 | Brs3 | Smim9 | Rgag4 | Arll3a | Wnk3 | |
| Gml4635 | Ndufal | Htatsfl | F8 | Pin4 | Trmt2b | A230072 ElORik | |
| Atp6ap2 | Rnfl 13al | Vglll | Fundc2 | Ercc61 | Tmem35 | Faml20c | |
| 1810030007Rik | Gm9 | | Cmc4 | Rps4x | | Phf8 | |

| | | | | | | |
|---|---|---|---|---|---|---|
| Medl4 | Rhoxl | Gml4718 | Mtcpl | Citedl | Cenpi | Huwel |
| Usp9x | Rhox2a | Cd401g | Brcc3 | Hdac8 | Drp2 | Hsdl7bl0 |
| 2010308F09Rik | Rhox3a | Arhgef6 | Vbpl | Phkal | Taf71 | Ribcl |
| Ddx3x | Rhox4a | Rbmx | Gml5384 | Gm91 12 | Timm8al | Smcla |
| Nyx | Rhox3a2 | Gm364 | Rab39b | Dmrtclb | Btk | Iqsec2 |
| Cask | Rhox4a2 | GprlOl | Gml5063 | Dmrtclcl | Rpl36a | Kdm5c |
| Gpr34 | Rhox2b | Zic3 | Pls3 | Dmrtclc2 | Gla | Kantr |
| Gpr82 | Rhox4b | 4930550L24Rik | Gml4715 | 170003 IF05Rik | Hnrnph2 | Tspyl2 |
| Gm5382 | Rhox2c | Fgfl3 | Gml4707 | Dmrtcla | Armcx4 | Gprl73 |
| Gml4505 | Rhox3c | F9 | Gml4717 | 170001 1M02Rik | Armcxl | Cldn34a |
| Drrl | Rhox4c | Mcf2 | Cldn34b3 | Napll2 | Armcx6 | Shroom2 |
| Cyptl | Rhox2d | Atpl 1c | Cldn34b4 | Cdx4 | Armcx3 | Gprl43 |
| Maoa | Rhox4d | Gm7073 | Cldn34d | Chicl | Armcx2 | Usp5 1 |
| Maob | Rhox2e | Gml4661 | Tbllx | Gm26952 | Nxf2 | Magehl |
| Ndp | Rhox3e | Sox3 | Prkx | Tsx | Zmatl | Foxr2 |
| Efhc2 | Rhox4e | Gml4662 | Gml4742 | Gm26992 | Gml5023 | Rragb |
| Fundcl | Rhox2f | Gml4664 | Pbsn | Tsix | Tceal6 | Klf8 |
| Dusp21 | Rhox3f | Cdrl | Gml4744 | Xist | Pramel3 | Ubqln2 |
| Kdm6a | Rhox4f | Ldocl | 5430402E1ORik | Jpx | Gm5 128 | Cypt3 |
| 4930578C19Rik | Rhox3g | 4933402E13Rik | Obpla | Ftx | Gm7903 | Kctdl2b |
| Gm26652 | Rhox2g | 493 14000 07Rik | Gm5938 | Zcchcl3 | AV320801 | RP23-106P7.5 |
| BC0497002 | Rhox4g | 1700019B21Rik | Obplb | Slcl6a2 | Nxf7 | 2210013021Rik |
| Chst7 | | Gm6760 | Gml4743 | Rlim | Prame | Spin2c |
| | | 3830417A13Rik | 4930480E11Rik | C77370 | Tcpl 1x2 | |
| | | | Prrgl | Abcb7 | Tmsbl5a | |
| | | | | Uprt | Armcx5 | |
| | | | | | Gpraspl | |

| XEN | | | | | | | |
|---|---|---|---|---|---|---|---|
| Dab2 | Pdgfra | Gata6 | Fxyd3 | Soxl7 | Lamal | Gata4 | Krt8 |
| Fst | Pthlr | Foxql | Tet3 | Foxa2 | Lambl | | |

| Trophoblast | | | | | | | |
|---|---|---|---|---|---|---|---|
| Ascl2 | Cdx2 | Esrrb | Grn | Lipg | Smad3 | Tfap2c | Gata3 |
| Bmp4 | Elf5 | Ets2 | Igf2 | Pcsk6 | Snai1 | Vav1 | Krt7 |
| Bmp8b | Eomes | Fgfr2 | Jade1 | Ptpra | Tead4 | Yap1 | Krt18 |

| Trophoblast progenitors | | | | | | | |
|---|---|---|---|---|---|---|---|
| Rhox6 | Hmgn2 | Tuba1b | Immt | Rps21 | Ccnd3 | Mrpl54 | Ruvbl2 |
| Rhox9 | Odc1 | Cenpw | Smagp | Pdlim2 | Rpl5 | Rps26 | Ndufv1 |
| 3830417A13Rik | Klhl13 | Cct7 | Hnrnpa2b1 | Rpl24 | Nip7 | Ndufb9 | Polr2l |
| Gjb3 | Ncl | Sfn | Cox7b | Asf1a | Psma5 | Arpc1a | Asns |
| Gm9112 | Tyms | Fkbp4 | Snx10 | Eif4a3 | Spc24 | Rps28 | Prkrip1 |
| Hspb1 | Prss8 | Ndufb6 | Stip1 | Ssb | Mdh2 | Prpf31 | 1700021F05Rik |
| Nup62cl | Atp5g3 | Snrpe | Rnf4 | Timm17a | Cep164 | Mrpl12 | Aimp1 |
| Ldoc1 | Dusp9 | Cenph | Gm648 | Mrpl18 | Cs | Epop | Rps7 |
| Hspe1 | Gmnn | Rad51 | Cct6a | Cenpk | Zc3h15 | Cct5 | Tra2b |
| Rhox12 | Rrm2 | Set | Snrpd2 | Dcakd | Pea15a | Pdap1 | Cox17 |
| Tex19.1 | Tbrg1 | Cd164 | Psmg2 | Hikeshi | Tsen15 | Ezh2 | Mrpl19 |
| Gjb5 | Cct3 | Cox6b1 | Tk1 | U2af1 | Ippk | Gpbp1 | Chchd4 |
| Sin3b | Nhp2 | Hnrnpdl | Rps5 | Acp1 | Thoc3 | Psme3 | Polr1d |
| 1700086L19Rik | Ppid | Lsm2 | Mtx2 | Tipin | Pithd1 | Ube2c | Ubfd1 |
| Ldhb | Ccna2 | Exoc3l4 | Phb | Fkbp3 | Pak1ip1 | Cbx1 | 2410015M20Rik |
| Krt19 | Anp32b | Dut | Hspa8 | Cdca3 | 1110038B12Rik | Gata2 | Tbcb |
| Hmgn5 | Cacybp | Pramef12 | mt-Nd5 | Tubb4b | Wdr18 | Nxf7 | Chchd1 |
| | Chchd2 | | | Mycbp | | Smc4 | |

| Trap1a | Phb2 | Cd320 | Orc6 | Apip | Nol7 | Tfap2c | Serbp1 |
|---|---|---|---|---|---|---|---|
| Plac1 | Snrpf | Snrpd3 | Dctpp1 | Mdk | Tomm70a | Creb3 | Hsph1 |
| Cdkn1c | Ran | Psmb7 | Sugt1 | Rpl14 | Snu13 | Clns1a | Xpo1 |
| Bex1 | Gale | Mcm7 | Wdr77 | Cox7a2 | Psma2 | 1810022K09Rik | 2310033P09Rik |
| Fthl17a | mt-Nd4 | Taf1d | Suclg1 | Hnrnpc | Eif2s2 | Eif2b1 | Prpf19 |
| Dbi | Birc5 | H2afz | Ddx39 | Sdr39u1 | Usmg5 | Idh3a | Apoo |
| Ube2a | Tpm2 | Ndufb2 | Polr2f | Slc25a3 | Eif3e | Sae1 | Hagh |
| Dnaja1 | Hsd17b4 | Lyar | Rpl38 | Psma7 | Cops5 | Eif5a | Ndufa9 |
| Phactr1 | Rpl22l1 | Rbms2 | Rpa2 | Psmd12 | Mrpl3 | Fhl2 | Mrpl2 |
| Phlda2 | Snrpd1 | Eif5b | Fmr1nb | Cyc1 | Mybbp1a | Lap3 | Ndufb7 |
| Hand1 | Hspa14 | Rbm8a | Gng12 | Apex1 | Elp2 | Ncbp2 | Psmb1 |
| Selenoh | Wfdc2 | Dynll1 | Tuba1c | Rad23b | 1110004F10Rik | Eps8l2 | Txndc9 |
| Rhox5 | Rfc4 | Stmn1 | Aasdhppt | C1qbp | St13 | Cdk4 | Hnrnpa1 |
| Atp5g1 | Rgcc | Got2 | Pfdn6 | Cox6c | Tbca | Rfc3 | Ndufs7 |
| Hmgn1 | Mfsd2a | Cox7c | Hspa9 | Txn1 | Snrpa1 | Cdk1 | Farsb |
| Hat1 | Cct8 | Lsm6 | Eif1a | Med19 | H2afv | Mrps25 | Cycs |
| Plet1 | Ubxn1 | Ccne2 | Pop5 | Slirp | Mcm5 | Coq3 | Tmem11 |
| Gm9 | Ddt | Sap18 | Nasp | G3bp1 | Tcp1 | Med10 | Rps17 |
| Rbbp7 | Dtymk | Liph | Xlr4b | Ak2 | Atp1b1 | Emd | Mrpl14 |
| Hspd1 | C430049B03Rik | Pa2g4 | Snrpb2 | Krt18 | Aprt | Ptrh2 | Diablo |
| Mrfap1 | Magoh | Slc38a4 | Nop58 | Rsl1d1 | Nup37 | Mrps18c | Cox4i1 |
| Krt7 | Calm2 | Irx3 | Uqcrc2 | Csrp1 | Hebp1 | Med4 | Pkp2 |
| Esam | Mrps22 | Srsf3 | Cfdp1 | 1600025M17Rik | Lsm8 | Fam133b | Psmc2 |
| Krt8 | Impdh2 | Dpy30 | Hn11 | Rpp30 | Mbd3 | Crip2 | Psmc1 |
| Fstl3 | Brd3 | Hmgc1 | Tsn | Mrpl38 | Gtf3c6 | Ndufaf3 | Slc25a4 |
| Ghrh | Fscn1 | Cenpa | Psma6 | Emg1 | Rpa3 | Thap4 | Eloc |
| Ranbp1 | 2610528J11Rik | Mgll | Ssrp1 | Cebpzos | Cdc34 | Mrps16 | Vma21 |
| Npm1 | | Eef1g | | | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| H19 | Zwint | Atp5cl | Acaala | Nsmce4a | Ndufb8 | Uchl3 | Mif |
| Sdcl | Tmem37 | Imp4 | Rpf2 | Cct2 | Naplll | Meal | Timml3 |
| Rps41 | Ndufa5 | Cks2 | Lgalsl | Rpsl6-ps2 | Adgrf5 | Psma3 | |
| mt-Ndl | Eif2sl | Rnd2 | Psmd6 | Ruvbll | Ptges3 | TimmlO | |
| Hsp90aal | Hsdl7b2 | Knstrn | Aplm2 | Arppl9 | Polr2j | Rrml | |
| Mbnl3 | Galkl | Atp5fl | Plppl | Rpl27 | Ndufal2 | Hnmpd | |
| Htatsfl | Cct4 | Skpla | Ndufaf2 | Dcunld5 | Cyb5b | Tomm22 | |
| Hsp90abl | Cox5a | Igf2bpl | Cull | Rpll8 | Tmod3 | Ndufabl | |
| Las 11 | Dkkll | Mrpl21 | Ndufal 1 | Mrpll5 | Ndufv2 | Aifml | |
| Ptma | Hmgb2 | Srsf7 | mt-Col | Psmal | Ash21 | Tfam | |
| mt-Cytb | Tubb5 | Psipl | Tomm40 | Baspl | Spc25 | Rrpl5 | |
| Snrpg | Med21 | Llph | Ndufs8 | Tead2 | Dnajc2 | Rps2 | |
| Fdxl | Nmel | Erdrl | Derl3 | Prmtl | 4921524J17Rik | Tinf2 | |
| Glrx5 | Cdca8 | Atp5k | mt-Nd2 | Esfl | Gins4 | Lypla2 | |
| Alpl | Tsen34 | Rmdn3 | Ckslb | Banfl | Naa38 | Ppmlg | |
| Elf3 | Oaf | Peg 10 | Eif3g | Pinl | Pole3 | Dars | |
| Ndufa4 | Ccnbl | Ccnel | Nop 16 | Mta3 | Nucb2 | Ingl | |
| Dynll2 | Ascl2 | Rps271 | Itpa | Priml | Tomm7 | Psmb2 | |
| Hsp25-psl | Lsm4 | Ezr | Mat2a | Ppih | Erh | Fcfl | |
| | Ahsal | Psmd7 | Gnl3 | Eif3i | Rps8 | Rpl30 | |
| | | | Pdcd5 | | Samm50 | | |

| Spiral Artery Trophpblast Giant Cells | | | | | | | |
|---|---|---|---|---|---|---|---|
| Car2 | Psg22 | Rgsl7 | Psipl | Eif31 | Got2 | Rpsl8 | Cct6a |
| Set | K1M13 | Mpzl2 | Tnfaip8 | Fscnl | Hnmpa2b 1 | Actr3 | Nectin 2 |
| 1500009L16Rik | Ldocl | Liph | Trap 1a | Ehdl | Prl7dl | Anxa7 | Grhpr |
| Serpinb9e | Galkl | Ddbl | Tubalc | Pramefl2 | 1110008P14 Rik | Cfll | Cct7 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Prl2al | Aφς1b | Irs3 | Cd82 | Eiflb | Rackl | Gtf2e2 | Chordcl |
| S100a6 | Anxa4 | Bexl | Gjb5 | Mxd4 | Rps7 | Parva | Vma21 |
| Plac8 | Cdx2 | Lysmd2 | Seφiηe2 | Rapla | Pdcd5 | Eeflg | Rpl39 |
| Serpinb9g | Tpm4 | Rpl2211 | Tubala | Borcs7 | Cct4 | Cct2 | Ccnbl |
| Prl6al | Anxa2 | Rhox5 | Txnl | Torlaip2 | Mif | Rpl9 | likGm2000 |
| Lgals9 | Seφinb9b | 2310030G06:likRalbpl | C430049B03¾kAvpil | Krtl9 | Csφï | 0610007P14I:lik | Snrpf |
| Prl7bl | Derl3 | Pdlim2 | H2afz | Actgl | Cox5a | Nmrkl | Aamp |
| Ada | Tfap2c | Nostrin | Pdcd4 | Cdkn2aipnl | Rpl27 | Eny2 | Smarcbl |
| Aldhla3 | Baspl | Glrx5 | Jup | Bex3 | Npml | Epop | Prelid1 |
| Serpinb6b | Rbbp7 | Tpml | Morf412 | Dnajc8 | Ppdpf | Ran | Paklip1 |
| Sri | Caldl | Cnn2 | Pfnl | Ubfdl | Ets2 | Krtl8 | Hmbs |
| Fstl3 | Laspl | Grb2 | Actnl | Cfap20 | Nrk | Kat7 | Polr2j |
| Serpinb9d | Hmgn5 | Fbliml | Aifll | Zwint | Gga2 | Exosc8 | Calm3 |
| Prl2c5 | Spata21 | Uppl | Cdh5 | Rps4x | Krt7 | Rpl23a | Ezr |
| H19 | Tbrgl | Ppplrl4b | Eif4ebpl | Mycbp | Ranbpl | Rps8 | Rps3al |
| Aprt | Dusp9 | Cdknlc | Erccl | Ndufaf3 | Rps41 | Rps3 | Elovl5 |
| Seφinb9c | TmsblO | Tfpi | Mvp | As3mt | Ywhab | Rrm2 | Rpsl7 |
| Ascl2 | Dynll2 | Fermt2 | Ndufal1 | Hatl | Fkbpla | Dtymk | Rps5 |
| Placl | Ctnnbipl | Palm | Ugp2 | Rps20 | Pdcl3 | RpllOa | |
| Mt2 | Sin3b | Tubb5 | Prmt5 | Myl6 | Rpsl6 | Actr2 | |
| Fthll7a | Igfbp7 | SlOOall | 1700086L19I:likPygl | K1M22 | Gnai3 | Olal | |
| Tφ53i11 | Mpzll | Krt8 | 1600025M17¾lRpp21 | Cetn3 | Eif4e3 | Cklf | |
| Mrfapl | Olrl | Zyx | Aφς2 | I12rg | Rpll2 | Cfdpl | |
| Phactrl | Mbnl3 | Alad | Abracl | Pletl | Tipin | RpslO | |
| Tnfrsf9 | Myll2a | Faml62a | Vasp | Gm9112 | Aφς5 | Rpl36a | |
| Lgalsl | Nek6 | AA467197 | Gngl2 | | Eif2sl | Rpsl9 | |
| Pitrml | Sbsn | Rps271 | Sqstml | | Chpl | Snφg | |
| Ncmap | Copz2 | Ncaml | | | Cepl64 | Clqtnf6 | |

| | Eif2s2 | Dcakd | Tpm2 | Eifla | Rpsa | Atpifl | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |

| Spongiotrophoblasts | | | | | | | |
|---|---|---|---|---|---|---|---|
| Phlda2 | Cs | Pttg1 | Cops5 | Lsm8 | Impa2 | Drg1 | Mrto4 |
| Dio3 | Lgalsl | Trappc5 | Psmd12 | Gadd45g | 2010107E04Rik | Nae1 | Rnf128 |
| Dkkll | Hagh | Eif3g | Panx1 | Med7 | Ndufb5 | Hspa8 | Wdr77 |
| Hspb1 | Npm1 | Gpx4 | Dld | 2310033P09Rik | 0610007P14Rik | Dars | Pepd |
| Tmem14c | Tex30 | Gtf2h5 | Ppid | Atp11a | Gtf3c6 | Ubald2 | Ddx18 |
| Cidea | Mfge8 | Magoh | Dnajc2 | Skp1a | Dnajc19 | Hnrnpk | Lrrfip2 |
| Tfrc | Usp1 | Fam50a | Hspd1 | Eloc | Atp5k | Idh3a | Psmb7 |
| Batf3 | B3gnt7 | Cct3 | Hmgb2 | Nsmce2 | Tubb2a | Plekhf2 | Erdr1 |
| Sin3b | Mageh1 | Srsf3 | Uaca | Slc25a3 | Slirp | Vps35 | Rps28 |
| Prss8 | mt-Nd4 | Rfc4 | Wwtr1 | Gadd45b | Phb2 | Mrpl47 | Fnta |
| Ldoc1 | Emc8 | Eif1a | Psmd6 | Cfdp1 | Psmc1 | Birc5 | Rtn3 |
| Maoa | mt-Nd5 | Marcksl1 | Hnrnpc | H2afz | Folr1 | Unc50 | Idh3b |
| Cdkn1c | Commd4 | Serpinb9e | Mrps23 | Ppa1 | Bax | Dut | Elob |
| Las11 | Dnaja2 | Apoo | Nap1l1 | Atp5b | Rmdn3 | Cdc34 | Pfdn6 |
| Rhox6 | Tbca | Slc2a1 | Tead2 | Polr2e | G3bp1 | Nabp1 | Sugt1 |
| Tex19.1 | Ndufb2 | Vdac3 | Cd164 | Clns1a | Trim27 | Hadhb | Dstn |
| 2610528J11Rik | Tubb4b | Cox5a | Pparg | Dnajb6 | St13 | Aimp1 | Smarcb1 |
| Gkap1 | Sct | Ppp1r3g | Rpl22l1 | Rnf181 | Slc38a2 | Fus | Coq3 |
| Cldn7 | Ing2 | Cct5 | Rhox5 | Rnf4 | Dusp9 | Etfb | Igsf8 |
| Slc22a18 | Cd320 | Anxa4 | Psmd7 | Hdac1 | Cggbp1 | Hnrnpab | Tomm22 |
| Rhox9 | Hsd11b2 | Nsmce4a | Ndufa4 | Prpf19 | Ptma | Ndufb4 | Hmbs |
| Mrps6 | Vamp8 | C430049B03Rik | Ndufb6 | Nsmce1 | Chchd1 | Exosc8 | Cyc1 |
| Serpinb9g | Tbrg1 | Tmem147 | Tma7 | Gm11361 | Rpl18 | Rplp1 | Txnl1 |
| Aqp3 | mt-Nd2 | Pa2g4 | Med21 | mt-Rnr1 | Psmc6 | Cox7b | Fam104a |
| | Gm9 | | Cox6b1 | | | Mrpl19 | Hn1 |

166

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| mt-Cytb | Slc38al | Tyms | Tardbp | Ncbpl | Atp5cl | Nsfllc | Ctnnal |
| Hsp25-psl | Rbbp7 | Eif4al | Uqcrc2 | Blvra | Eroll | Timml7g | Ndufs8 |
| Rdh12 | AtxnlO | Snrpe | Psma6 | Pφ sap 1 | Hspa9 | Pigp | Bsg |
| Krtl8 | Hsp90aal | Smul | Larp7 | Ube2el | Anapcl5 | Ndufsl | Gskip |
| Pfdnl | Calml | Tbcb | Ranbpl | S100al6 | Rps8 | Appbp2 | Cnihl |
| Tulpl | Hspel | Baspl | Mφ 14 | Serbpl | Seφ inb9d | Zwint | Rbm8a |
| Selenoh | Faml36a | Fam90alb | Suclgl | RablO | Cotll | Duspl 1 | Gm2a |
| Dynll2 | Elf3 | Nup85 | Pgrmcl | Rala | Ash21 | Mcm2 | Eif3e |
| Glrx5 | Prkd2 | Lonp2 | Mdh2 | Psmdl3 | Arl6ipl | Set | Erh |
| Slcl6al | mt-Col | Mφ s22 | Rpl5 | Pmpca | Borcs7 | Scarb2 | Naa35 |
| Krt8 | Ncl | Lyar | Ndufa5 | Serpinb9b | Psmc2 | Smc4 | Mφ 13 |
| Tmeml50a | Hadh | Fermt2 | Gucdl | Ppa2 | Zcchcl7 | Ywhaq | Mapllc3b |
| Stx3 | Cisdl | Srsf6 | Car2 | Hebpl | Ncbp2 | Cdca8 | Tcpl |
| Gjb2 | Snrpg | Nxf7 | Dnajc9 | Mφ 115 | Psmbl | Hmgcl | SrsflO |
| Nudt22 | Syngrl | Rad23b | Wdrl8 | Rrm2 | Priml | Tra2a | Psma3 |
| Mbnl3 | Chchd2 | Fkbp3 | Cox7c | Ccnbl | Thoc3 | Npepll | Ndcl |
| Gm9 112 | Ubqlnl | Atp5o | Ssb | Gprl37b | Nop58 | Med28 | Mtch2 |
| Cd9 | Fbxll9 | Cct8 | Ran | Idh3g | Polrld | H2afv | Psmdl 1 |
| Rbpl | Pphlnl | Snx5 | Emd | Srsf7 | Sap 18 | Sdhb | Rpl27 |
| Rps41 | Slc25a5 | Clqbp | Hsp90ab1 | Slc25a4 | Gmfb | Uqcrcl | E2f5 |
| Eif2s2 | Ccdc5 1 | Bglap3 | Hnmpal | Gata2 | Lsm4 | Nsφ ī | Pitpnb |
| Ugp2 | Mpdul | Atp5fl | Atp5al | Nhp2 | Rps5 | Snipf | |
| Zfp655 | Eif2sl | ChchdlO | Psmg2 | Rars | Cdipt | Snφ d2 | |
| mt-Ndl | Hspal4 | Olrl | Pdcd5 | Snx6 | Uspl4 | Rabif | |
| Tdφ | Prkcz | Cenph | Cacybp | Dpy30 | Psme3 | Commd5 | |
| Urod | Tafld | Uchl3 | Lsr | Ube2c | Lamtorl | Smiml 1 | |
| Hmgn5 | Mipll6 | Cenpk | Ttc4 | Ahsal | Cycs | Cox4il | |
| | 170002 1F0 | Paklipl | | Peg 10 | Ndufb8 | | |

167

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Car4 | 5Rik | Gml5536 | Cox7a2 | Eif3i | Imp4 | Cetn3 | |
| Krtl9 | Rap2c | Naa38 | Lsm6 | Mφ 155 | Mφ s25 | Ruvbl2 | |
| Rassf6 | Acvr2b | Tφ ïï | Stmnl | Rfc5 | Nop 16 | Strap | |
| Tfeb | Irx3 | Psmc5 | Ccna2 | Cystml | Eif3d | Txnl | |
| Hbegf | Placl | Got2 | Uchl5 | Ndufaf2 | Sael | Cyb5r3 | |
| Rab9 | Abhd5 | Syce2 | Gadd45gipl | Cox 14 | Uqcrfsl | Szrdl | |
| Dnajal | Seφ iηe2 | Atp5g3 | Epop | Usp39 | Ilf2 | Eeflg | |
| Fhl | Snrpd3 | Atplbl | Ndufb9 | Hatl | Rad5 1 | Ndufs7 | |
| Atp6v0dl | Prss36 | Maea | Txndc9 | Lysmd2 | Psmc3 | Mφ 145 | |
| Impdh2 | Perp | Psmal | Slc38a4 | Psma7 | Hnrnpdl | Samm50 | |
| Aplm2 | Tmeml09 | Ddx39 | Rbbp4 | Pole3 | Brixl | Fdxl | |
| Sod2 | Cct6a | Tmeml 16 | Lgalsl | Renbp | Cox6c | Ndufvl | |
| Slc26a2 | 3830417A13Rik | Nasp | Psmfl | Mrpl41 | Ddt | Siupal | |

| **Oligodendrocyte precursor cells (OPC)** | | | | | | | |
|---|---|---|---|---|---|---|---|
| Sppl | Mcm3 | S100a3 | Rassf4 | Adam9 | Irfl | Col23al | Mmp2 |
| Ccnbl | Pgcp | Creb5 | Nt5dcl | Mnsl | Kif20b | Col4a5 | Plekhbl |
| Pdgfra | Neu4 | Tram2 | Kif23 | Bean | Tcn2 | Cdldl | Slc7al 1 |
| Den | Emp3 | Serpinfl | Troap | Zfp3611 | Rnfl80 | Pcdhga5 | Cenpl |
| Rlbpl | Slc6a20a | Enppl | Slc25a29 | Ssfa2 | Slc38a3 | Gal3stl | 1118 |
| Slc6al3 | Igf2 | Tacc3 | Epn2 | Tnfrsfl 1b | Lgals2 | Ddah2 | Alpl |
| Inmt | Kif2c | Spry4 | Qpct | Gpr81 | 17001 12 E06Rik | Alx3 | Cede 18 |
| Pnlip | Zcchc24 | Loxl3 | Gml9705 | Tmeml46 | Neil3 | 4921530 L18Rik | Fam35a |
| Lum | Mxra8 | Cyplbl | Timp4 | Kctdl2b | 2900005J 15Rik | Frmd8 | 20103 17 E24Rik |
| Cmbl | Ampd3 | Htra3 | Jun | Col9a3 | Clgn | Gprl46 | Fdxr |
| Pcolce | Ccnb2 | Ccl5 | Cxcll2 | Ostfl | Cercam | Phldb2 | Medl8 |
| Postn | Chstl 1 | Ezh2 | Col3al | D2Ertd75 Oe | 6720463 M24Rik | Itf**g3** | MtmrlO |
| Apod | Kif20a | Agbl2 | Rfx4 | Fbxo7 | LOC6266 93 | Trim45 | E130309 F12Rik |
| Ednrb | Musk | Maml2 | Ppfibpl | Clecla | Ehd2 | Cdk4 | 111003 11 02Rik |
| Scrgl | SlOOb | Klhl5 | Cyr61 | Gpx7 | Thbsl | Itga9 | Hells |
| Tmem45a | mt_AK13 1586 | Frmd7 | Zebl | Atp6v0e | | Prtg | Tφv 4 |
| Fam70b | Efempl | Ccl2 | Ppic | Cdkl | | Cdk5rap2 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Cspg4 | Gpc5 | Fam70a | Rhoc | Pcyoxll | Cd302 | Arhgapl9 | Cyp20al |
| Cacng4 | Tmeml76b | Abtb2 | Abhd2 | Caprin2 | Coll5al | 4930517E1 IRik | Col4al |
| Fabp7 | Shc4 | Fkbp9 | Traf4 | Pabpc5 | Plekhg6 | Rasll 1a | Antxrl |
| Pbk | Gm2a | Cenpe | Tspan4 | Fzd6 | Creb313 | Tubalc | Aldhlal |
| 1110015018Rik | SlOOal | Slc2al2 | Cpxml | Gm5089 | Map3k8 | Islr | Gabl |
| Emidl | Galnt3 | Slc22a8 | SoxlO | Cenpf | Timp3 | Prrxl | 13000141 06Rik |
| Serpingl | S100al6 | Ladl | E1301 14P18Rik | Mmpl l | Akapl3 | Rrm2 | 9930021D14Rik |
| Oligl | Clqtnf6 | Clqtnf2 | Mfsd2a | Rasa3 | Arhgap29 | Pars2 | Tmem220 |
| Vtn | Afap 112 | Ccndl | Lrp4 | Gsn | Melk | Cftr | Rhpnl |
| Prcl | Lbp | Lamal | Fos | Gm9839 | Antxr2 | Slcl3a5 | Tmeml98b |
| Faml80a | Cdkn2c | Smc4 | Tpx2 | Sall3 | Bmp7 | Lgals3bp | Ebfl |
| E130306D19Rik | Vipr2 | Adamtsl3 | Cenpi | 1810034E14Rik | Rabl3 | Cklf | Ssl8 |
| Bgn | Chst5 | Vegfc | Lame 3 | Gpr3711 | Tsgal4 | Col4a2 | E2f8 |
| Lmcdl | Gpx8 | S100a6 | Mapk7 | Tril | Smpd2 | Vamp5 | Faml 11a |
| Colla2 | Pdpn | Kankl | Lama2 | Jam2 | Abca6 | Rassf8 | Tgfbr3 |
| Spc25 | Lims2 | Irak4 | Fosb | Evi51 | Gatm | Faml32a | Sema5b |
| Calcrl | Mavs | Sh3bp4 | Susd5 | Dna2 | Slitrk6 | Rftn2 | Ifitm3 |
| Itih5 | Aurka | Btd | Dpyd | Serpina3n | Snx22 | Dill | Gdpd2 |
| TmemlOO | Empl | Mc5r | Uhrfl | Cdc20 | Mpzll | Caldl | Cfh |
| Adm | 01ig2 | Rnf43 | Plekho2 | Sulfl | Prkcq | A430107013Rik | Nnat |
| Tmeml76a | Aox3 | Collal | Tmc6 | P2rx7 | 4933425H06Rik | Fam82al | D930014E17Rik |
| 0610040JOIRik | Mytl | Bcasl | Apobec3 | Map3kl | Gprc5a | Tcirgl | Mcm9 |
| Pmel | Fignll | Plkl | Faml 14al | Dab2 | Pcca | Nusapl | Gins2 |
| A930009A15Rik | Pcdhgc3 | Notchl | Birc5 | Clqtnf7 | Prelp | Gprl82 | Slcla5 |
| Cavl | Gpsm2 | Angptll | B3gnt5 | Kif22 | Gnb4 | Serpindl | Ptgds |
| Nuprl | Mir568 | Cdca8 | Itgb8 | Xlr3b | Cyp2j6 | Mcm7 | Tnpol |
| Gstm2 | Cd9 | Mc4r | Stonl | Kifl8a | Ctdspl | Sgk3 | Ifitm2 |
| Ckap2 | Fanci | Gpt2 | Kcnj 1O | Zfp3612 | Rab34 | Lekrl | Notch2 |
| Spryl | Fam64a | mt_AKl43357 | 363245 l006Rik | S100a4 | Fzd9 | Srpx2 | Luzp2 |
| Top2a | Zic4 | Hapln3 | Socs3 | Seel | Msh6 | Gpldl | Mure |
| 1190002F15Rik | Cd40 | Lpo | Tmeml44 | A330041J22Rik | Cep72 | 1700013G23Rik | |
| Ube2c | Meoxl | Hpsl | Ptgfr | Plat | Otos | Icaml | |
| Ccl7 | Ect2 | Boll | Slcl6al2 | Fam71f2 | Anxa2 | Jam3 | |
| Cp | Rcn3 | Sema3d | Chaflb | Smocl | Ftsjdl | mt_AK159184 | |
| | Cyp2j9 | S100al3 | Dbi | Sox8 | Saal | Coblll | |
| | 1190002H23Rik | Nuf2 | Gfral | Hmgb2 | Sh3tc2 | Trafl | |
| | Wipfl | Ggt5 | Cdca2 | Bmp6 | Rnpepll | | |
| | | Meisl | | | Atpla2 | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Vcan | Poldl | Cenpn | Gpr82 | Pomtl | Pion | Mmd2 | |
| Ugdh | 1810010H24Rik | Spsb4 | Nhsll | Orail | Ppplrl4b | Sulf2 | |
| Mdk | Cdcl4a | Cks2 | Zfp41 | Frrsl | Myll2a | Cnn2 | |
| Gprl7 | Tgfa | Fkbp7 | Cyp4v3 | Shmtl | Ndc80 | Ror2 | |
| Tnfrsfla | Tnr | Pmp22 | Mtssll | Plscrl | mt_AK140174 | Rsul | |
| Ptpizl | Phxr4 | Cdca3 | Slc22a6 | Car8 | AI8545 17 | 1700018G05Rik | |
| Cdc25c | Pllp | Frk | Derl3 | Srebfl | Matn4 | Rab3 1 | |
| Pcdhl5 | Arhgap3 1 | Kcnj 16 | Limal | Plekha2 | Foxcl | Dynltlc | |
| Ckap21 | Kcnh8 | Ltbpl | Ecil | Txlna | Vcaml | Sfmbt2 | |
| Pdgfrl | Tbxl8 | Cdol | Selenbpl | Epasl | Cpa4 | Nkiras2 | |
| Lhfpl3 | Seφ iηe2 | | Stk32a | 4933406J10Rik | Mdfic | Wnt7a | |
| Ogn | | | | | Cspg5 | Mpzl2 | |
| Itih2 | | | | | | | |

| Astrocytes | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gjal | Gramd3 | Slc7al 1 | Btd | Zfyve21 | Aldh6al | Alpl | Neu4 |
| Gjb6 | Slc7al0 | Phkal | Gpldl | Lgr4 | Pou3f4 | Gludl | Ugtla2 |
| CldnlO | 3 110082J24Rik | Id4 | Ccdcl41 | Tmeml76a | Clmn | Tsc22d3 | BC013529 |
| F3 | Hsd3b7 | Agmo | cx_tRNA-Ala-GCG | Sycp2 | Timp3 | Ccbl2 | Zfp783 |
| Slcla3 | Mtl | Fermt2 | Tomlll | Cptla | Slc6a20a | Tnfaip8 | Fjxl |
| Slc39al2 | Bean | Crot | Scrgl | Mettll 1b | Mif4gd | Zfp438 | Rasl2-9-ps |
| Sdc4 | Appl2 | Elovl2 | Smpd2 | Loxl3 | Plscr2 | Hesl | Suclg2 |
| Acsbgl | Chi311 | FkbplO | Bdh2 | Abhd4 | Pnp | A130022J15Rik | GdflO |
| Mfge8 | Adhfel | MegflO | Elovl5 | Papss2 | Btbdl7 | Slcl3a3 | Atp6v0e |
| Ntsr2 | Pxmp2 | AA387883 | Cd38 | Pdgfrl | Pdk4 | Cklf | Csgalnact 1 |
| Lcat | Tlr3 | Oaf | Ttyhl | Retsat | Fzd2 | Egfr | 1700003M07Rik |
| Cml5 | Vcaml | 1118 | Ccdc90a | Tcf712 | Slc7a2 | Ghr | Pyroxd2 |
| Aqp4 | Ctso | Pmp22 | Crlf3 | Sema4b | Tubb2b | Slc25a35 | Efemp2 |
| Pla2g7 | Agxt211 | Fabp7 | Slc26a6 | Rnasel2 | Rapgef3 | Ephx2 | Afap 112 |
| Ppap2b | AI46413 1 | Faml63a | Lxn | Fgfrl | Prkdl | Rbpl | Dbi |
| Ppplr3c | Maob | Satl | Pcsk6 | Igf2 | Adora2b | Pdlim5 | Gml073 1 |
| Slprl | Rfx4 | Kirrel2 | Paqr8 | Nat2 | Aoxl | Cdc42epl | 1190005106Rik |
| Slc25al8 | Acat3 | Serhl | Luzp2 | Mirl l92 | Hist2h3cl | Qk | Abhdl4b |
| Plcd4 | Mmd2 | Gstkl | Egfl6 | Dcxr | Cyp7bl | Faφ 1 | |
| Chrdll | Ugtla6a | Zfp3612 | Fgd6 | Apln | Arsk | 2210417K05Rik | |
| Faml07a | Gdpd2 | Arhgef26 | Hgf | Nrarp | Dhrsl 1 | | |
| Dio2 | | | | | S100al3 | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gpr3711 | Bmprlb | Slc4a4 | Cibl | S100a4 | Histlh2bq | Arapl | Trip6 |
| Mt2 | Prelp | Cyp4fl3 | Hspb8 | Sfxn5 | Histlh2br | Calml4 | Lama2 |
| Entpd2 | Pon2 | Emp2 | Acssl | Dok7 | Gng5 | Chst2 | Gml7660 |
| Gstml | Tril | Gm973 | Acsl6 | Plscrl | Acsl3 | Emx2 | Rin2 |
| Cbs | Gpc5 | Agt | Pion | Den | Sultlal | Slc22a6 | Fndc4 |
| Tst | Nat8 | Lixl | Notch2 | Ddo | Maml2 | Parp3 | Slc30al0 |
| Prodh | C030037D09Rik | Uppl | Ppil6 | 1810014BOlRik | Echdc2 | Gml0052 | Scg3 |
| Slcolcl | Cyp4fl4 | Naaa | Tcn2 | Nwdl | Tmem229a | Cede 18 | Abcd4 |
| Gfap | Nkain4 | Nfe212 | Renbp | Ugp2 | c2_tRNA-Ala-GCG | Tifa | C230035I16Rik |
| Tlcdl | Gml l627 | Steap3 | Pax6 | Myo6 | Notchl | Triml2a | Ptplad2 |
| Mlcl | Slc27al | Ptpizl | Cyr61 | Gpt | Slcl2a4 | Serpine2 | Rasa2 |
| Apoe | Natl | Cd63 | Gpam | Cst3 | Agpat5 | Mro | Acadl |
| C030018K13Rik | Mertk | Cmtm5 | Klfl5 | 01fr287 | Rlbpl | Vcl | Lrrc9 |
| Slc38a3 | Fmol | Gabrgl | Swap70 | Kctdl4 | LOC433374 | Per3 | 1700040N02Rik |
| Aldoc | 2900052NOlRik | Phkgl | Slc6al 1 | Zbtb20 | Kctdl2b | Taf4b | Zfp521 |
| Timp4 | Cth | Gasl | Lgals4 | Ddhdl | Ecil | I113ral | Prkcd |
| Cyp2d22 | TmemlO0 | Selenbpl | Psd2 | ZnrG | Tex 11 | 1190002H23Rik | Ranbp31 |
| Slcl5a2 | Cideb | Gpx8 | Pnpla7 | Olfmll | Lmcdl | Gypc | Npcl |
| Htral | Cmll | Soatl | Sall3 | Rmst | Cbr3 | Kcnj 13 | Hif3a |
| Atpl3a4 | Efempl | SlOOal | MyolO | Tmem5 1 | Zic5 | Gabrbl | **Pfkfbl** |
| Atpla2 | Mdk | Thrsp | Elmod3 | Hsdl lbl | Calr4 | Cmtm3 | Fcgr2b |
| Prdx6 | Kcnj 16 | A330048O09Rik | Histlh2bc | Rdh5 | Lhx2 | Itga7 | Rdml |
| 2010002N04Rik | Daam2 | Sc4mol | Smox | Eyal | Atplb2 | Angptll | Mmpl4 |
| Fgfr3 | Scara3 | Rfx2 | Ndel | Odf311 | Sox21 | Stkl7b | Grtpl |
| Pdpn | Mfsd2a | Phgdh | A330076C08Rik | Kankl | Gjb2 | Hacll | Wnt7b |
| Sox9 | 1700084COlRik | Hopx | 2610034M16Rik | Paqr6 | Dera | 01fr288 | Trp53bp2 |
| Fxydl | Rftn2 | Naprtl | Gml303 1 | Utpl4b | Hsdl2 | Faml81b | C2 |
| Itih3 | Prex2 | Ndrg2 | Enho | Histlh4h | Lpin3 | Ccdc77 | Lgals3bp |
| Faml76a | Dhrs3 | Acaa2 | Tnfsfl3 | Lpcat3 | Vgll4 | D630033O1lRik | |
| Cyp4fl5 | Grm3 | Slcla2 | Plxnbl | Aldhla2 | Zcchc24 | Phxr4 | |
| Gldc | 1700019G17Rik | B230209KOlRik | Cdkn2c | Lum | Slc22a4 | Nek3 | |
| Cml3 | Hepacam | S100al6 | Gem | A2m | Kcnj1O | 1700084J12Rik | |
| Ndp | Pgcp | Pbxipl | Tmeml76b | Rpe65 | Vav3 | Asrgll | |
| Cyp2j9 | Clu | Spatal7 | Nudt7 | Rcn3 | Gli3 | Gprc5d | |
| Slcl4al | Smpdl3a | Lpar4 | E030003E18Rik | Gnal3 | Akt2 | Decrl | |
| E1301 14 P18Rik | | Gpr56 | Cnn3 | Cyp2j6 | Eps8 | Lonrf3 | |
| | | Aass | | Fpgs | Nfia | | |
| | | | | Plodl | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Pdlim4 | Fam20a | Hadh | 4932438H23Rik | Fgfr2 | Tsc22d4 | Rnfl82 | |
| Aldhlll | Gm5083 | Acotl 1 | Lrp4 | Dockl | Lrrc5 1 | Mmgt2 | |
| Mgstl | Abhd3 | Pax6osl | Id3 | Frrsl | Grhll | Paqr7 | |
| Dbx2 | Ednrb | Ttpa | Aqp9 | Fads2 | Tnfrsfl9 | Haplnl | |
| Ezr | St3gal4 | Gstt3 | Histlh4i | Seppl | Adrbk2 | Cox6b2 | |
| Slc9a3rl | Rarres2 | Cdhl9 | Tdo2 | Trp63 | 2810055G20Rik | Sohlh2 | |
| | Glul | Nrlh3 | Gstm5 | | | Nphp3 | |
| | Faml98a | | Slcolb2 | | | Idh2 | |
| | Gm5089 | | | | | Btgl | |

| Cortical Neurons | | | | | | | |
|---|---|---|---|---|---|---|---|
| Nos1 | Scrt2 | Neurod2 | Serpinil | Nedd4l | Gstm7 | Elavl4 | Cdk2ap1 |
| Fam84a | Cdh4 | Srrm4 | Ttc28 | Fam114a2 | Emx1 | Scg5 | Cplx2 |
| Unc5d | Slc17a6 | Adgrl2 | Epha5 | Cux1 | Tmem108 | Scrnl | Efnb2 |
| Rnd2 | Osbpl6 | Jarid2 | Ankrd6 | Mta2 | Dbn1 | Ptprs | Klhdc2 |
| Pou3f2 | Sema3c | Pou3f3 | Tmeml58 | Acly | Mytll | Midn | Ccng2 |
| Pdzrn3 | Kif21b | Cttnbp2 | Plxna4 | Baz2b | Cul1 | Kdm2b | Parp6 |
| Hs3st1 | Wnt7b | X6330403K07Rik | Nfasc | Phf21b | H1f0 | Laptm4a | Nipsnap1 |
| Sstr2 | Tbr1 | Nav2 | F2r | Phip | Kif21a | Fam49a | Tax1bp3 |
| Pcp4 | Chga | Pantr1 | Fmnl2 | Tmeff1 | Ilf2 | Acinl | Ezr |
| Meis2 | Tenm4 | Lrpap1 | Cbfa2t2 | Ddah2 | Rpf1 | G3bp2 | Nol4 |
| Lrrc16b | Lmo1 | Trim2 | Lztsl | Grina | Ing4 | Mdk | Elavl2 |
| Plekhf2 | Tsc22d1 | Nek6 | Sorbs2 | Smim18 | Hist3h2a | Sbk1 | Arhgef2 |
| Sorl1 | Igfbpl1 | Ldhb | Frmd4a | Rbfox1 | Bcl7a | Auts2 | Nsg2 |
| Ppp2r2b | Nrn1 | Lhx2 | Plxna2 | Sncaip | Hivep3 | Kdm5b | Pbx1 |
| Trim9 | Wbscr17 | Tagln3 | Foxgl | Lrp8 | Hbb.bs | Ap3s1 | 43346 |
| Pou3f1 | Itpk1 | Mn1 | Cdknlb | Avl9 | Gdaplll | Basp1 | Zfp462 |
| Frmd4b | Sox5 | Vopp1 | Luzp2 | Nfix | Fam107b | Tmem57 | |
| Mllt3 | Prex1 | Gm17750 | Dpyl911 | Tnrc18 | Podxl2 | Peli1 | |
| Plcb1 | Rcor2 | Nfib | Rbfox3 | Znrf2 | Setbp1 | Cux2 | |
| Ppp2r1b | Kctd4 | Neurod6 | Cd24a | Adgrg1 | Wbp1 | Ttc9b | |
| Lsamp | Cited2 | Rasgef1b | Cdldl | Abracl | Ip6k2 | Rundc3a | |
| Enc1 | Epha3 | Hs6st2 | Cyth2 | Mpped1 | Igsf3 | Mpped2 | |
| Robo2 | Palmd | Insm1 | Negri | Gria2 | Gm14964 | Mkrn1 | |
| | Bcar1 | Tmem178 | Hist3h2ba | Zbtb18 | Nrp1 | Akap9 | |

| RadialGlia-Id3 | | | | | | | |
|---|---|---|---|---|---|---|---|
| Id3 | Heyl | Efcabl | Add3 | Morn2 | Slc25a25 | Pex7 | X2810417 H13Rik |
| Idl | Aldoc | Nes | Lrp4 | Nafl | Pmp22 | Galkl | Extl |
| Foxj l | Anxa2 | Mest | Ifitm3 | Cripl | B9dl | Hsdl7b7 | Tancl |
| Mtl | Atplb2 | Slc6al l | Tspanl 5 | GrblO | Purb | Anxa5 | Lhfp |
| Mt2 | Ncan | Glul | Slc27al | Itm2c | Ctso | Ift22 | Amot |
| Pla2g7 | Atpla2 | Faml81b | Gludl | Sparc | Axl | Sgcb | F3 |
| Hes5 | Cybrdl | Camk2d | Timp3 | Mmd2 | Dhcr24 | 43358 | Pmfl |
| Hesl | Tmeml07 | Zfp3612 | Hopx | Mcm3 | Tppl | Tmem218 | Stat3 |
| Mia | Lgalsl | Gjal | Cav2 | Acyp2 | Stxbp6 | Slcla2 | Ppplrla |
| Egrl | Slcl4a2 | X2810459 M l lRik | Arl4a | Adcyaplrl | Rasa3 | Rbpl | Gprc5b |
| Metrn | Rhoq | Spry2 | Chptl | S100al3 | Cbfb | Arhgef26 | Dhfr |
| Fos | Tlcdl | Vim | Fhll | Eif4ebpl | Pacsin2 | Dnajcl5 | Lyrm5 |
| Tmem47 | Rhoc | Acadl | Tst | Irsl | Gcsh | Pmml | Cdk2 |
| Ednrb | Sox9 | Igfbp2 | Plpp3 | Cibl | Parva | Cfap36 | Nfkbia |
| Tppp3 | Ccndl | Ckb | Spal7 | Afapll2 | Zebl | Etfa | Cntln |
| Clu | X1500015 OlORik | Paqr8 | Tomlll | Ttyh3 | Nkain4 | Pidl | Gasl |
| Serpine 2 | Bhlhe40 | Gng5 | 43352 | Notch2 | Snx5 | Ctdspl | Pfnl |
| Riiad l | Zfp3611 | Hspa2 | Msn | S100a6 | Ormdl2 | Ecil | Prdxl |
| Gfap | Ddit41 | Lrigl | Pttgl | X2610301 B20Rik | Adgrvl | Plxnbl | Golph3 |
| Sparcll | Nimlk | Erf | Ninj l | Magtl | Stard4 | Klf6 | Cy stall |
| Apoe | Nme5 | Zic5 | Fkbp9 | Itgb5 | Car2 | X I500009 L16Rik | Kcnip3 |
| Slcla3 | Lfng | X1810037I 17Rik | Ctsc | Kbtbdl l | Sox21 | Emc7 | Prdx4 |
| Nlrxl | Tagln2 | Bcl2 | Rrbpl | SlOOal | Slprl | Dennd2a | Rad23a |
| Selm | Mfge8 | Ier2 | Prkcdbp | Mif4gd | Slcl2a4 | Zdhhc21 | Traml |
| Ttyhl | Stom | Vcaml | Gnai2 | Tnfaip8 | Hacdl | Plcel | Dclkl |
| Gstml | Pbxipl | Ptn | Nr3cl | Pcx | Cd9 | Oat | Hspa5 |
| Lxn | Empl | Nkdl | Ldha | Dnajc3 | Wwpl | MyolO | Gm2a |
| Cyr61 | Mpp6 | Trim47 | Slc38a3 | Dagl | Jun | Phyhipl | Smo |
| Fbxo2 | Pdpn | Ptprzl | Zcchc24 | Rgs20 | K1M13 | Maml2 | Spcs3 |
| Mlcl | S100al6 | Krccl | ZnrG | Tapbp | Gabrbl | Irs2 | AI8545 17 |
| Enkur | Tspan33 | Scd2 | Akrlbl0 | Hmgcsl | Msi2 | Msmol | Flna |
| Mlfl | Aldhlll | Tnfrsfl9 | Hadh | Nudt4 | B2301 18 H07Rik | Mras | Csrpl |
| Mgstl | Fam212b | Zfp36 | Myo6 | Mlec | Eef2kmt | Mtssll | Gpt2 |
| Slc9a3r1 | Fzd9 | Idil | Kcnj lO | Degsl | Nr2c2ap | Asrgll | Ift74 |
| Bean | Pdlim5 | Serpinhl | Acadm | Abhd4 | Dpcd | Faml95a | Sytl 1 |
| | Eepdl | Ntrk2 | | Sp3os | I16st | Socs2 | Clicl |
| | Ier3 | | | Sashl | Rgcc | Fadsl | 1118 |

| Fabp7 | Fbln2 | Suclg2 | Psph | Fjxl | Rnftl | Trip6 | Myll2a |
|---|---|---|---|---|---|---|---|
| Dbi | Junb | Metrnl | Psatl | Uhrfl | Rasll 1a | Rexo2 | Scrgl |
| Emp2 | Peal5a | Rgma | Prrxl | Slcl5a2 | Ak3 | Ptgfrn | Nphpl |
| Ppp1r3c | Kcnell | Rcnl | Tns3 | Cenpw | Echdcl | Sri | Proml |
| Igfbp5 | Etv4 | Axin2 | Slc39al | X I 110004 E09Rik | Nr2f6 | Nfe212 | Ctnnal |
| Wis | Rampl | Klf9 | Itgav | Cebpb | Vamp3 | X23 10022 B05Rik | Pde4b |
| Tpbg | Sfxn5 | Klfl5 | Gm5617 | Tspanl2 | Arhgef40 | Snx3 | Ligl |
| Fgfr3 | Egfr | Npas3 | Ccpglos | Tribl | Ifngrl | Thbs3 | Itgb8 |
| Hepacam | Klf4 | Satl | Notchl | Pcgf5 | Phxr4 | PcdhlO | Sox8 |
| Aqp4 | Gpx8 | Chst2 | Prrl8 | Pnp | Tm7sf2 | Elofl | |
| Oligl | Cpne2 | Paqr4 | Cbs | Faml20a | Mvk | Tctexld2 | |
| Tnc | ChchdlO | Cd63 | Rest | Gmnn | Dnajc24 | Fgfr2 | |
| Mt3 | Ndrg2 | Spryl | Anxa6 | Polr3h | Hsdl2 | 43345 | |
| Slc4a4 | Rmst | Dkk3 | Insigl | Creb5 | Bola3 | Betl | |
| Gngl2 | Nebl | Bmprla | Nrarp | Pygb | Wwtrl | Spsb4 | |
| Pacrg | Jam2 | Epdrl | Emc2 | Trim9 | **Traf3** | Lss | |
| Rspo3 | Acsbgl | Yapl | Thrsp | Ppargcla | Spata24 | Phlda3 | |
| Phgdh | Pon2 | Adamtsl | Efemp2 | Grm5 | Bakl | E2f5 | |
| Tril | Fosb | Mnsl | Acotl | Rab3 1 | Tspan7 | Nrcam | |
| Qk | Smpdl3a | Aldoa | Bphl | Grhpr | Lppos | Ddahl | |
| Ccdc80 | Fatl | Ccnd2 | Nr4al | Btg2 | Nab2 | Klhdc8b | |
| Aard | Sema6a | Slcla4 | Ppic | Gale | Mcee | Plin3 | |
| Plat | Gdpd2 | Nog | Cxxc5 | Tjpl | Chsyl | KlflO | |
| 01ig2 | Tsc22d4 | SlOOal 1 | Ill lral | Cnp | Dusp6 | Klf3 | |
| Rfx4 | Sall3 | Itga6 | Gins2 | Donson | Midlipl | Gltp | |
| Cmtm5 | Gsta4 | Fgfbp3 | Rorb | Cst3 | Cetn2 | Ccdc8 | |
| Id4 | Cspg5 | Duspl | Sox2 | Hspa41 | Dtd2 | Speccl | |
| Socs3 | Neatl | X 3 110082J 24Rik | Rabl3 | Cln5 | Trpsl | X4933434 E20Rik | |
| Scdl | | X1700088 E04Rik | Nacc2 | | | | |
| | | | Ung | | | | |

| RadialGlia-GdflO | | | | | | | |
|---|---|---|---|---|---|---|---|
| GdflO | Assl | Pdpn | Arhgef26 | Gmnn | Ligl | Rfcl | Msi2 |
| Id3 | Htral | Dkk3 | Rcnl | Pdcd4 | Prps2 | Glol | Tyms |
| Tesc | X2810459 M 1 lRik | Col9a3 | Noval | Cdl64 | Gstm5 | Tpx2 | Spg20 |

174

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Thrsp | Bcl2112 | Mgstl | Appl2 | Maml2 | Naa50 | Atxn7 | Fut9 |
| Tnfrsfl9 | Gjal | Lrp4 | Mki67 | Scrgl | Sypl | Cenpw | Proxl |
| Frzb | E1301 14P1 8Rik | Foxol | Phxr4 | Kcnmb4 | Krccl | Ddahl | Pmp22 |
| Idl | Nkdl | Dmd | Anxa6 | Ccna2 | Eci2 | Proxlos | Ccdc34 |
| Sdpr | Ninj l | Entpd2 | Nr2f6 | Kbtbdl 1 | Jam2 | Torlb | Sntal |
| Emidl | Enpp2 | Dmrt3 | Gli3 | Lap3 | Cisd3 | Asahl | Cdv3 |
| E330013P04Rik | Fzdl | Chst2 | Tgifl | Knstrn | Fezf2 | Ndufc2 | Tmem256 |
| Hspb8 | Selm | Gpx8 | Pygb | Gng5 | Lhfpl2 | Bmprla | Ssl8 |
| Pdlim3 | Hadh | Tsc22d4 | Tspanl5 | Chptl | Mcm5 | Crip2 | Aamdc |
| Den | Psph | Isocl | Sdc2 | Snx5 | Nadk | Cpne3 | 43345 |
| Gfap | Sfxn5 | FkbplO | Tspanl2 | 4335 1 | Tjpl | Lysmd2 | Sox6 |
| X15000150lORik | Aard | X 1110015 018Rik | Fatl | Slit2 | Cxxc5 | Sat2 | Arhgap5 |
| Mt2 | Lrrcl | Gngl2 | Zfp3612 | Itgb8 | Proml | Abhd4 | Paics |
| Lefl | Dbi | Epdrl | Hells | Mcm3 | Pacsin 3 | Faml20a | Snap23 |
| Rmst | Frasl | Cpne2 | Hmgb2 | Prdx4 | Pankl | Rcn3 | Scd2 |
| Gasl | Slc9a3rl | Ptgfrn | Cdca8 | Litaf | Dennd 2a | Ckslb | Ctdspl |
| Tst | Ltbpl | Mt3 | Cst3 | Ctdsp2 | Rdml | Kpna2 | Gsr |
| Mgll | Dmrta2os | Zicl | Aifll | Kcnipl | Uspl | Evi5 | Fkbp9 |
| Zic5 | Notchl | Lmcdl | Itga6 | Hnll | Cmc2 | Pmfl | X493343lE20Rik |
| Sp5 | Lhfp | Notch2 | Lockd | Gcsh | Nit2 | Dpysl4 | Atplbl |
| Hopx | Emx2 | Id4 | Gstml | Hs2stl | Adgrb 1 | Ifitm2 | Exosc5 |
| Prex2 | Bcl2 | Msn | Acotl | Cdkl | Nme4 | Bach2 | Mettll |
| Eyal | Axin2 | Mlcl | Ube2c | Slcla4 | Echdc 1 | Slc35a4 | Atplal |
| X0610040JOlRik | Etv4 | Qk | Pttgl | Dhcr24 | Apoe | Kcnell | Syce2 |
| Cavl | Sez61 | Smco4 | Lixl | Arl4a | Mcm6 | Cdol | Ost4 |
| Mtl | Efcabl | Eepdl | Btg3 | Dhfr | Smc2 | Sival | Actnl |
| Adamtsl9 | Fos | Myl9 | Otxl | Shisa4 | Dclkl | Pcna | Rangrf |
| Wnt8b | Mro | Cdkn2c | Cbfb | Tmeml07 | Dtymk | Efemp2 | Hmgn3 |
| Nme7 | Tnc | Tspan7 | Pnp | Pcx | Jam3 | Cntln | Nrarp |
| Cripl | Rhoc | Cd9 | Tgif2 | Ldha | Pax6 | X23 10022B05Rik | Carnmtl |
| Zfp3611 | Rfx4 | Gabra4 | Cks2 | Slc39al | Paqr4 | Acadm | Hmbs |
| Cyplbl | Rgma | Dtl | Pbk | Serpinhl | Stard4 | Ier2 | Rnftl |
| Lhx9 | GrblO | Gnai2 | Rpa2 | Tcfl9 | Elavil | Cdc42sel | Sytl l |
| Vim | Ung | Plpp3 | Limdl | Bola3 | Vcan | Adrbk2 | Fuz |
| Rgs20 | Atpla2 | Cenpf | Idil | Ndel | Histlh 1e | Mvk | Tspanl8 |
| Hes5 | St3gal4 | Klf9 | Cyba | E2f5 | | Rragd | Fam96a |
| | X2700046A07Rik | Faml67a | Top2a | Camk2d | | D8Ertd82e | |
| | | Gldc | Sesn3 | Cdk2 | | Nudt4 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Tpbg | Fbln2 | Paqr8 | Csrpl | Ccnb2 | Tulp3 | Csad | Dennd5 a |
| Slcla2 | Vephl | Rftn2 | Tancl | SlOOal l | Mcee | Purb | Nudcd2 |
| Aldoc | Tmeml32c | Stxbp6 | Erf | Tmem97 | Nudt5 | Rpl2211 | Dnphl |
| Slcla3 | Dmrta2 | X23 10009 B15Rik | Sox8 | Rabl lfip2 | Ptprg | Fjxl | Ybx3 |
| Psatl | Col2al | Gins2 | Tex9 | Eefld | Histlh 2ap | Mpp6 | Speccl |
| Ttyhl | Emp2 | Uhrfl | Map3kl | Mcm4 | Decrl | Bcl7c | Tpil |
| Hesl | Nimlk | Ephbl | Fignll | Suclg2 | Higdl a | Stx4a | Akr7a5 |
| Tspan33 | Loxll | Clu | Sirpa | Gem | Ift74 | Mgatl | |
| Cpne8 | Pbxipl | Lrrc4c | Spc24 | Ehbpl | Lsm2 | 43358 | |
| Hepacam | Mfge8 | Gsap | Dnajcl | Insigl | Ldlrad 3 | X2810004 N23Rik | |
| Sox9 | Rest | X2810417 H13Rik | Ephb3 | Pdk3 | Cachd 1 | X150001 1 K16Rik | |
| Vcaml | Trip6 | Cdca3 | Atplb2 | Amot | Ppplr 1a | Anp32b | |
| Ccndl | Gabrbl | Socs2 | Mif4gd | Smo | Histlh 4i | Rpal | |
| Tmem47 | Fgfr3 | Adcyaplrl | Heyl | A730017 C20Rik | Acadl | Spredl | |
| Gludl | Pon2 | Ptn | K1M5 | Vamp3 | Mcm2 | Hspa41 | |
| Snedl | Tns3 | Yapl | Birc5 | Ramp2 | Nacc2 | Crot | |
| Ccdc80 | Tgfb2 | Cbs | Sapcd2 | Arhgef40 | Prdxl | Tmeml67 | |
| Fbxo2 | Fam49b | Sparc | Tead2 | Epsl5 | Fxyd6 | Echdc2 | |
| Lfng | Prkcdbp | Cenpm | Ecil | Wwtrl | Nr2el | Caldl | |
| Tfap2c | Cspg5 | Cyr61 | Chd7 | Rnf26 | Itgb3b P | Lhx2 | |
| Ndrg2 | Zcchc24 | Prdx6 | Npas3 | Vgll4 | Ckap2 | Nek6 | |
| Cthrcl | Slc27al | Vatll | Cenpa | Rexo2 | Vldlr | Lyrm5 | |
| Cav2 | Sashl | Sox2 | Hrspl2 | Btgl | Tipin | Toporsos | |
| Mmd2 | Gas6 | Ttyh3 | Klf4 | Cdon | Homer 2 | Arl6 | |
| Phgdh | Adgrvl | | | | Kctdl 2 | | |
| | | | | | Dagl | | |
| | | | | | Rpe | | |

| RadialGlia-Neurog2 | | | | | | | |
|---|---|---|---|---|---|---|---|
| Neurog2 | Kif26b | Wasf2 | Dnajb2 | Echdcl | Asahl | Hyal2 | Ndufaf7 |
| Eomes | Tmem98 | Ecil | Asnsdl | Elavil | B230354 K17Rik | Nrnl | Gm8730 |
| Gadd45g | Fam53b | Mmpl4 | Zbed3 | Akr7a5 | Acadvl | Shmt2 | Dexi |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Rhbdl3 | Dhx32 | Ckb | Vps37b | Ift22 | Cnih4 | Zfp62 | Pno1 |
| Ptgds | Abcd2 | Gadd45gip1 | Fubp3 | Ctnnb1 | Yif1a | Svip | Gspt1 |
| Btbd17 | Lzts1 | Ddah1 | Dcaf8 | Azi2 | Ift52 | Ubxn2a | Fxn |
| Snhg18 | Dll3 | Glo1 | Tbrg1 | Ece2 | Srsf6 | Rad23a | Snhg6 |
| Lima1 | Aif1l | Ccs | Ufm1 | Pmepa1 | Hibadh | Golim4 | Ccdc86 |
| Tfap2c | Cbs | Ift74 | Wscd1 | Bphl | Foxp4 | Scrn1 | Bola3 |
| Mfng | X1500015O10Rik | Slc25a5 | Lta4h | Fundc2 | Gnpda2 | Vrk3 | Kti12 |
| Btg2 | Gpx8 | Sfxn5 | Idh2 | RP23.207N5.2 | Cpne3 | Urod | Pou2f1 |
| Myo10 | Cmc1 | B230118H07Rik | Gstm5 | Paics | Lamp2 | Taf10 | Mrpl24 |
| Csrp1 | Slc1a2 | Pam | Sema5b | Rbpj | Itgb3bp | Pdcd4 | Rit1 |
| Tead2 | Bcl2l12 | Lzts2 | Hadh | Rangrf | Rcor2 | Rbfox3 | Lztfl1 |
| Pax6 | Rnaseh2b | Hmgn2 | Ftsj3 | Rpl22l1 | Cplx2 | Mphosph10 | X1810058I24Rik |
| Celsr1 | Mcm2 | Ddr1 | Pyurf | Ptbp1 | Cadm3 | Emg1 | Swt1 |
| Gm29260 | Ezr | Ninj1 | Eci2 | Nedd4 | Ankrd6 | Smarcad1 | Eif3i |
| Chd7 | Gng5 | Srek1ip1 | Paqr8 | Aco1 | Myl12a | Rrp15 | Spata2 |
| Acads | Tank | Adk | Fam96a | Flna | Lman2 | Ldha | Tef |
| Heg1 | Apool | Snx5 | Atf5 | Nkain4 | Cnpy2 | Ppib | Vamp3 |
| Dll1 | Spsb4 | Acot1 | Rps18.ps3 | Rprm | Mrpl17 | Cdk4 | Ift43 |
| Gamt | Hrsp12 | Zfand1 | Cdca7 | AI854517 | Trp53 | X1500011K16Rik | Guf1 |
| Kcne1l | Cd63 | X2610301B20Rik | Rexo2 | Polr3k | Mrps14 | Tmed1 | Gm10020 |
| Tox3 | Ccdc136 | Serpinh1 | X2810004N23Rik | Hsd17b4 | Fars2 | Cdk5rap3 | X2310011J03Rik |
| Rcn1 | Ddit4 | Cib1 | Prdx1 | Trap1 | Serinc2 | Acly | Setbp1 |
| Gfap | Grb10 | Fbln1 | Efs | Mcee | Prdx3 | Lyrm4 | Rnf13 |
| Igfbp5 | Pttg1 | Syne2 | Golph3l | Npc2 | Fam162a | Slc48a1 | Mccc1 |
| Hes6 | Nr2e1 | Nrg1 | Echs1 | D10Jhu81e | Atp5g2 | Mt2 | Akr1b3 |
| Efhd2 | Tmem218 | Ncald | Ormdl2 | Mettl1 | Sp3os | X1110012L19Rik | Hspe1 |
| Inpp1l | Btg3 | Elavl2 | Exosc3 | Dazap2 | Mettl5 | Fam174b | Ralgds |
| Lrrn3 | Zeb1 | Phgdh | Ccdc58 | Ino80b | Clic4 | X1810037I17Rik | Hmgn5 |
| Sfrp1 | Eef1d | Ly6e | Anp32b | Rbbp9 | Twf1 | Hnrnpf | Immp1l |
| Nme4 | Sstr2 | Insm1 | Cul1 | Prdx6 | Lap3 | Tpm4 | Carnmt1 |
| Sox21 | Thrsp | Abca1 | Sox6 | Elp4 | Creb5 | Mt1 | Iscu |
| Loxl1 | Sema5a | Slc1a3 | Hdac1 | H1f0 | Emx1 | Acvr2b | Isca2 |
| Fam210b | Gas1 | Ttc8 | Tmem33 | Exosc5 | Rrs1 | Gcsh | Tspan3 |
| Dbi | Slco1c1 | Phyh | Limd1 | Sipa1l1 | Cdkn2c | Ift57 | Gkap1 |
| Tgif2 | Rcn3 | Ccdc167 | Tor1aip1 | Sesn1 | Rps27l | X2310039H08Rik | Actl6a |
| Ccnd2 | Ctnna1 | Dnajc15 | Por | Gm14305 | Ebpl | | Pdia6 |
| Vim | F2r | | Adcyap1r1 | Pbdc1 | Timm21 | | Ppie |
| Mfap4 | Zfp703 | | | | Nsmce4a | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Mdk | Mdgal | Lyrm5 | Cyba | Wdr61 | Dhx40 | Rpe | Sod2 |
| Notchl | Inhbb | Smpd2 | Hadha | Adgra3 | Mmd2 | Zbtb38 | Odcl |
| Gem | Pnpla2 | Litaf | Teadl | Pabpcl | Rhoc | Crnkll | Fucal |
| Magil | Zfp3611 | Nudt5 | Calu | Llgll | Ppp2r3d | Aamdc | Polr3c |
| Corolc | Sufu | Krccl | Ndufc2 | Clicl | Spire 1 | Gnpat | Med9 |
| Mfap2 | Smco4 | Scp2 | Etfa | X2210016 F16Rik | H2afv | Pfkl | Pex2 |
| E1301 14 P 18Rik | Rab8b | Ube2g2 | Dync21il | Draxin | Mφ 154 | Gml0073 | |
| Dleu7 | Dmrta2 | Betl | TmedlO | Ginml | Tlel | Mybbpla | |
| Ascll | Ndrg2 | Trappc6a | Snapin | Ddx52 | Tpcnl | Capn2 | |
| Igdcc4 | Cdk2apl | Tsc22d4 | Lrp8 | Msi2 | Igbpl | Eiflb | |
| Tmeml3 2b | Ehbpl | Actr3b | Hdhd2 | Zfp219 | Dszf5 | Ntrk2 | |
| Myo6 | Echdc2 | Dnajc24 | Cdk6 | Ppp2r3c | Sec23b | Pgaml | |
| Uaca | Egrl | Sdc3 | Ssl8 | Rcn2 | Chracl | Josd2 | |
| Slc30al0 | Hs3stl | Sox2 | Ctage5 | Arl6ip6 | Smim20 | Tφ c4aρ | |
| Gml 162 7 | Msn | Fezf2 | Pcbd2 | Tmed4 | Gpil | Ctsz | |
| Pdlim4 | Hmg20b | Gtf3c6 | Fam58b | Stx4a | Pts | Ubxn4 | |
| Zhx2 | Cbfa2t2 | Emidl | Qars | Klf3 | Plagll | Lengl | |
| Jam3 | Rgs3 | Pcmtd2 | Tfdp2 | Ivd | Rcbtb2 | Tmem230 | |
| Zfp423 | Elavl4 | Aldh6al | Aldh7al | Fgd4 | MφIIO | Tmeml78 | |
| Cdl64 | Aldh2 | Prmt8 | Kat6b | Bbx | Pgap2 | Sat2 | |
| Pgpepl | Chn2 | Smiml 1 | Nit2 | Ssbpl | Zmizl | Cd320 | |
| Dhrs4 | Rabl3 | Kdm7a | Tcf3 | Hadhb | Slc35b2 | Dermd5a | |
| Igsf8 | Fdxl | Qsoxl | Adgrgl | X2810006 K23Rik | Morn2 | Ost4 | |
| | Mfge8 | Nrarp | Acadm | Bckdha | Zfp664 | Nabp2 | |
| | | Pex7 | Glrx2 | Efnbl | | Nudcd2 | |
| | | B9dl | | | | Faml20a | |
| | | | | | | Mrfapl | |

| Long-term MEFs | | | | | | | |
|---|---|---|---|---|---|---|---|
| Rps3a3 | Ckslb | Utfl | Crabpl | Nop 16 | Manf | Rplp 1 | Cox6al |
| Timpl | Pinl | Trappc4 | Pfdnl | Tacc3 | Psmc2 | Srsf3 | Ppmlg |
| Bexl | Ccngl | Vdac2 | Atp5b | Ncl | Dnlz | Psma 5 | Nosip |
| Rhox5 | Tpil | Mφ s6 | Hspa9 | Naca | Rps25 | Polr2 e | Olal |
| Gml5459 | Eif4ebpl | Gml0039 | Nedd8 | Hintl | Pdrgl | Eif31 | Gtf2f2 |
| S100a6 | Tubb6 | Srupe | Ube2a | Rcn2 | Steapl | Srupa | Hprt |
| Gml032 | Txnl4a | Ruvbl2 | Nsmce 1 | Pgd | Snx5 | | Sec 13 |
| | Cdkn2a | Txnrdl | Rpl23a | MφII I | Rtn4 | | Ndufs6 |

| 0 | Npm1 | Actb | Psmd12 | Rps17 | Csnk2b | Rps4x | Eif3g |
|---|---|---|---|---|---|---|---|
| Gsto1 | Cenpa | Snrpa1 | Dynll1 | Ftl1 | Nab2 | Farsa | Brix1 |
| Gm11942 | Tagln | Mrto4 | Rps20 | Strap | Hcfc1r1 | Rpl17 | Timm10 |
| S100a4 | Lgals1 | Abracl | Rhoc | Atp5f1 | Eif1a | Mrps15 | Mrps14 |
| Gm10260 | Tmsb4x | Pgk1 | Pdlim1 | Idh3a | Cap1 | Cisd1 | Sf3b4 |
| Mif | Hmgn1 | Ngf | Cct5 | Ctxn1 | Fhl2 | Eif2s2 | Prps1 |
| Esd | Atp5g3 | Cct3 | Phf5a | Avpi1 | Pam16 | Arpc5 | Emc8 |
| Gm15772 | Acot7 | Hbegf | Glrx3 | Rps8 | Psmb5 | Mrpl42 | Ndufs4 |
| Anxa1 | Ranbp1 | Rack1 | Sh3bgrl3 | Stip1 | Chchd1 | Noct | Uba3 |
| Ctgf | Plaur | S100a11 | Pomp | Cdca8 | Dtymk | Txndc9 | Srm |
| Rps271 | Vim | Eno1b | Nudcd2 | Mdm2 | Bud31 | Mrpl35 | Gtf2h5 |
| Pkm | Cnih4 | Cox5a | Apoc1 | Eif2b3 | Rassf1 | Nt5c | Mrpl17 |
| Bex3 | Anxa3 | Timm17a | Nmd3 | Arl6ip1 | Rbm8a | Snrpg | Selenof |
| Txn1 | Tnfrsf11b | Eloc | Rpl19 | Rps3 | Snu13 | Eif3i | Praf2 |
| Tagln2 | Dctpp1 | Mtch2 | Cacybp | Capg | Snrpd2 | Rpl7l1 | Med7 |
| Tnfrsf12a | Cnn2 | Fkbp3 | Ddx39 | Hspe1 | Mthfd2 | Tgif1 | Tuba1a |
| Ldha | Eif5a | 2810025M15Rik | Hnrnpc | Edf1 | Gins2 | Rab11a | Tspan4 |
| Selenoh | Ass1 | Slc25a3 | Spp1 | Calr | Hsd17b12 | Nip7 | Degs1 |
| Serpinb2 | Krt18 | Rps13 | Cstb | Spc24 | Rplp0 | Plp2 | Rps26 |
| Gm28438 | Cdc20 | Rpl7a | Cox7b | Rps24-ps3 | Bzw1 | Vps29 | Ppil3 |
| Tex19.1 | Psma6 | Gm11273 | Tes | Prdx2 | Psmd13 | Dph3 | Dnaja2 |
| Gm10263 | Ccnb2 | Pa2g4 | Lxn | Shmt2 | Denr | Ndufb6 | Itgb1bp1 |
| Tubb5 | Prelid1 | Thyn1 | Nasp | 2810004N23Rik | Atpif1 | Lap3 | Cldn4 |
| Birc5 | AA465934 | Cdk4 | Atp5o | Lamtor1 | Cox7a2 | Naa38 | Commd2 |
| Ran | Cct8 | Eif1ax | Rpl39 | 2010107E04Rik | Ptrh2 | Zyx | Nol7 |
| Anxa2 | Ppia | Serpine1 | Eif4a3 | Yrdc | Mybbp1a | Sae1 | Cops5 |
| Gsta4 | Bola2 | Psma1 | Gars | Commd3 | Nsun2 | Rpl30 | Txndc17 |
| Nme1 | Eef1b2 | Cct7 | Gjb3 | Pebp1 | Mrpl30 | Tpm2 | Txn2 |
| Trap1a | Dut | Btf3 | Mrpl20 | Ccna2 | Aimp1 | | Prdx4 |
| Rrm2 | Ap1s1 | Hspd1 | Elob | Perp | Emc6 | | Wdr12 |
| Prdx1 | Rpsa-ps10 | Gng2 | Ptgr1 | Tmem126a | Arpp19 | | Prdx5 |
| Il11 | Psma2 | Mtpn | Acta2 | Rps5 | Snx3 | | Vta1 |
| Tm4sf1 | Cct4 | Tomm40 | Eif3d | Fcf1 | Coq7 | | Alad |
| Tuba1c | Hmga2 | Ccnb1 | Bdnf | Atp6v1g1 | Tmco1 | | Imp4 |
| Tuba1b | Psmd8 | Slc25a5 | Cops6 | Dars | Rars | | Exosc8 |
| | Pclaf | Psmb3 | Pno1 | Lsm5 | Phb2 | | Mrpl39 |
| | Snrpd1 | Tyms | Fam162a | Tpm4 | 1810022K09Rik | | Rpl22 |
| | | Rpl13a | Hnrnpab | | Apex1 | | Nras |

| | | | | | | |
|---|---|---|---|---|---|---|
| Enol | Bax | Tbca | Mɸ 113 | Cct6a | Tpml | Uqcrb |
| Cks2 | Rpl27 | Sgkl | Rpsl 2 | Rpl34-psl | Rslldl | Cede 58 |
| Psatl | Inhba | Aldoa | Rpll l | Mɸ 128 | Rɸ 9 | Rpl6 |
| Ube2c | Psph | Mtap | Fkbpla | Ssscal | Psmb6 | Gpxl |
| Cldn3 | Gml673 | Actgl | Eefld | Hspbl | Bag2 | Ppplr11 |
| Fabp3 | Naplll | Rps41 | Rplp2 | Rgsl6 | Psmcl | Thoc 7 |
| Hatl | Pttgl | Gmnn | Nme4 | Rpl9 | Nup35 | Cdc37 |
| Mrpll2 | Eeflel | Prdx6 | Aurka | Paics | Psmbl | Polr2f |
| Eif2sl | Sɸ 14 | Med21 | Aaas | Ciapinl | Prss23 | Nradd |
| Cfl1 | Psmdl4 | Dnphl | Fosll | Mɸ 151 | Ndufa8 | Aɸς 2 |
| Myll2a | Bri3bp | Pfdn4 | Ndufb8 | Elofl | Akl | Mɸ ï 57 |
| Tubb4b | Asns | 1110008F13Rik | Lsm8 | Mɸ s18a | Bcap3 1 | Gnl3 |
| Clicl | RpslO | Lsm2 | Timm50 | Tcpl | Sigmarl | Vbpl |
| Cdkl | Clqbp | Pfnl | Hnl | Tkl | Ak6 | Pmm 1 |
| Aprt | Cnihl | Slcl6a3 | 2200002D0IRik | Phlda3 | 1500009L16Rik | Rpsl 5a |
| Gm4366 | Rpll2 | Psmc6 | Serbpl | Zwint | Tipin | Mob 4 |
| Hmgal | Nhp2 | Capzb | Ankrdl | Rheb | Slirp | Atxn 10 |
| Vmpl | Cct2 | Txnll | Rbxl | Chmp6 | Snx7 | Usp3 9 |
| Crlfl | Cdkn2b | Uqcrq | Itga5 | Ndufa7 | Pmfl | Zfp5 93 |
| Gapdh | Rpl2211 | | | Cox6bl | | Hikeshi |
| Banfl | | | | | | Tars |
| Rpll8 | | | | | | Rpl2 8 |
| Galkl | | | | | | Erh |
| | | | | | | Rpsl 5 |
| | | | | | | Phgdh |
| | | | | | | Krt8 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| | | | | | | Cox17 | |
| | | | | | | Fez2 | |
| | | | | | | Tbpl1 | |
| | | | | | | Arhgdia | |
| | | | | | | Ddal | |

| Embryonic mesenchyme | | | | | | | |
|---|---|---|---|---|---|---|---|
| Matn4 | S100b | Hmgn1 | Pdap1 | Prelid1 | Bub3 | Peg3 | Rpl31 |
| Matn1 | Crabp1 | 1110004F10Rik | Sdhaf2 | 2210013O21Rik | Psmb6 | Atp5g1 | Rps11 |
| Col9a1 | Fibin | Gm1673 | Hpf1 | Serf1 | Thoc3 | Slc25a4 | mt-Nd1 |
| Col9a3 | Siva1 | Psmd6 | Rer1 | Pdxdc1 | 2310036O22Rik | Nop58 | Rpl10 |
| Cnmd | Gpc3 | Ssr2 | Tmed1 | Srsf3 | Rpl36al | Chchd2 | Rps5 |
| Asb4 | Cthrc1 | Sub1 | Mif | Gnl3 | Limd2 | Arf1 | Rpl26 |
| Col9a2 | Tpi1 | H19 | Hnrnpm | Ndufa4 | Hnrnpa2b1 | Ier3ip1 | Rps8 |
| Wwp2 | Hnrnpd | Grb10 | Gars | Meg3 | Snx17 | Rps27a | Rps15a |
| Sox9 | Col11a1 | Prpf19 | Capn6 | Fkbp4 | Elp2 | Calr | Rplp0 |
| Col2a1 | Cpe | Elovl6 | Fus | Rcn1 | Atp5a1 | Swi5 | Rpl13 |
| Nnat | Fgfr3 | Dek | Psma7 | Itm2a | Slirp | Rps9 | Rps25 |
| Hapln1 | Eno1 | Pkm | Gstm5 | Hsp90b1 | Atp5k | Cox5a | Rpl18a |
| Cytl1 | Ccnd1 | Snrpd3 | Fkbp11 | Ugdh | Blmh | Rpl18 | Rps14 |
| Cd24a | Rflna | Ptov1 | Skp1a | Ddx39b | Nasp | Ndrg2 | Dlk1 |
| Mest | Rangap1 | Psmc4 | Apex1 | Hspe1 | Hint1 | Usmg5 | Rpl41 |
| Mia | Maged2 | Nop10 | Papss1 | Sec61b | Ddx39 | Rps2 | |
| Bex2 | Mlf2 | Tial1 | Cct3 | Ptma | Ap1m1 | Tmem258 | |
| Mpz | Snrpa1 | Lman1 | Mrpl15 | Atxn10 | Eif5a | Serbp1 | |
| Cdkn1c | H2afx | Tceal9 | Nsfl1c | Ranbp1 | Galk1 | Rps13 | |
| Papss2 | Cacybp | Hspd1 | Anapc11 | Cct6a | Polr2i | Elob | |
| Stmn1 | Gale | Eef1g | Mcm7 | Mrpl34 | Tspan4 | Dad1 | |
| Ldha | Pdrg1 | Krtcap2 | Npm1 | Serpinh1 | Atp5f1 | Rpsa | |
| Plod2 | P4hb | Snap47 | Snhg6 | Dcakd | Rpl11 | Gapdh | |
| Cdk4 | Ldhb | Cks1b | Rnf7 | Atp5j | Rpl14 | Gnas | |
| Slc26a2 | Srm | Tmem97 | Ssrp1 | Tecr | Luc7l3 | | |
| Bex3 | Susd5 | Kdelr2 | Cnpy2 | Serp1 | Ube2e3 | | |
| Epyc | Ltv1 | Selenoh | Tfg | Nme1 | Ywhab | | |
| Pdia6 | Tubb5 | Vdac3 | Lrrc59 | Hnrnpc | Akr1a1 | | |
| Ss18l2 | Gadd45gi | | Mdk | | | | |

| | p1<br>Ccnd2 | Srsf2<br>*Srpl2* | K1M13 | Atp5o<br>Siupa | Rps26<br>Ndufcl | Tsc22d1<br>Igf2<br>Id3<br>Cfll<br>Hsp90abl<br>Rpsl7 | |
|---|---|---|---|---|---|---|---|

| Cxcl12 co-expressed | | | | | | | |
|---|---|---|---|---|---|---|---|
| Il1r1 | Il13ra1 | H6pd | C1ra | Gas6 | Itga11 | Serpina3g | Pkdcc |
| Col3a1 | Apln | Isg15 | C1s1 | Sfrp1 | Col12a1 | Serpina3n | Epas1 |
| Col5a2 | Hs6st2 | Steap4 | P3h3 | Slc7a2 | Selm | Ghr | Colec12 |
| Igfbp5 | Bgn | Emilin1 | Fxyd1 | Comp | Ebf1 | Osmr | Egr1 |
| Sned1 | Slc16a2 | Htra3 | Rcn3 | Bst2 | Slfn2 | Lifr | Lox |
| Ifi203 | Capn6 | Nsg1 | Fcgrt | Rnf150 | Col1a1 | Snhg18 | Iigp1 |
| Nenf | Gpm6b | Sod3 | Saa3 | Ier2 | Igfbp4 | Ly6e | Synpo |
| Pfkfb3 | Cp | Pdgfra | Prss23 | Nfix | Mrc2 | A4galt | Pdgfrb |
| 1110008P14Rik | Dclk1 | Cxcl5 | P2ry6 | Junb | Timp2 | Fbln1 | Efemp2 |
| Lcn2 | Mme | Cxcl1 | Adm | Mmp2 | Lgals3bp | Pdzrn4 | Pcsk5 |
| Serping1 | Ptx3 | Plac8 | Il4ra | Mt2 | Sfrp4 | Rtp4 | Ifit3 |
| Ube2l6 | Tbx15 | Spp1 | Ifitm2 | Mt1 | Aspn | Mylk | Ifit1 |
| Fibin | Slc16a1 | Pkd2 | H19 | Cdh11 | Ogn | Fstl1 | |
| B2m | Vcam1 | Tgfbr3 | Igf2 | Hp | S1pr3 | Nfkbiz | |
| Eid1 | Penk | Oasl2 | Rspo3 | Stc1 | Cxcl14 | Abi3bp | |
| Fgf7 | Svep1 | Col1a2 | Bicc1 | Pdlim2 | Gas1 | Tmem45a | |
| Cpxm1 | Ugcg | Ptn | Col6a1 | Slc39a14 | Vcan | Col8a1 | |
| Ism1 | Plpp3 | Rarres2 | Aes | Tsc22d1 | Pik3r1 | Adamts5 | |
| Cst3 | Podn | Tmem176a | Igf1 | Mmp13 | Il6st | Kcnj15 | |
| Lbp | Hivep3 | Loxl3 | Dram1 | Mmp3 | Stxbp6 | Fndc1 | |
| Wisp2 | Col8a2 | Cyp26b1 | Dcn | Clmp | Hif1a | Sod2 | |
| Zbp1 | Nbl1 | Antxr1 | Lum | Nnmt | Zfp36l1 | Thbs2 | |
| Srpx | Mfap2 | Slc6a6 | Ndufa4l2 | Islr | Npc2 | Angptl4 | |
| | Dhrs3 | Cxcl12 | Lrp1 | Loxl1 | Ltbp2 | Cyp1b1 | |

| Ifitml co-expressed |
|---|

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1500015010 Rik | Seφ i¾1 | Cp | Ifitm2 | 1500009L16Rik | Ctsh | Tgfbi | Apod |
| Crocc2 | Cst3 | Gperl | Ifitml | Scara5 | Zicl | Hifla | Abi3bp |
| Snedl | Ptgis | Gngl1 | H19 | Zic5 | Zic4 | Aspg | Epha3 |
| Fmod | Slcl6a2 | Cemip | Akapl2 | Mmpl3 | Ebfl | Fblnl | Smoc2 |
| | Fabp5 | Adm | Gjal | Clmp | Sfφ4 | Kng2 | Thbs2 |
| | | | | | | | Epasl |
| | | | | | | | Prdm6 |

| Matn4 co-expressed | | | | | | | |
|---|---|---|---|---|---|---|---|
| Spats21 | Kcns1 | Penk | Eln | Pdgfrl | Mfap4 | Igfbp4 | Nov |
| Igfbp5 | Matn4 | Mfap2 | Cpxm2 | Igfbp3 | | | |

| 2-cell | | | | | | | |
|---|---|---|---|---|---|---|---|
| Tcl1b1 | Pxt1 | Omt2b | Inpp4a | Stbd1 | Ampd3 | Stk36 | Rnf182 |
| Dusp7 | Smad3 | Obox5 | NA.15103 | NA.13579 | NA.15121 | Sytl4 | NA.12407 |
| Zbed3 | B4galt6 | Itga9 | Mllt3 | Man1c1 | Angel2 | Tmem92 | Ptpre |
| Tcl1b2 | X7420426K07Rik | Ptprr | Mcc | Sh3bp1 | Sipa1l1 | Akt3 | Zcchc2 |
| Gm839 | Creld1 | NA.15153 | Slc15a5 | Kit | Gm21762 | X9130023H24Rik | Tcstv1 |
| NA.13991 | Lbx1 | Hmces | Fam167a | Nos1ap | NA.9588 | Hoxa7 | Spesp1 |
| Gm1965 | Gad2 | Mfsd2a | Pip5k1b | Mvb12b | Gm13023 | Coro2b | Ppp1r3d |
| Phf1 | Mn1 | Tgfb2 | Bmp5 | Prr5l | Olfr288 | NA.15065 | Grip1 |
| Tcl1b3 | Ccdc69 | Plekhg1 | NA.15072 | Adm2 | Gm12735 | Ctdspl | Hsd17b13 |
| Siah2 | Pak7 | Mcu | Oosp1 | Igsf11 | H2.Q6 | AU015836 | Tet3 |
| Tcl1b4 | Stradb | Myo3a | Vil1 | Aida | NA.15138 | Cngb1 | Wdr25 |
| Phc2 | Rfpl4 | Gm11131 | NA.2207 | Rimkla | Wasf3 | NA.10579 | Mapkbp1 |
| Tcl1 | Fam43b | Zscan4d | Bcorl1 | Jazf1 | Polm | Usp46 | Fchsd2 |
| Tbx19 | Gli3 | Bmp2k | Zfp513 | Tshz1 | Man2a2 | Cdc42se2 | Fam19a2 |
| Obox3 | Grm2 | Btg4 | Plxnc1 | Gng3 | Gm9125 | Gyg | Ssh1 |
| NA.6855 | Parp12 | Fyn | F2r | Dpysl3 | Usp21 | Igdcc3 | Errfi1 |
| Gm12789 | D6Ertd474e | NA.13288 | Kcnk18 | Gfod1 | Tmc8 | Plag1 | Fbxw22 |
| Wee2 | Reep2 | Pik3cd | Klhl8 | Tesc | Ccdc92 | Arntl2 | Ajap1 |
| Bcl2l10 | Btbd2 | Adcy5 | Cby3 | Oosp2 | Lrrc4 | Fbxw14 | Gm20767 |
| Rph3a | Gpr68 | Smpd3 | | Syt11 | NA.10324 | Catsperg1 | |
| | Slc45a3 | Pld1 | | Tmcc3 | Sipa1l2 | Itpk1 | |
| | | NA.80 | | | Nlrp4e | | |

183

| Gm6507 | Iqca | AUO 16765 | Cpal | Elavl2 | Gja3 | Prss46 | Epha3 |
|---|---|---|---|---|---|---|---|
| Th | Tubg2 | Oasld | Sbkl | Plek | Ramp3 | Spire 1 | DpplO |
| Musk | Kcnhl | Gml775 1 | Zscan4c | Spocdl | Orail | Nlgnl | Slc30a3 |
| NA. 103 66 | X2210019I 11Rik | Krt84 | Slcla4 | Dennd3 | Sufu | Dbnddl | Gm2807 8 |
| Tmcc2 | Accsl | Uncl3c | Ablim2 | Lip lb | Lefl | A630095E1 3Rik | Itga8 |
| Fa2h | X2010107 G23Rik | Fmn2 | Manse 1 | Pcdhl5 | NA. 15 19 | Nr2el | NA. 15 1 23 |
| Spry4 | B4galt2 | Angptl2 | NA. 15 1 14 | Nav2 | Nav3 | Gml3 103 | Taf9b |
| Tbxa2r | AC126035. 1 | X9530082P 21Rik | Peakl | NA. 107 49 | Gstm5 | Lhx8 | Plxna4 |
| Rims1 | Uspl71c | Pdgfrl | Colgalt2 | D6Ertd5 27e | Smox | Nrep | Mfsd6 |
| NA.406 2 | Rab3d | Rasd2 | Zfp30 | Timd4 | X4933404 012Rik | Pla2g4c | Pou4fl |
| Papd7 | NA. 10463 | Per3 | Rapgef5 | Efna5 | Vps9dl | Rasa4 | Fgfrll |
| NA. 142 00 | Eif4e3 | Smiml4 | Ctif | Rspo2 | Sortl | AI987944 | Evl |
| NA.729 4 | Prkaca | Hipk2 | Eif4elb | Mamll | Shank2 | NA. 12447 | Gdf9 |
| Gml 182 7 | NA. 12521 | Slc24a3 | Ifitm6 | LsmlO | X4933415 A04Rik | Prmt2 | Dnasell 3 |
| NA.553 9 | Mmp2 | AA415398 | Cobl | Slc6a7 | Faml l7a | Dact3 | Shroom 4 |
| NA.354 1 | Axin2 | St6gall | Zfp46 | Gml566 8 | Jade2 | Magil | Fbxo43 |
| Uspl71b | Fzd2 | Ctdspl | Ppplr9b | Lrrc8a | Ptcra | Gml3 191 | Unci 3b |
| Bmpl5 | Cbx2 | Adarb2 | Mypop | Txndc2 | Dpfl | Emilin2 | Scg3 |
| Tfap2e | Fmnl3 | Foxml | Mlltl 1 | Gm2878 4 | Pld6 | Smagp | Fgf7 |
| Rbm38 | Hpcall | Adamtsll | Cdh4 | Efcabl2 | Ets2 | Spinl | C87499 |
| Zdhhc8 | Prrgl | Arhgap20os | Ccnjl | Tef | Elmod3 | Tbcld8 | Tubb3 |
| Lztsl | Sebox | Lingo2 | Midn | Nhsll | Acot3 | Gphn | NA.232 |
| Tcllb5 | Oboxl | Tox3 | Tspan5 | Glis3 | Apol7b | Synm | Limdl |
| Slco3al | Zfp957 | Bmp6 | Gbas | Mark2 | Pacs2 | Tmem72 | Esytl |
| Dclk2 | Taar2 | Fsdl | Ttbkl | Apela | Tmeml08 | Fkbp5 | AF0670 6 1 |
| Tulp3 | Rassf5 | Gm21818 | B4galnt 4 | Adam33 | Dmwd | Clvs2 | Trakl |
| NA. 189 1 | Afapll2 | Tcf20 | Gml 138 1 | Cacnalh | Ubash3b | Rnf220 | Slc22a2 3 |
| NA. 15 1 24 | Tmeml84b | E330012B0 7Rik | Rragc | AI85470 3 | X23 100611 04Rik | Platr22 | |
| Rgsl7 | Omt2a | Tob2 | Nrpl | Zfp703 | Fbxw24 | B4galt4 | |
| Zfp352 | Trim75 | X4933427 D06Rik | AU0227 5 l | Creb314 | Ccno | Sgms2 | |
| NA. 104 33 | Pcdh9 | Dnah7c | Ncehl | Fzd7 | Acox3 | Aicda | |
| Cmya5 | Foxj2 | Angel 1 | Lrrcl6a | Mmpl9 | BC147527 | Glisl | |
| | Tmtcl | Prlr | Oosp3 | Khdclb | NA.3893 | E330021D1 6Rik | |
| | Prkdl | Ccdc6 | Faml99 x | Prrx2 | Eef2k | Oogl | |
| | Ppmlh | Shb | Myadml | Kmt2d | Farpl | Sh3rf3 | |
| | NA.95 12 | NA.7047 | | | E330034G1 9Rik | Ttyh3 | |
| | | | | | | C330021F2 | |

184

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Cdr2 | Nrsn2 | Ybx2 | 2 | Prss45 | Fbxwl8 | 3Rik | |
| Mfap2 | Trim60 | Kifl7 | Ms4al | Trim7 | Kpna7 | N4bpl | |
| Gnal2 | Slc25a48 | Lmxla | Diras2 | 117 | NA.613 1 | Dcakd | |
| Cntnapl | Snph | Pou2f2 | Pde4c | Sbf2 | Tbcld2b | Obox2 | |
| NA. 102 80 | Antxrl | Ninj l | Pptc7 | Tcf7 | Fhod3 | Gramd2 | |
| Mesp2 | B020004C17Rik | Cables 1 | D13Ertd 608e | Ksrl | Pygol | Tmeml80 | |
| Vrtn | Derl3 | Meis2 | Gml605 0 | Rundc3 b | Ap3m2 | Prr32 | |
| ParplO | Ahdcl | | Faml3 1 a | NA. 157 9 | | Ccdc88a | |
| Fam222 a | | | Obox7 | Lmol | | | |
| Pkd212 | | | Cythl | | | | |
| SamdlO | | | Rnf26 | | | | |
| Tbx4 | | | Nobox | | | | |

| 4-cell | | | | | | | |
|---|---|---|---|---|---|---|---|
| X1700019 E08Rik | Esam | Otopl | NA. 15084 | Tmem210 | E030044 B06Rik | Ptdss2 | NA.9870 |
| Gcml | Tmc5 | Caapl | Eif4e | Pdlim4 | Arrdc3 | Vmnlr90 | Toporsl |
| Gm26815 | Kcne3 | Tc2n | Ttc30al | Lamp2 | Spink2 | Cracr2b | Mlfl |
| Handl | Dnmt3bos | Kcnfl | Ccr4 | X181003 4E14Rik | Rhoq | P3h4 | Gm26745 |
| Esxl | Nags | Slc38a2 | Hoxb9 | Pcolce2 | Ddx60 | Gm26632 | X1700092 M07Rik |
| NA. 13936 | Zfp644 | Gm9918 | Tmem5 | NA.55 1 | Cdkn2a | Clec2g | Akapl2 |
| Mbnl3 | Tspan6 | Spata25 | Zfp273 | Pgm211 | Psma8 | Gml6302 | Cnnml |
| Tgfbl | Gm9732 | Myc | Nabpl | Chicl | Best2 | Elf4 | Tmem63a |
| NA. 11398 | Sycpl | C2cd4b | Adaml9 | Trim40 | Gml5 12 8 | Slc25a46 | 01fr815 |
| Ltb | NA.965 1 | Gm595 | Ythdc2 | Rmdn2 | Dppa2 | Tmem47 | Tacr2 |
| X1700003 E16Rik | AI606181 | Rbm41 | Gramdla | Ddit4 | Mettl20 | Sowahc | Adamtsl4 |
| Pil6 | Foxal | NA. 1261 1 | Rnfl 1 | Tramlll | Ei24 | Mxra7 | RdhlO |
| Calm5 | Ccdc89 | Cacng7 | AC133 10 3.1 | Ptprcap | Nr2c2 | Apls3 | Pxdcl |
| Tmem37 | Nrg2 | Jakmipl | Ctsl | Epm2a | D930016 D06Rik | Hfml | Cyr61 |
| 01fr836 | Eidl | NA.5 175 | Crabpl | H3f3b | X493050 3E14Rik | Ccdc57 | Prpf4b |
| Map7dl | Rtn4r | Zswim5 | Uhrf2 | Agbl2 | Soxl5 | Wipfl | X1700123 IOlRik |
| Tceal8 | P4ha3 | Obox8 | NA.556 | Igfbp3 | Six4 | NA. 11442 | NA. 1350 |
| Nfatcl | Cavl | Syne3 | Faml22b | Upk3b | Ramp2 | Wdr5b | NA.9846 |
| Wbp5 | NA.7320 | Lrrcl5 | Cbfb | X603044 3J06Rik | NA.44 | Plin5 | Unc5cl |
| NA.7187 | Tex 15 | Iraklbpl | Lpar6 | Robo4 | | Dixdcl | Zfp948 |
| Tcf23 | Rbml2 | Kcnk5 | Gm6871 | | | Gml l23 | NA. 13261 |
| | Bexl | Pdlim3 | | | | Brwd3 | |

| Noto | NA.8609 | Mat2a | Gml6010 | Ddias | Gm5773 | Amigo2 | Tdpoz4 |
|---|---|---|---|---|---|---|---|
| Pet2 | Gml l961 | Gml4443 | Ahil | Gml5389 | Slcl2a2 | NA.5634 | Zfp799 |
| Nuprl | Fgr | Klfl7 | Spaca6 | Lame2 | Slc35f5 | AC125 14 9 . 1 | Nafl |
| 43353 | X3 110021 N24Rik | Lixll | Ube2e3 | Calb2 | Lbhdl | Ppwdl | NA.990 1 |
| Myh7 | X9030407 P20Rik | Trpd5213 | Xcrl | NA.337 | H2afx | Gm26522 | NA.7995 |
| Zfp457 | Tbcldl2 | Gml4124 | Zfp874a | Mtmr6 | Arl4c | Rasgefla | Gml0509 |
| Nxf2 | NA. 15089 | Fscnl | Cenpq | Fam65c | NA. 1005 8 | Zfp874b | Gm28875 |
| Prdml4 | NA.7248 | Platr25 | NA.3213 | Lrifl | FkbplO | Cyb561dl | Rnd2 |
| Dlx3 | Abcb5 | Trim2 | Ggt7 | Ehd2 | Krt28 | Ttc29 | Nudtl6 |
| X4930502 E18Rik | Sphkl | Tuba3b | Zfp85 | Chmbl | Set | Gm7334 | Rsrpl |
| X I700065 O20Rik | Hivep2 | Wnk3 | Ctsk | Cpz | Cbx3 | NA. 15 101 | Uty |
| WntlOb | Beanl | Map7d2 | Gm28043 | Prep | Sdc3 | Uaca | Vgf |
| Bbsl2 | Spsb4 | Morc4 | Ctag2 | Slc24a4 | Cyp2j6 | NA.8430 | NA. 12375 |
| Lrrcl9 | NA.9430 | Kalrn | 01frl43 | Zfp950 | Endog | Obox6 | NA.2730 |
| Phyhipl | Armcx4 | NA.93 16 | Mier3 | Mesdcl | X943002 OKOlRik | Nanos2 | Unc45b |
| Pla2g4a | Zfp758 | Platr3 | Isll | Zfp729a | Atp2c2 | X4930505 A04Rik | Pigw |
| Tceal7 | Tnfrsfl 1a | Cyplal | Pank3 | Gm8 104 | Gml0550 | Trpc5os | D730003I 15Rik |
| Siahla | NA.5916 | Sox30 | Ap4b l | NA.539 | Coll7al | Rnpc3 | Gm4285 |
| Trim56 | NA. 15077 | X3222401 L13Rik | Pik3c2a | NA. 1506 4 | Wsbl | A930003 A15Rik | Slfn9 |
| Magea8 | Pkdll3 | Gml6185 | Capn9 | Hmhal | Slcl9al | Pnn | Edaradd |
| Hesl | Hicl | NA.264 | Foxfl | Wdr54 | Rsph9 | NA.4962 | Slc5a3 |
| Btgl | Chrnd | Gml7056 | Tnfsfl3b | Jrkl | Zfand5 | Hnrnpll | L3mbtl3 |
| Zfp239 | NA.407 | Hsdl7bl4 | NA. 1494 | Pax6 | Seppl | NA. 186 | Pin |
| Gml0226 | Magea5 | Tmem229 b | Rnftl | Etnkl | Relb | Ctsb | Gml 508 |
| P2iy4 | X1700019 B21Rik | Usp44 | Notch4 | Cebpa | Gm2399 | NA. 10139 | NA.4305 |
| Usp9y | Crybal | Gml23 15 | Hsf3 | Atg3 | X4930447 C04Rik | | |
| Gm5930 | Pm20d2 | Gbxl | Aebpl | Fzd4 | Prss36 | NA. 10456 | |
| Sox21 | Sec 16b | Gm8126 | Tex37 | Hkdcl | NA.222 | Gabra4 | |
| Selenbpl | Mastl | Nufip2 | Rhox9 | CldnlO | Elovl3 | Col5a3 | |
| Gm6526 | NA. 1742 | Ubaly | X4930432 K2 1Rik | SmimlOl 1 | Npas2 | Pbld2 | |
| NA. 15085 | Nrxn2 | Irf2bpl | Soat2 | Gm2678 2 | Nme5 | Cd81 | |
| X I700049 G17Rik | Acsl4 | Aim2 | Hesxl | Zfp945 | Mysml | Lrrc46 | |
| Gm53 | B230219 D22Rik | NA.4044 | Vatl | Slc26al0 | C130026 I21Rik | Gm7073 | |
| Mycn | Gml55 18 | Ranbp6 | Nlrp6 | Gm6268 | NA.6224 | Fam228b | |
| Gml5097 | Ptprzl | Id4 | Hrk | NA. 180 | Lrrc58 | Ctsc | |
| NA. 10436 | NA. 15 112 | Platr23 | Prrtl | Cardl4 | NA.7446 | Mrap | |
| | A930017 | Spic | Zfp40 | | | | |
| | | | Argl | | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Fbnl | K11Rik | Gml7404 | Man2clos | Rimklb | Bhlhb9 | Grikl | |
| Adgrbl | NA.4501 | Chadl | Gm5532 | Zfp953 | Mplkip | Rblccl | |
| Klf2 | Mbnll | Cede 152 | Hnrnpal | Fgf4 | Sparcll | NA.7081 | |
| Fam212a | B3gnt8 | 01fm3 | Tnfrsfla | Tenm3 | NA.7433 | Dgat2 | |
| Fgf3 | Gm29087 | NA. 12133 | E112 | Mirl7hg | Cfap73 | AC133 10 3.5 | |
| Tcp1 ll2 | Dsc3 | Ffar4 | Dszf5 | Ambn | Gml4168 | Lcat | |
| Sema6b | Irf7 | | | Btbd3 | Slcl6al4 | NA.4426 | |
| | Plek2 | | | Fbln2 | Avl9 | | |
| | | | | Per2 | Ogn | | |
| | | | | | X170001 9G24Rik | | |

| 8-cell | | | | | | | |
|---|---|---|---|---|---|---|---|
| NA.71 10 | Xist | Lif | BC052040 | Zfp936 | Slc7a7 | NA. 13976 | NA.3445 |
| Cyp2d9 | Arhgefl6 | Qpct | Ly6a | NA.5874 | Gml4582 | Arfip2 | Plekhfl |
| Ackr3 | NA.689 | NA.88 | Prdx6 | Vpreb3 | Adgrg3 | NA.9630 | Cd59a |
| Perp | Kcnv2 | Nr4al | Chmp4c | Vsxl | NA.6826 | Pmaipl | Tfcp211 |
| Cstl3 | Fkbp9 | Grinl | X2410141 K09Rik | Kctdl | Rpl39 | Gcfc2 | Gml3212 |
| NA.9215 | Gas6 | Nup62cl | Fbxl20 | Ccdc84 | Nog | Gml305 1 | Parpl6 |
| Cpne3 | H60b | TrmtlOb | Tyms | Gstal | Gm26584 | Gml9667 | Nln |
| Dok2 | Gm26692 | Exoc314 | Eps812 | Zfp275 | Fbp2 | NA. 10925 | NA. 1527 |
| Cd28 | Slcl2a7 | I830077J02Rik | A230083G16Rik | Hopx | Clcnka | NA.5489 | NA.4804 |
| Phlda3 | Plagll | NA.7942 | Prkra | NA.3556 | Gml4401 | Lrpapl | NA.3235 |
| Cartpt | Ppmlk | Hsh2d | Gm9776 | NA.3384 | Mef2d | Regl | Esrp2 |
| Cthrcl | Ppfibp2 | Cd300a | Laspl | Vgll4 | Myo 15b | Golga7 | Ly96 |
| Msc | Gml2705 | Ptpn6 | Cstf3 | Ptdssl | Cdc42ep3 | Chordcl | X9030624J02Rik |
| Stxbp6 | Vavl | Gm6020 | Akrlc21 | NA.6297 | NA.2700 | I122ra2 | NA.3453 |
| NA.810 | NA.8401 | Siglecg | Hoxa9 | Plcdl | Hhex | Gml 630 | Mfs d8 |
| Stfa211 | Pla2g7 | Prrg3 | Ecell | Gm265 14 | Gml2289 | Ehdl | Slc45a4 |
| Pdzd3 | Dkkl | Zfp932 | NA.4219 | NA.4998 | Hmga2 | Pkp2 | Urgcp |
| Gm27204 | Sbp | Gm21060 | X9430060 I03Rik | NA.7408 | Zfp429 | Pdcd6 | Igbpl |
| Anxa3 | Hsdl7bl | X1010001 N08Rik | Mocos | Gml6503 | Pou5fl 4335 1 | Efnal | Lgals8 |
| NA. 1015 | Rragd | Rnfl38 | Slc6al4 | NA. 10479 | AdgrG | Ttc39b | NA.4193 |
| Vrk2 | Tmem8 1 | Sync | Smpdl3a | Plxnb2 | Faml98b | Cyba | Atp6v0e 2 |
| Npy | H60c | Xkr9 | Nudtl 1 | Slcl0a4 | Hprt | NA. 14015 | Chptl |
| Tspanl | Svil | Gml7655 | Krt7 | Salll | NA.71 1 | Cd209e | NA.588 |
| Stard4 | Pramel5 | | | | Grk6 | NA.9466 | |
| Lectl | Irf5 | | | | | Gm205 15 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gyltllb | Deaf 1211 | Eno2 | NA.5 168 | NA. 1214 8 | Atp2bl | A530040E 14Rik | Adam4 |
| Nxpe5 | Gm413 1 | Amph | Oimdll | C3arl | Satl | NA.443 1 | Zfp607 |
| Dynap | X4930550L24Rik | Cede 150 | NA.4188 | Gml306 2 | Fam2 17b | Rnf32 | Atp6v0a 4 |
| Gml5446 | Zfp52 | Cdc42epl | Hspa8 | Fndc3cl | Etohil | Ly6g6e | Arhgap2 7 |
| Zfp934 | NA.3646 | NA.4813 | Rassf7 | Dpyl912 | G430049 J08Rik | Ldbl | Cdh1 |
| PlatrlO | X4930522L 14Rik | Eda2r | Star | Ano2 | Fam83b | Gml 1541 | Il17re |
| Amot | Slco2al | Hes2 | Pkdlll | NA. 1390 0 | Pde7a | Gm2366 | NA.3823 |
| Id3 | Gm26836 | Etl4 | D930020B 18Rik | Iqgap3 | NA.4566 | Prrl9 | NA.4035 |
| Amotl2 | Ap3b2 | Vangll | Arhgapl8 | Sh3d21 | Cldn4 | Cmtm5 | NA.4009 |
| Gm26740 | NA.41 12 | Atp8b4 | Ppp2r2c | 43 160 | Foxf2 | Tmem45a | Lpin1 |
| Abcbla | NA. 10665 | Cav2 | Denndlb | Akp3 | Pank4 | NA.9621 | Atg4c |
| Diaph2 | Tmem245 | Slc29a3 | BC05 1665 | Glt28d2 | B930036 NIORik | NA.336 | Alg13 |
| Akrlcl4 | Pik3r6 | Nradd | Diiall | Gm | NA.7030 | Gml0687 | Rad23a |
| Ciyab | Tsix | Tmem253 | Klf8 | NA.2621 | Gm26668 | Zfp4 18 | Gm2653 8 |
| 1133 | Hsdl7bl 1 | NA. 1630 | Gml3235 | Cwh43 | Gabrd | Gml976 | Prr15l |
| Slcl9a2 | Zfp354a | Ddahl | B4galtl | NA.7337 | Tbx3 | NA. 1763 | NA.7290 |
| Epasl | Gml 1O | Ano9 | NA.5 135 | Sh3tcl | X9430002A10Rik | NA.7085 | Upf3b |
| NA. 1618 | Bves | Acp5 | NA. 1892 | Pinlrtl | Ctsf | Acyp2 | Slco4c1 |
| Pcdhbl6 | Xlr | B2303 12C 02Rik | Ckslbrt | C030039L03Rik | NA.6 | Oxctl | NA.5912 |
| Bex4 | AI467606 | LITC23 | LiTc37a | Caldl | Gm27206 | Pigz | Emilin1 |
| Tmem64 | Mtml | Cux2 | Krt27 | Akap2 | Rnf208 | Tpd52 | NA.5335 |
| Bmp8b | Ccngl | NA.9543 | Wnt3a | I113ral | Bhmt2 | NA.47 | Tmem14 4 |
| Gml0139 | Arhgdib | Gm6712 | Smocl | NA.9845 | NA.293 1 | Mllt6 | Zfp599 |
| Gpc4 | Faml24a | N A .7720 | Igsfl | Sbpl | NA.691 | Plcgl | |
| Vnnl | Slc52a3 | Faml29a | NA.5696 | NA. 1027 | Adam21 | Pnpla2 | |
| Rbmsl | Gml3 154 | NA.2889 | K cnh7 | NA.3 116 | Serine 1 | Gml5 137 | |
| Apob | Suox | Gml0324 | Gm 13242 | Alcam | NA. 1264 9 | Dnajc6 | |
| X9330185C 12Rik | NA.2957 | Slc29a4 | Sema5b | NA. 1390 6 | Mybpc2 | X2410018L13Rik | |
| Camk4 | Fgfl3 | NA.2540 | NA.9923 | Imnt | Runxl | Actnl | |
| NA.559 | Parva | Gml25 14 | NA.5 13 | Card 11 | Vtn | NA.223 | |
| Mpped2 | Casc4 | Cd53 | GrM3 | Asap2 | Fancb | Rbks | |
| Poflb | X9230009I02Rik | Msmol | Lparl | Smim22 | KlflO | Nrtn | |
| Papss2 | F12 | Rampl | NA.3947 | Sycn | Gm26624 | Fut9 | |
| Tbx20 | X22 10404O09Rik | Postn | Isl2 | Ak7 | NA. 1030 3 | Ednrb | |
| Gng2 | SlOOal 1 | Havcrl | Fes | Nprl2 | NA.7385 | Zfp458 | |
| Nr2f2 | X5430403 | Ttpa | Nap 112 | | | Itpkb | |
| Rarb | | Gjb3 | Sh3glb2 | | | NA. 11397 | |
| Gml0772 | | Ahsg | Nck2 | | | | |

| Zfp157 | G 16Rik | Strada | Gata6 | Zfp422 | NA.487 | NA. 1522 | |
|---|---|---|---|---|---|---|---|
| | Steap3 | Reep1 | Slc36a3os | Alg6 | NA.2929 | NA.991 1 | |
| | Matn3 | Ncf2 | NA. 14579 | Npnt | Rdh5 | NA.2756 | |
| | Slc22al3 | NA.4991 | Bok | NA.424 | NA.5637 | | |
| | Fgd4 | | | Psrcl | Vps33b | | |
| | | | | Sfrpl | | | |
| | | | | Ace2 | | | |

| 16-cell | | | | | | | |
|---|---|---|---|---|---|---|---|
| Gm2245 | H2afy | Khdc3 | Tbca | Erlec 1 | Adam9 | NA. 12986 | Nipal |
| Fabp5 | Rhob | X4930558J18Rik | Mycl | Slc7al5 | Pomtl | Egfl7 | Tppl |
| Gml7067 | Trip6 | Gml4409 | Phlppl | Vcpkmt | Gjb3 | Ormdll | Gm4673 |
| Apoal | Tmsb4x | Top2b | Sqstml | Trim47 | Acad 12 | B3gnt3 | Slc35al |
| Stat6 | Slc6al3 | Ank2 | Hbegf | Bcl91 | Tmeml35 | BC052040 | NA.5230 |
| Capn6 | Plk5 | NudtlO | Serpinb6a | Evpl | X2610528Jl lRik | Paqr5 | Hdac3 |
| Abcal | Col4al | Pvrll | Acpl | Actgl | BC029214 | Pfn2 | Whamm |
| Gml4305 | Shkbpl | Anxa9 | Nanog | AU021092 | Them5 | Gml4403 | Gpx2 |
| Eomes | Mgst2 | Hal | Reml | Cdkl8 | Atp8al | Vmn2r29 | Trappc 1 |
| Zfp3611 | Cdcl23 | Slc2al | Sppl | Dok2 | Psmg2 | Gstpl | Tmeml98 |
| Sox2 | Dsg2 | Acaa2 | Texl9.2 | Cldn23 | Sik2 | Gml7087 | NA.4039 |
| Sh3bp5 | Mpzl2 | Lyrm9 | Pdzklipl | Nsmaf | Wnt6 | Slc5a2 | X311005 2M02Rik |
| Ptgdr | Glrx | Slprl | X1700095 A21Rik | Cpxml | Bre | NA.73 16 | Adprh |
| As3mt | Frrsll | Pgapl | Camkl | Impadl | Elf3 | Npcl | Thrsp |
| Pmaipl | Gss | E130012A19Rik | Gml4327 | Crip2 | Pigz | Pmsl | NA. 10775 |
| Dokl | Hebpl | Bend7 | Lamcl | Itga7 | Sccpdh | Gm26578 |
| Slc37a2 | Sox7 | Xbpl | Alg8 | NA.61 14 | Lrmp | Spcs3 | NA.385 1 |
| Tinagll | Cbx4 | Zcchcl6 | Napll3 | Eps811 | Vapb | NA.499 | Aasdhppt |
| Aldhlb1 | Fbxo3 | Mapt | Vpsl3c | Camk2d | Bhlhal5 | Slc4a2 | Pkp2 |
| Mafb | Pnma2 | Arl6ip5 | Epcam | Alcam | Gml0605 | Gatadl | Plgrkt |
| Lypd8 | Fam92a | Pou2fl | Dpysl4 | Assl | Hsp90aal | Atp2a3 | NA. 14210 |
| BC048679 | Ddx3y | Cited4 | Fas | Mospd2 | Nsdhl | Fancb | Itm2b |
| Gml4412 | Wfdc2 | Tbxl | Tgfbr2 | Lip 11 | Sdcbp2 | Rac3 | Duspl 1 |
| Otx2 | Msx2 | Zip 119b | Dmcl | Trim21 | Faml32a | Mthfsd | Lgals9 |
| NA. 186 | X5730507 COlRik | X1700086 P04Rik | Ctgf | Slc24a5 | X2700068 H02Rik | Acadvl | Sdhaf4 |
| | Herpudl | Cstal | Sult4al | Csf3r | Kbtbdl3 | NA. 10404 | Emp2 |
| | Hspalb | Efnbl | Zfp459 | 43352 | | Tfcp211 | Idhl |
| | AdamtslO | Hemkl | Zfp688 | Lrrc75b | | | |
| | | | | NA. 13 142 | | | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| 6 | Mdhl | X4930522 L14Rik | Cgrefl | Map2k3os | NA. 102 | NA. 1896 | Zfp850 |
| Oxt | Rhoc | Hormad2 | NA.92 | Prkce | Gimap9 | X101000 lB22Rik | Txndcl7 |
| BC05 1 142 | Ier2 | Cd82 | Naal 1 | X4930563 D23Rik | Gm4262 | Erf | Apeh |
| Kcnn4 | Slfn3 | Map3kl | NA.388 | Ank | NA.6479 | Slc28a3 | Gml0439 |
| Zfp93 1 | Zfp759 | X I500009 L16Rik | Tdrp | Dact2 | Ralb | Junb | NA. 1925 |
| Pletl | B3galt2 | Phfl 1d | Pcbdl | Pacsin3 | Tmeml7 | Zfpl l9a | Cnn3 |
| Ppl | Laccl | NA. 13623 | Slco2al | Hmcn2 | NA. 1999 | Perp | Mmpl5 |
| Chpf | Tnsl | Trim38 | Cyb5rl | Eef2kmt | Leprot | NA.369 | Cxcr6 |
| Tspan3 | Tmem45b | Vps29 | Magea2 | Chchd7 | Ube2q2 | Calcoco2 | Foxb2 |
| Hyal2 | Tapl | Tbllx | Prokrl | Zfp248 | Lmfl | Gm28085 | Lama5 |
| Fstl3 | Slc38a4 | Lsr | Mbnl2 | NA. 10780 | Tmeml4 7 | Clqa | X170008 0O16Rik |
| Slfn2 | D10Jhu81 e | I117rc | Mex3b | Clecl 1a | Sh3bgrl3 | NA. 1618 | Gml6136 |
| Dusp6 | Srxnl | Aqp3 | Gml6712 | Sgpll | Tradd | Zfp81 | Asap3 |
| Cat | Spata9 | Zfp429 | Zfp395 | Xlr3a | 111 Orb | Ntf5 | Syngr4 |
| Nppb | Pmepal | Ggtl | Krt8 | Msc | X170008 6O06Rik | Oaslg | Zdhhcl5 |
| Tpcn2 | Gm26853 | Tcea2 | Tceall | Zfp442 | Sdhaf3 | Appl2 | Fam83b |
| Cede 16 9 | Pfkfb4 | Gm5 141 | Gata3 | Gml4418 | Galnt9 | Gnal5 | Rnase4 |
| Elovl5 | Zfp266 | Tmem5 1 | Serinc2 | Usp25 | Ogdhl | Gm6169 | Fbxl21 |
| NA. 122 39 | Cdc42ep5 | Stx7 | Rgsl4 | Ntpcr | Pearl | Cmal | Hdx |
| Zfp326 | Magea3 | A530017D 24Rik | Mocsl | Prosl | Fezfl | Lrrn2 | |
| AI3 173 95 | Chrna3 | X1700003 M07Rik | Tmeml3 1 | Lpp | Svbp | Acot6 | |
| AA467 197 | Gm26624 | Ladl | Vps45 | Trp53il 1 | Larplb | Dmrta2 | |
| NA. 113 35 | Elovl7 | Hint2 | Plpp2 | X2610008 E l lRik | A730015 C16Rik | Skidal | |
| Ptges | Nkx6.2 | Exph5 | Mogat2 | Akrlel | Gm26779 | Ccngl | |
| Smiml | Ctam | Sfrpl | NA. 12035 | Pla2g7 | Cryzll | Trabd | |
| Kirrel | Nfkbiz | Hspel | B2301 18H 07Rik | NA.4703 | Stl4 | X241002 2Ml lRik | |
| Gbp9 | Cyp4fl4 | X9430065 F17Rik | Serpinb6c | Gmpr2 | Egr4 | Tet2 | |
| Ckap4 | Tnfrsflb | Ahcy | Fos | StardlO | Hmgal.rs 1 | Cetn3 | |
| Naps a | Dsp | Magee2 | P2ry2 | Enpep | Lcpl | Sri | |
| Gjb5 | Khnyn | Mageb4 | Lgal s4 | Prss35 | Hadh | Vill | |
| Clic3 | Rndl | Gm7325 | Epb4111 | NA.2001 | Sec 1414 | Msantd4 | |
| Marcks | Hnf4a | Tmem266 | Snrk | Eml2 | Txndcl2 | Abhdl4a | |
| NA.724 9 | Adat2 | Txnl | X2410018 L13Rik | Ghdc | NA.7425 | Gm413 1 | |
| Scd2 | X2200002 DOlRik | Rec8 | Rims4 | X2610301 B20Rik | Histlh2b c | Pnpla6 | |
| | Gabarapll | Tgm2 | Gchfr | Pdzd3 | P2rx3 | NA.413 1 | |
| | NA. 12352 | | Nrgl | Gm5424 | | Smapl | |
| | Shc2 | | Skil | | | Lysmd2 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Adgre5 | | Xkr6 | | | Arhgef5 | Xrcc4 | |
| Faml29b | | Egln3 | | | Sfmbt2 | | |
| Pycr2 | | Man2al | | | Btg2 | | |
| Dcafl211 | | | | | Ndufc2 | | |
| Barx2 | | | | | | | |
| I14ra | | | | | | | |

| 32-cell | | | | | | | |
|---|---|---|---|---|---|---|---|
| Lrp2 | Ezr | Oc90 | Ptprn | Baiap211 | Plod2 | Tcn2 | Fez2 |
| Fhl2 | Fam213b | Mapre3 | Gpr4 | Cdc42ep5 | Phfl 1d | Rnaset2b | Rap2b |
| Capn2 | Xbpl | Gm364 | Ptgrl | Etfb | Pdgfa | Aldh2 | Prkce |
| Sppl | CeacamlO | Gstol | 43352 | Gml2169 | SlOOalO | Dab2ip | Gm2381 |
| BC053393 | NA.5461 | Nanog | Nrl | Mdhl | Tpm4 | Actb | Gucylb2 |
| Hspb8 | Msn | Eml2 | Optn | Pletl | Pgm2 | Cck | NA.7242 |
| Cdx2 | Frmd4b | Lsr | Slc25al3 | Wdrl | Gml4326 | Efhd2 | Histlhie |
| Krtl8 | Glrx | Stl4 | Dqxl | Zfp37 | Xrcc5 | Pank4 | Gmpr |
| Enpep | Gapdh | Nfic | Gm26579 | Histlh3c | Esd | Arvcf | Pla2g6 |
| Elf3 | Gstpl | B2301 18H07Rik | Tmeml25 | H2afy | Actr3b | Gml4327 | NA.2972 |
| Vgll3 | Seφ iηb6c | Gm6169 | Cmip | NA. 148 | D630003M21Rik | Wdr6 | NA.7262 |
| Wnt7b | Epb4111 | Gm7325 | Gml4325 | X1700042G15Rik | Ppplrl4d | Abcg2 | Anxa6 |
| Akrlb8 | NA. 123 12 | Gm26917 | Dtd2 | Adrb3 | Mkrn3 | Mgstl | Fthll7e |
| C2cd4a | Lgalsl | Zfp93 1 | Tspan3 | Gml4399 | Adgrl2 | Aldh3a2 | Cdc42ep3 |
| Bglap3 | Ptges | Rp2 | Srxnl | Fthll7a | NA. 101 14 | Omd | Tradd |
| Rab 17 | D10Jhu81e | Tat | Huslb | H2.D1 | Sox6 | Chrnal | Sccpdh |
| Serpinb9b | StardlO | Epcam | Slc6al3 | Cat | Tnsl | Tdpl | Xlr3b |
| Bmyc | Apoal | Rnfl30 | Adaml5 | NA. 1550 | Emp2 | Sgpll | Figla |
| Cmbl | Cela2a | Gml4403 | Vill | Fgfbpl | Col4al | Ttf2 | NA.14180 |
| Klf6 | Tuba4a | Tmeml39 | Sult6bl | Lgals4 | Ndrgl | Faml29b | Dap |
| Krt8 | H2.K1 | Pycr2 | Mecp2 | Trim50 | Dap3 | Emc9 | |
| Nppb | Hint2 | Plscrl | Tarml | Prkcdbp | Capzb | Tmeml7 | |
| Tppl | Cubn | Mfi2 | Camkl | Tφπι6 | Fhl4 | NA. 102 | |
| Tmem9 | Rnfl28 | Adad2 | Mgl2 | NA. 1546 | Wfdc2 | Vps29 | |
| Dppal | Dusp4 | Dsp | Chstl3 | Cidea | Anp32a | AU021092 | |
| Rhox5 | Ogdhl | Mbp | Myhl3 | Nagk | X23 10015 AlORik | Pard6g | |
| Gm5424 | X1500009L16Rik | Chrnb4 | Barx2 | Slc38a4 | Hist3h2a | Kcnkl2 | |
| | Tet2 | Tfcp211 | X1810030O07Rik | Serinc2 | | X8030474 | |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Id2 | Chmp2b | Exph5 | Ccdc43 | Rgsl4 | Slc37a2 | K03Rik | Hspd1 |
| Gjb5 | Lama3 | Rcanl | Ppmlm | Tpil | Gml4418 | Atplbl | Efcab10 |
| Nek6 | Fbxo3 | X9530059014Rik | Slc24a5 | Gstzl | Hsdl7b4 | A330050F15Rik | Tubb2a |
| Oasla | Elovl7 | Eef2kmt | Xlr3a | Ggtl | Sergef | Hdac3 | Gprc5d |
| Scd2 | Patl2 | Mucl | Tmeml98 | Insig2 | Psme2b | Ftx | Smim12 |
| Atp12a | Cede 13 | Efcab5 | BC05 1019 | Ly6a | Ill lral | Fthll7d | Mtmr7 |
| Gstp2 | Col4a2 | **Nyniin** | Erbb2 | X23 10039 H08Rik | Tpcn2 | NA.4386 | Gsta3 |
| Ngfn pl | Acaa2 | Gm26603 | Cnpy2 | NA.2957 | Sh3bgrl2 | Arl2 | Skida1 |
| Pycard | Acaala | Nlφ4c | Idh3a | Carl2 | Asic3 | Apeh | Idh1 |
| Pafah2 | Apbblip | Susd2 | Dab2 | F2rll | Lurapll | Slc2al2 | Hlf |
| Cstal | Tmx4 | Tst | Mksl | Zfp454 | Plau | Zfp850 | Tcea3 |
| Fam213a | Snai2 | Khdc3 | Gimap9 | Eci3 | Fam83h | Iftl40 | Znrd1as |
| Binl | AI662270 | Plbl | NA. 1892 | Gjb3 | Tφ 53i 11 | Slc2al | Pkm |
| Gm694 | Sox9 | NA.5999 | Hk2 | Ly6f | AA46719 7 | Prkx | Map3k15 |
| Dsg2 | Tes | Tdφ | Marcks | Pηlipφ 2 | Gml4409 | X1700086O06Rik | Ak4 |
| Assl | Trim38 | Gale | Gm773 | Praf2 | NA.5 13 | Cox7b | Gml2828 |
| Gm4737 | Cryz | Gml4322 | NA.83 | Gml4393 | Mettl7al | Faml36a | Myo1e |
| Slc38al | Anxa2 | Cpxml | Cdk5 | Abcb8 | Clic4 | Pwwp2b | Slc4a5 |
| Slc38al 1 | Sft2d2 | Tmprssl2 | Gstm6 | Mras | Acol | Cyb5r3 | Slc2a3 |
| Camk2d | NA.388 | SlOOal l | AtxnlO | Gml4444 | Sh3bp5 | Mapt | Sdr42el |
| Bex2 | X2610528J 11Rik | Hoxd3osl | Smco2 | Bckdhb | NA. 1866 | Vpsl3c | Slc7a6 |
| Sdc4 | Gsn | A230005M 16Rik | Enolb | NA.9436 | Gm4779 | Abcal | Snx19 |
| Rfx4 | Hadh | Hnf4a | Pir | Tbxl5 | Cbr4 | Hibch | Ndufaf3 |
| NA.7440 | X0610009O20Rik | Histlh3d | Gpx2 | Acsf2 | NA.6249 | Micall | Plin2 |
| Tinagll | Plp2 | Bdnf | Csf3r | Slcl8al | MyhlO | Adat2 | Gipc1 |
| Col7al | Abcc4 | Ppp4rl | Atg4c | Hdx | Crip2 | Lpp | Pla2g4f |
| Kng2 | Lcpl | Lta4h | Uhrfl | Apocl | Psmb9 | Srebfl | |
| Adgre5 | Actgl | Dpysl4 | Clic3 | Seφ iub6a | Gm4926 | Arhgap9 | |
| Tnfrsf9 | Fam25c | Tmeml02 | Gstm7 | Zyx | I117rc | NA. 14050 | |
| Mmell | Xk | Trhr2 | Coasy | Rec8 | Sdhaf4 | Tctnl | |
| Lgals9 | NA.92 | Tbllx | Tmem256 | Ppplrl8 | Dokl | Tubalb | |
| Texl9.2 | Fabp3.ps 1 | Kremen2 | NA.529 | Cyb5a | Slc25a39 | Whamm | |
| Gata3 | Ube216 | D130040H 23Rik | Tmem45b | Fblnl | Ccdc42 | Smyd4 | |
| Atxn711 | Nsmaf | Cyp4f39 | Krt23 | Dpyl911 | Aφ 8aï | Cbfa2t3 | |
| Txndc 12 | Cited4 | Tmem266 | Mpzl2 | Tpml | Echsl | Arhgef25 | |
| Clcnkb | Fabp3 | NA.5910 | Sqstml | Gdf3 | Akrlel | Nbll | |
| | As3mt | | Zfp780b | Map2k6 | Nudtl 1 | Mgat4b | |
| | | | | | Gcat | | |

| Trp53bp 2 | | Gss | | | | Adh4 | |
|---|---|---|---|---|---|---|---|
| | | | | | | | |

[0251]    In a nutshell, and further discussed below, we identified notable features within the landscape, including sets of cells classified as pluripotent-, epithelial-, trophoblast-, neural-, and stromal-like based on strong expression of signatures related to these cell types and a set of cells (FIG. 24E, purple) that appeared poised to undergo a mesenchymal-to-epithelial transition (MET) following withdrawal of dox (FIG. 24E, orange). The relative proportions of these subsets at different times differed between serum and 2i conditions (FIG. 24G).

[0252]    Using Waddington-OT, we calculated the ancestor and descendant distributions for all cells and determined the trajectories to/from various cell sets (FIG. 24F, arrows). Briefly, the time course began with MEFs at day 0 in the lower right, proceeded leftward to day 2, and then upward over the subsequent week toward two destinations: the MET Region and the Stromal Region. The cells in the MET Region were predicted to give rise to the pluripotent-, epithelial-, trophoblast-, and neural-like cells, with this last class seen in serum but not 2i conditions. By contrast, the Stromal Region appeared to be terminal: cells entered the region, but our model predicted that they did not leave (FIG. 3 IE).

[0253]    The optimal-transport analysis provided insights into when cell fates emerged. As early as 1.5 days, cells' fates began to concentrate toward either the MET Region or Stromal Region, and the distinction sharpened over the next several days (FIG. 25G). The fate of pluripotent-, epithelial-, trophoblast-, and neural-like cells did not appear to be determined until after withdrawal of dox on day 8. That was, the ancestor distributions of these cell types were indistinguishable on and before day 8.

[0254]    <u>The model was predictive and robust</u>

[0255]    Before analyzing the cell sets and trajectories in greater detail, we assessed the accuracy and robustness of our model. Because current experimental approaches for tracing cell lineage did not provide a rich description of the full transcriptional state of a cell set's ancestors, we developed a computational approach to test the model. Specifically, we used optimal transport between the distribution of cells at times t1 and t3 to predict the distribution of cells at an intermediate time t2 and compared this prediction to the observed distribution at t2.

**[0256]** Our predicted trajectories were accurate, such that the distance between the computational prediction and experimental observation at t2 was similar in magnitude to the distance between the two experimental replicates taken at t2, confirming that the prediction is roughly as good as could be expected given experimental variation (FIG. 24H, FIGs. 30A-30G, Methods).

**[0257]** The optimal-transport analysis was also robust to perturbations of the data and parameter settings. We down-sampled the number of cells at each time point, down-sampled the number of reads in each cell, perturbed our initial estimates for cellular growth and death rates, and perturbed the parameters for entropic regularization and unbalanced transport. In all cases, we found that the interpolation results above are stable across wide range of perturbations (STAR Methods).

**[0258]** In initial stages of reprogramming, cells progressed toward stromal or MET fates

**[0259]** Reprogramming began with all cells exhibiting rapid changes. By day 1, cells showed an increase in cell-cycle signatures and a decrease in MEF identity. MEF identity continued to fall through day 3, by which point nearly all cells showed lower signatures than the vast majority of MEFs at day 0 (FIG. 24D). Over time, cells assumed either Stromal or MET identities (FIGs. 25A-25H).

**[0260]** Cells in the Stromal Region showed distinctive signatures, which fully emerged after withdrawal of dox at day 8; these signatures included a secretory phenotype (SASP), extracellular matrix (ECM) rearrangement, senescence, and cell cycle inhibitors (FIG. 25A). By contrast, the MET Region contained cells with increased proliferation and loss of fibroblast identity (FIG. 25E).

**[0261]** Mapping signatures of distinct stromal cell types obtained across mouse tissues from a mouse cell atlas (Han et al., 2018) showed that the most widely expressed stromal signatures corresponded to embryonic mesenchyme and long-term cultured MEFs (FIG. 31A). Yet, the Stromal Region did not simply reflect "MEF reversion." The gene expression profiles were distinct from (FIG. 3IF) and more heterogeneous than day 0 MEFs, with clusters of cells with signatures that more closely correspond to other stromal cell types, such as those found in neonatal muscle and neonatal skin (p-values < 0.01) at levels 20- to 30-fold higher than day 0 MEFs.

[0262]    The proportion of stromal cells peaks several days after dox withdrawal (at -64% of cells at day 10.5 in 2i conditions and day 11 in serum conditions) and then declines through day 18, consistent with the low proliferation signature relative to other cells in the landscape (FIG. 24G). A subset of stromal cells expresses an apoptosis signature starting on day 9, which peaks at day 14.5 in -14% of stromal cells in serum conditions and at day 13 in -3% in 2i conditions.

[0263]    Our trajectory analysis allowed us to trace how these fates were gradually established: we found that the ancestor distributions of cells in the Stromal and MET Regions differred by 30%, at day 3 and by 60%> at day 6 (FIG. 25H). A powerful predictor of a cell's fate was its expression level of the OKSM transgene, with high values predictive of MET fate and low values predictive of stromal fate (FIG. 31C); the expression level statistically explained ~50%> of the variance in the logarithm of the fate ratio (MET Region fate probability divided by Stromal Region fate probability) by day 2 and ~75%> by day 5 (FIG. 31C). Importantly, the divergence was gradual and could not be described by a simple graph with a sharp (that was, zero-dimensional) branch point. Indeed, our optimal-transport analysis indicated that a significant minority of cells that were on the trajectory to the MET region continues to switch to the trajectory to the Stromal Region (FIG. 25G).

[0264]    Regulatory analysis identified TFs associated with the two trajectories. Three TFs (Dmrtc2, Zic3, and Pou3fl) were induced in all cells (from undetectable levels at day 0), but showed higher expression along the trajectory to the MET Region (FIG. 25E, 25F). Zic3 was required for maintenance of pluripotency (Lim et al., 2007), Pou3fl was required for self-renewal of spermatogonial stem cells (Wu et al., 2010), and Dmrtc2 was involved in germ cell development (Gegenschatz-Schmid et al., 2017; Yamamizu et al., 2016). Four TFs (Id3, Nfix, Nfic, and Prrxl) were upregulated in all cells (from basal levels at day 0) but showed higher expression in cells with a stromal fate (FIGs. 25E, 25F). (Analysis of subsequent time points showed that, following withdrawal of dox, these genes maintained high expression in stromal cells but shut off in cells along the trajectory to iPSCs.) Nfix was reported to repress embryonic expression programs in early development, while Nfic and Prrxl were associated with mesenchymal programs (Froidure et al., 2016; Messina et al., 2010; Ocana et al., 2012). Id3 was known to inhibit transcription through formation of nonfunctional dimers that were incapable of binding to DNA. Higher expression of Id3 along the trajectory toward stromal cells may seem

somewhat surprising, because forced expression of Id3 was shown to increase reprogramming efficiency (Hayashi et al., 2016; Liu et al., 2015). However, Id3 might cause increased efficiency via its activity in stromal cells, which secreted factors that enhance iPSC reprogramming (Mosteiro et al., 2016) (see below), or via activity in non-stromal cells, in which it was expressed through day 8, albeit at lower levels.

[0265]    There has been much interest in finding early markers of successful reprogramming—namely, genes whose early expression was correlated with a cell's descendants being enriched for iPSCs. Our analysis suggested that it would be more precise to define "early markers of successful MET", because the iPSC, trophoblast and neural fates did not appear to be established until after withdrawal of dox at day 8.

[0266]    Trajectory analysis revealed early markers of successful MET, including known markers such as Fut9 (which synthesizes the glyco-antigen SSEA-1) and novel candidates such as Shisa8. Shisa8 was the most differentially expressed gene at day 1.5. When we sorted cells based on the ratio of their likelihood of transition to the MET Region vs Stromal Region, we found Shisa8 expressed in 50% of the top quartile but only 5% of cells in the bottom quartile. (Table 16). Shisa8 was a little-studied mammalian-specific member of the Shisa gene family in vertebrates, which encoded single-transmembrane proteins that played roles in development and are thought to serve as adaptor proteins (Pei and Grishin, 2012; Polo et al., 2012). (Analysis of subsequent time points showed that Shisa8 and Fut9 also showed similar patterns following dox withdrawal: both were expressed strongly in cells along the trajectory toward successful reprogramming, and lowly expressed in other lineages (FIG. 3ID).)

Table 16 - Differential genes between top ancestors of MET vs. top ancestors of stromal cells.

| Differential genes between top ancestors of MET vs. Stromal cells at D1.5 | | | | | |
|---|---|---|---|---|---|
| Gene | p-value | Average logFC | Fraction expressed in top ancestors of MET | Fraction expressed in top ancestors of stromal cells | Adjusted p-value |
| Shisa8 | 2.37E-56 | 0.439583976 | 0.505 | 0.051 | 4.52E-52 |
| Anpep | 1.24E-44 | 0.399501581 | 0.548 | 0.141 | 2.37E-40 |
| Gch1 | 5.09E-37 | 0.381008072 | 0.607 | 0.245 | 9.71E-33 |
| Gpm6b | 1.24E-29 | 0.275486032 | 0.538 | 0.209 | 2.37E-25 |

| Npnt | 3.61E-30 | 0.382743398 | 0.714 | 0.395 | 6.89E-26 |
|---|---|---|---|---|---|
| Dsp | 9.36E-34 | 0.290320422 | 0.389 | 0.072 | 1.79E-29 |
| Rbl | 1.12E-25 | 0.280506707 | 0.616 | 0.315 | 2.13E-21 |
| Dgat2 | 5.18E-28 | 0.349298687 | 0.524 | 0.225 | 9.88E-24 |
| Carl2 | 1.06E-23 | 0.299588702 | 0.552 | 0.254 | 2.02E-19 |
| Lrp4 | 9.73E-27 | 0.247967802 | 0.405 | 0.11 | 1.86E-22 |
| Clql3 | 2.93E-26 | 0.325323868 | 0.45 | 0.155 | 5.60E-22 |
| Sgol2a | 1.65E-25 | 0.33023125 | 0.685 | 0.395 | 3.16E-21 |
| Gm26737 | 2.93E-25 | 0.534938533 | 0.656 | 0.368 | 5.59E-21 |
| Lepr | 1.15E-22 | 0.588193067 | 0.695 | 0.417 | 2.19E-18 |
| Nol4l | 1.78E-21 | 0.374175462 | 0.65 | 0.374 | 3.40E-17 |
| Gm29666 | 1.49E-20 | 0.279383915 | 0.511 | 0.237 | 2.84E-16 |
| Pfkp | 8.34E-30 | 0.316216243 | 0.796 | 0.524 | 1.59E-25 |
| RP23-4H17.3 | 4.98E-21 | 0.441940336 | 0.695 | 0.425 | 9.51E-17 |
| Ralgps2 | 4.40E-22 | 0.217741022 | 0.38 | 0.117 | 8.40E-18 |
| Xafl | 1.12E-18 | 0.328905337 | 0.564 | 0.307 | 2.14E-14 |
| Zdhhc2 | 2.08E-17 | 0.200585787 | 0.519 | 0.264 | 3.97E-13 |
| Ppmlk | 1.38E-22 | 0.307219164 | 0.658 | 0.411 | 2.63E-18 |
| McmlO | 1.99E-16 | 0.230302782 | 0.593 | 0.348 | 3.80E-12 |
| Gml3075 | 1.33E-27 | 0.861118262 | 0.771 | 0.528 | 2.53E-23 |
| Repl5 | 2.80E-18 | 0.29626083 | 0.658 | 0.423 | 5.34E-14 |
| Pola2 | 3.37E-23 | 0.311939681 | 0.748 | 0.519 | 6.44E-19 |
| Trim37 | 7.52E-17 | 0.218079056 | 0.583 | 0.358 | 1.44E-12 |
| Rtkn | 3.27E-18 | 0.287996995 | 0.382 | 0.16 | 6.24E-14 |
| Ppif | 1.58E-21 | 0.252798031 | 0.767 | 0.548 | 3.02E-17 |
| Rsfl | 2.84E-15 | 0.229977128 | 0.591 | 0.374 | 5.42E-11 |
| Ptcra | 5.85E-13 | 0.417578437 | 0.413 | 0.2 | 1.12E-08 |
| Nmrkl | 4.51E-13 | 0.528279491 | 0.554 | 0.344 | 8.61E-09 |
| Perp | 4.55E-65 | 0.656396496 | 0.963 | 0.753 | 8.69E-61 |
| Chmp2b | 1.29E-30 | 0.335057338 | 0.849 | 0.64 | 2.46E-26 |
| Pcgf2 | 5.58E-15 | 0.541239697 | 0.591 | 0.387 | 1.07E-10 |
| Gmcll | 4.30E-14 | 0.523834071 | 0.544 | 0.344 | 8.21E-10 |
| Pacsl | 1.50E-18 | 0.251074727 | 0.785 | 0.587 | 2.87E-14 |
| Wdr35 | 3.75E-14 | 0.224471336 | 0.656 | 0.464 | 7.15E-10 |
| Ppat | 2.16E-16 | 0.243243284 | 0.708 | 0.517 | 4.13E-12 |
| Slamfl | 5.19E-11 | 0.228267013 | 0.468 | 0.28 | 9.90E-07 |
| Homer2 | 6.66E-14 | 0.236094482 | 0.624 | 0.438 | 1.27E-09 |

| | | | | | |
|---|---|---|---|---|---|
| Cenph | 7.86E-14 | 0.206088745 | 0.72 | 0.538 | 1.50E-09 |
| B930036N1 0Rik | 2.34E-10 | 0.518225771 | 0.544 | 0.368 | 4.46E-06 |
| Hpcall | 8.65E-13 | 0.208476389 | 0.613 | 0.438 | 1.65E-08 |
| H2-T23 | 8.64E-11 | 0.235054556 | 0.337 | 0.164 | 1.65E-06 |
| Sgoll | 2.01E-16 | 0.266408936 | 0.853 | 0.683 | 3.83E-12 |
| Ccdcl37 | 2.58E-20 | 0.287870449 | 0.793 | 0.624 | 4.93E-16 |
| Exosc2 | 9.42E-37 | 0.652481854 | 0.933 | 0.765 | 1.80E-32 |
| Gkapl | 1.74E-23 | 0.397791708 | 0.781 | 0.613 | 3.31E-19 |
| Agl | 1.58E-16 | 0.495744367 | 0.798 | 0.63 | 3.01E-12 |
| Ckap2 | 8.06E-12 | 0.205735226 | 0.796 | 0.632 | 1.54E-07 |
| Nt5dc3 | 1.29E-10 | 0.200909668 | 0.638 | 0.481 | 2.46E-06 |
| Tapbpl | 7.86E-09 | 0.226071905 | 0.315 | 0.164 | 0.000150089 |
| Shoc2 | 9.21E-15 | 0.231434184 | 0.751 | 0.601 | 1.76E-10 |
| Faap24 | 3.98E-11 | 0.2159197 | 0.642 | 0.495 | 7.60E-07 |
| Haus8 | 2.63E-16 | 0.634579918 | 0.744 | 0.599 | 5.01E-12 |
| Cenpf | 7.61E-11 | 0.214446511 | 0.908 | 0.763 | 1.45E-06 |
| Mrpsll | 3.66E-41 | 0.430516438 | 0.906 | 0.763 | 6.99E-37 |
| Aldh3al | 8.14E-08 | 0.221022512 | 0.456 | 0.313 | 0.001554728 |
| Gm7120 | 8.12E-08 | 0.306764672 | 0.311 | 0.168 | 0.001550761 |
| Lpgatl | 4.28E-16 | 0.244225687 | 0.806 | 0.665 | 8.17E-12 |
| Topbpl | 5.86E-12 | 0.224664357 | 0.734 | 0.593 | 1.12E-07 |
| Mrps6 | 3.39E-43 | 0.396132536 | 0.939 | 0.798 | 6.47E-39 |
| 1700047117 Rik2 | 5.69E-09 | 0.200128893 | 0.521 | 0.382 | 0.000108639 |
| Myc | 4.08E-26 | 0.347729368 | 0.898 | 0.763 | 7.80E-22 |
| TimmlO | 4.34E-14 | 0.223178202 | 0.845 | 0.71 | 8.28E-10 |
| Mrpl9 | 9.74E-09 | 0.222293218 | 0.503 | 0.368 | 0.000185972 |
| Famll4a2 | 2.19E-18 | 0.23879583 | 0.83 | 0.697 | 4.18E-14 |
| Rm3 | 1.49E-11 | 0.228168673 | 0.724 | 0.591 | 2.84E-07 |
| Dcafl7 | 2.63E-08 | 0.521823548 | 0.487 | 0.354 | 0.00050265 |
| Asph | 2.31E-14 | 0.224904909 | 0.787 | 0.656 | 4.42E-10 |
| Abcblb | 6.60E-40 | 0.441369564 | 0.947 | 0.818 | 1.26E-35 |
| Ctnnbll | 2.19E-11 | 0.207192935 | 0.777 | 0.648 | 4.18E-07 |
| Slbp | 1.84E-15 | 0.374861946 | 0.873 | 0.748 | 3.52E-11 |
| TexlO | 3.22E-15 | 0.251420666 | 0.8 | 0.677 | 6.14E-11 |
| Dennd5b | 3.94E-11 | 0.298384346 | 0.755 | 0.632 | 7.52E-07 |
| Lrrc42 | 3.19E-14 | 0.250507008 | 0.748 | 0.626 | 6.09E-10 |

| Paip2b | 6.60E-09 | 0.233070859 | 0.691 | 0.571 | 0.000126059 |
|---|---|---|---|---|---|
| 1700037H04Rik | 3.73E-13 | 0.21591323 | 0.777 | 0.663 | 7.12E-09 |
| Noal | 1.13E-34 | 0.490924229 | 0.9 | 0.787 | 2.17E-30 |
| Gtf2hl | 5.71E-19 | 0.253937461 | 0.843 | 0.738 | 1.09E-14 |
| Ndcl | 4.28E-18 | 0.25208573 | 0.89 | 0.785 | 8.16E-14 |
| Ddx42 | 1.64E-13 | 0.213024231 | 0.83 | 0.726 | 3.13E-09 |
| Golga3 | 9.43E-07 | 0.495832978 | 0.595 | 0.491 | 0.018003133 |
| Pop5 | 1.28E-28 | 0.301595886 | 0.949 | 0.847 | 2.44E-24 |
| Tgfbi | 1.63E-09 | 0.200070657 | 0.828 | 0.726 | 3.11E-05 |
| Hells | 3.70E-13 | 0.222587886 | 0.949 | 0.851 | 7.06E-09 |
| Plk4 | 1.42E-23 | 0.57479234 | 0.922 | 0.826 | 2.72E-19 |
| Ezh2 | 1.90E-18 | 0.236909466 | 0.906 | 0.81 | 3.64E-14 |
| Naa20 | 8.41E-18 | 0.270587809 | 0.806 | 0.714 | 1.61E-13 |
| Epnl | 1.54E-14 | 0.209191303 | 0.902 | 0.812 | 2.94E-10 |
| Smnl | 9.92E-38 | 0.401700379 | 0.941 | 0.853 | 1.89E-33 |
| Mcm7 | 1.42E-16 | 0.229113377 | 0.955 | 0.867 | 2.72E-12 |
| Enah | 1.19E-12 | 0.207086155 | 0.828 | 0.742 | 2.27E-08 |
| Mrps25 | 2.24E-16 | 0.238478878 | 0.863 | 0.783 | 4.27E-12 |
| Carnmtl | 7.08E-15 | 0.213768504 | 0.871 | 0.791 | 1.35E-10 |
| Zfpl06 | 4.55E-12 | 0.206955912 | 0.943 | 0.863 | 8.69E-08 |
| Hmgb3 | 4.37E-16 | 0.244565953 | 0.879 | 0.802 | 8.34E-12 |
| PsmblO | 8.45E-25 | 0.305887579 | 0.937 | 0.861 | 1.61E-20 |
| Scp2 | 7.16E-12 | 0.211532788 | 0.883 | 0.808 | 1.37E-07 |
| Histlh2ap | 1.60E-27 | 0.599321987 | 0.978 | 0.904 | 3.05E-23 |
| Limk2 | 1.79E-12 | 0.34639987 | 0.81 | 0.738 | 3.42E-08 |
| Dbf4 | 5.21E-15 | 0.209332579 | 0.922 | 0.851 | 9.95E-11 |
| Bazla | 2.09E-20 | 0.276857187 | 0.881 | 0.812 | 4.00E-16 |
| Ifrd2 | 4.47E-21 | 0.25780276 | 0.908 | 0.84 | 8.53E-17 |
| Ccdc50 | 1.00E-25 | 0.293196782 | 0.955 | 0.888 | 1.92E-21 |
| Pbdcl | 3.94E-14 | 0.228782894 | 0.875 | 0.808 | 7.52E-10 |
| Wdr45b | 8.91E-11 | 0.203638926 | 0.832 | 0.769 | 1.70E-06 |
| Noc2l | 8.02E-21 | 0.235002625 | 0.951 | 0.89 | 1.53E-16 |
| Ruvbll | 3.88E-11 | 0.20097654 | 0.828 | 0.767 | 7.41E-07 |
| Prmt5 | 1.96E-13 | 0.20762784 | 0.888 | 0.832 | 3.74E-09 |
| Tmem245 | 1.26E-32 | 0.731436804 | 0.963 | 0.908 | 2.40E-28 |
| Pnol | 1.18E-22 | 0.284205102 | 0.894 | 0.84 | 2.25E-18 |
| Chchd7 | 1.97E-33 | 0.376522958 | 0.92 | 0.867 | 3.76E-29 |

| Yif1b | 2.51E-12 | 0.204286063 | 0.91 | 0.857 | 4.80E-08 |
|---|---|---|---|---|---|
| Nip7 | 1.61E-09 | 0.317643192 | 0.896 | 0.843 | 3.07E-05 |
| Stmn1 | 7.91E-13 | 0.214767905 | 0.926 | 0.875 | 1.51E-08 |
| Rtcb | 3.23E-21 | 0.248019171 | 0.933 | 0.885 | 6.16E-17 |
| Nmt2 | 9.69E-54 | 0.59549564 | 0.988 | 0.941 | 1.85E-49 |
| Fnta | 2.30E-11 | 0.208830016 | 0.824 | 0.779 | 4.40E-07 |
| Snhg9 | 4.41E-41 | 0.578853339 | 0.971 | 0.928 | 8.42E-37 |
| Tax1bp1 | 1.04E-11 | 0.20563376 | 0.855 | 0.812 | 1.98E-07 |
| Cdk6 | 9.45E-13 | 0.216050004 | 0.935 | 0.896 | 1.80E-08 |
| Tcof1 | 3.45E-31 | 0.302647593 | 0.965 | 0.928 | 6.58E-27 |
| Cebpz | 1.09E-16 | 0.237798069 | 0.939 | 0.902 | 2.09E-12 |
| Loxl2 | 1.30E-17 | 0.571139295 | 0.89 | 0.857 | 2.48E-13 |
| Rangap1 | 2.34E-40 | 0.369409656 | 0.984 | 0.953 | 4.46E-36 |
| Dek | 1.64E-18 | 0.231074803 | 0.996 | 0.967 | 3.12E-14 |
| Nolc1 | 9.61E-30 | 0.309060428 | 0.986 | 0.959 | 1.83E-25 |
| Mybbp1a | 1.01E-15 | 0.209760443 | 0.969 | 0.943 | 1.92E-11 |
| Uchl3 | 4.63E-23 | 0.291386824 | 0.963 | 0.937 | 8.83E-19 |
| Mt2 | 2.21E-46 | 0.647830277 | 0.982 | 0.959 | 4.21E-42 |
| Fam177a | 7.40E-29 | 0.318947806 | 0.965 | 0.943 | 1.41E-24 |
| Ak2 | 2.85E-38 | 0.322110667 | 0.992 | 0.971 | 5.45E-34 |
| Pdcd11 | 1.06E-26 | 0.317776644 | 0.994 | 0.973 | 2.03E-22 |
| Clns1a | 7.78E-15 | 0.200963226 | 0.955 | 0.935 | 1.49E-10 |
| Nsun2 | 4.46E-23 | 0.25780744 | 0.965 | 0.947 | 8.51E-19 |
| Eif1ax | 6.10E-25 | 0.259171146 | 0.998 | 0.982 | 1.17E-20 |
| Utp11l | 2.11E-21 | 0.247732591 | 0.978 | 0.963 | 4.03E-17 |
| Nifk | 4.74E-16 | 0.25794523 | 0.973 | 0.959 | 9.06E-12 |
| Mrpl36 | 8.39E-15 | 0.203735334 | 0.963 | 0.949 | 1.60E-10 |
| Chchd4 | 3.75E-49 | 0.406592072 | 0.99 | 0.978 | 7.15E-45 |
| Mt1 | 1.69E-19 | 0.330543022 | 0.99 | 0.98 | 3.23E-15 |
| Mcm6 | 5.05E-14 | 0.203330997 | 0.93 | 0.92 | 9.64E-10 |
| 2810004N2 3Rik | 2.73E-25 | 0.282539829 | 0.982 | 0.973 | 5.21E-21 |
| Lmo4 | 1.74E-66 | 0.775349512 | 0.992 | 0.986 | 3.31E-62 |
| Sms | 1.65E-36 | 0.313663566 | 0.992 | 0.986 | 3.15E-32 |
| Tmem5 | 7.44E-27 | 0.31509393 | 0.949 | 0.943 | 1.42E-22 |
| Abcf1 | 4.64E-25 | 0.277959491 | 0.992 | 0.988 | 8.85E-21 |
| Sfxn1 | 6.98E-21 | 0.212944289 | 0.984 | 0.98 | 1.33E-16 |
| Gml6286 | 8.21E-20 | 0.224472114 | 0.988 | 0.984 | 1.57E-15 |

200

| Cox7a2l | 1.45E-19 | 0.200215258 | 0.994 | 0.99 | 2.77E-15 |
|---|---|---|---|---|---|
| Psatl | 2.81E-16 | 0.206124692 | 0.994 | 0.99 | 5.37E-12 |
| Zfosl | 5.30E-16 | 0.206256512 | 0.992 | 0.988 | 1.OlE-11 |
| Nhp2ll | 9.94E-34 | 0.239069695 | 1 | 0.998 | 1.90E-29 |
| Txn2 | 8.06E-23 | 0.202261807 | 0.994 | 0.992 | 1.54E-18 |
| Dctppl | 1.40E-22 | 0.221067567 | 0.992 | 0.99 | 2.67E-18 |
| Eif3jl | 8.55E-20 | 0.270419381 | 0.992 | 0.99 | 1.63E-15 |
| Nhp2 | 3.24E-68 | 0.348934627 | 1 | 1 | 6.19E-64 |
| Txnl4a | 6.38E-49 | 0.36485702 | 0.99 | 0.99 | 1.22E-44 |
| Naplll | 1.10E-46 | 0.276547552 | 1 | 1 | 2.10E-42 |
| Srm | 1.22E-45 | 0.356879476 | 0.992 | 0.992 | 2.32E-41 |
| Tomm5 | 1.65E-43 | 0.313429107 | 1 | 1 | 3.15E-39 |
| Dnajc2 | 4.24E-40 | 0.373302174 | 0.988 | 0.988 | 8.10E-36 |
| Ddx21 | 2.72E-35 | 0.383841731 | 0.996 | 0.996 | 5.18E-31 |
| Ncl | 6.24E-31 | 0.351868277 | 1 | 1 | 1.19E-26 |
| Serbpl | 1.10E-27 | 0.22648657 | 1 | 1 | 2.11E-23 |
| Naal5 | 1.44E-20 | 0.281257486 | 0.982 | 0.982 | 2.75E-16 |
| Maplb | 1.99E-11 | 0.211674236 | 0.949 | 0.949 | 3.79E-07 |
| Gngl2 | 3.44E-45 | 0.336166251 | 0.994 | 0.996 | 6.58E-41 |
| Bola2 | 1.95E-33 | 0.243627002 | 0.998 | 1 | 3.72E-29 |
| Ddxl8 | 1.13E-20 | 0.236133065 | 0.994 | 0.996 | 2.15E-16 |
| Calml | 4.37E-20 | 0.209338392 | 0.998 | 1 | 8.35E-16 |
| Llph | 2.37E-16 | 0.207946587 | 0.994 | 0.996 | 4.52E-12 |
| Hnrnpm | 1.63E-15 | 0.211499543 | 0.99 | 0.992 | 3.11E-11 |
| NoplO | 2.74E-32 | 0.258763009 | 0.996 | 1 | 5.23E-28 |
| Wdr43 | 1.46E-25 | 0.286052346 | 0.992 | 0.996 | 2.80E-21 |
| mt-Nd3 | 2.70E-23 | 0.241501548 | 0.994 | 0.998 | 5.15E-19 |
| Knopl | 1.42E-22 | 0.257948217 | 0.992 | 0.996 | 2.71E-18 |
| Dpy30 | 1.40E-15 | 0.206386698 | 0.971 | 0.975 | 2.67E-11 |
| Dph3 | 1.25E-33 | 0.288444631 | 0.982 | 0.988 | 2.38E-29 |
| Anp32b | 6.68E-20 | 0.23155113 | 0.99 | 0.996 | 1.28E-15 |
| Odcl | 2.58E-14 | 0.212362532 | 0.988 | 0.996 | 4.92E-10 |

[0267]  iPSCs emerge through a tight bottleneck from cells in the MET Region

[0268]  Trajectory analysis showed that cells from the MET region subsequently gained a broad epithelial identity and began to rapidly diverge to give rise the iPS-, epithelial-, trophoblast-, and neural-like cells (FIG. 26A). Importantly, the ancestor distributions of these

classes were not distinguishable before the withdrawal of dox at day 8, suggesting that the cells' fates did not appear yet to be determined at that point (FIG. 26B).

[0269] By day 11.5-12.5, the iPS-like cells began to show a clear signature of pluripotency, including canonical marker genes such as Nanog, Sox2, Zfp42, Otx2, Dppa4, and an elevated cell-cycle signature (FIGs. 26C, 26D). In 2i conditions, these iPS-like cells accounted for 12% of cells by day 11.5 and 80-90% from days 15 through 18. In serum conditions, the trend was similar, but the process was delayed by roughly one day and was far less efficient: the pluripotency signature was found in 3.5% of cells by day 12.5 and peaked at just 10-15% from days 15.5 through 18 (FIG. 24G). Notably, we found substantial heterogeneity among the iPSC-related cells. Recent studies reported that a small subset of cells in 2i conditions showed a signature characteristic of the embryonic 2-cell (2C) stage (Falco et al., 2007; Kolodziejczyk et al., 2015; Macfarlan et al., 2012). Scoring our iPS-like cells with signatures based on profiles from 2 cell-, 4 cell-, 8 cell-, 16 cell-, and 32 cell-stage embryos (Goolam et al., 2016) (Table 15, FIG. 32A, 32B), -20% of cells in both 2i and serum conditions showed a 2C, 4C, 8C, 16C, or 32C signature (with roughly half showing signatures for two consecutive stages).

[0270] Trajectory analysis suggested that successfully reprogrammed cells passed through a tight bottleneck in days 10-1 1. The ancestral distribution of iPSCs spanned -40% of all cells at day 8.5. It falls to -10% of cells at day 10 in 2i conditions and only -1% at day 11 in serum conditions. These results suggested that only a small and distinct subset of cells transitioning out of the MET Regions toward various fates had the potential to become iPS cells (below). These iPSC progenitors did not yet fully acquired the pluripotency signature but were changing rapidly toward this fate. They resided along certain thin 'strings' in the FLE representation (FIG. 24F, white arrow and 4C, green). iPSC ancestors then rose to -40% at day 14 in 2i (and 10% on day 14 in serum), reflecting rapid expansion of pluripotent precursors (FIG. 26C, yellow).

[0271] By clustering genes according to similar expression trends along the trajectories to successful reprogramming in 2i and serum conditions, we found induction of various groups of genes involved in regulation of pluripotency, and repression of genes involved in certain metabolic changes and RNA processing (FIG. 32C). Among the upregulated genes, 24 were preferentially expressed in the late stage of reprogramming on successful trajectories and were

mostly absent from other cell types; these included Ooep, Fmrlnb, Lncencl, and Tell (FIG. 32C, Table 17). These genes can be candidate markers for fully reprogrammed cells.

**Table 17** - List of genes for 15 groups of genes along the successfully reprogrammed trajectory reported in FIG. 32A

| Gene sets related to FIG. 32A | | | | | | | |
|---|---|---|---|---|---|---|---|
| 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 |
| Sbspon | Terf1 | Lypla1 | Lactb2 | Pnkd | Rpl7 | Tcea1 | Il1rl1 |
| Dst | 1700007K13Rik | Tceb1 | Igfbp2 | Ptma | Rpl31 | Mcm3 | Fhl2 |
| Nrp2 | Ass1 | Dnpep | Trip12 | Dtymk | H3f3a | Sgol2a | Col3a1 |
| Eef1b2 | Mdk | Tfcp2l1 | Marc2 | Dbi | Rpl7a | Psmd1 | Col5a2 |
| Serpine2 | Chchd5 | Kdm5b | Gm13580 | Snrpe | Rpl12 | R3hdm1 | Sdpr |
| Ephx1 | Praf2 | Swt1 | Hat1 | Cacybp | Zfos1 | Mcm6 | Fn1 |
| Nudt5 | Timm17b | Atp1b1 | Tfpi | Ndufs2 | Pcsk1n | Dhx9 | Col6a3 |
| Commd3 | Hdac6 | Phyh | Platr3 | F11r | Rpl10 | Gm2000 | Gpc1 |
| Ndufa8 | Ndufb11 | Wdr5 | Scand1 | Atp5c1 | Bex2 | Prrc2c | Serpinb2 |
| Ccdc34 | Uxt | Odf2 | Platr27 | Tubb4b | Ndufb5 | Parp1 | Ubxn4 |
| Nop10 | Klhl13 | Rif1 | Fthl17c | Spc25 | Rps3a1 | Nvl | Klhdc8a |
| Knstrn | Slc25a5 | AA467197 | Usp9x | 2700094K13Rik | Apoa1bp | Lbr | Ptgs2 |
| Dtd1 | Ube2a | Slc24a5 | Ndufa1 | Cd59a | Txnip | Enah | Rgs16 |
| Rbck1 | Upf3b | Mrps5 | Gm9 | Eif3m | Gstm1 | Cenpf | Ier5 |
| Nnat | Rhox6 | Eif2s2 | Rhox1 | Rad51 | Rpl34 | Dtl | Soat1 |
| Rbm3 | Rhox9 | Mybl2 | Rhox5 | Spint1 | Rps20 | Yme1l1 | Copa |
| Hmgb3 | Mcts1 | Gtsf1l | Thoc2 | Hypk | Gm11808 | Set | Grem2 |
| Fundc2 | Bcap31 | Wfdc2 | Rbmx2 | Dut | Rps6 | Prrc2b | Col5a1 |
| Slc7a3 | Idh3g | Ncoa3 | Usp26 | 1700037H04Rik | Rps8 | Rpl35 | Angptl2 |
| Hmgn5 | Lage3 | Sall4 | Hprt | Tpx2 | Laptm5 | Hnrnpa3 | Hspa5 |
| 2210013O21Rik | Pbdc1 | Tfap2c | 1700013H16Rik | Ube2c | Rpl11 | Nusap1 | Gorasp2 |
| Rnf13 | Bex4 | Ebp | Fmr1nb | Aurka | Rpl22 | Mga | Creb3l1 |
| Cks1b | Bex1 | Atp6ap2 | Dusp9 | Ppdpf | Rpl9 | Zfp106 | Rcn1 |
| Psmb4 | Wbp5 | Nono | Ssr4 | Plp2 | Rpl5 | Myef2 | Bdnf |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Bolal | Ngfrapl | Algl3 | Dkcl | NaalO | Rpl21 | Xrn2 | Thbsl |
| Gstm5 | Trapla | Gm8797 | Vbpl | Pdhal | Gapdh | Csnk2al | Fgf7 |
| Psrcl | Hsdl7bl0 | Tpd52 | Pdk3 | Exosc8 | Rps9 | Ubal | Dstn |
| Cth | Rab9 | Chmp4c | Lasll | Smc4 | Cox6b2 | Gnl3l | Rrbpl |
| Ndufb6 | Dnajcl9 | Lrrc31 | Ogt | Pmfl | Rpl28 | Huwel | Thbd |
| Cdc26 | Lamtor2 | Actl6a | Pin4 | Rab25 | Rps5 | Smcla | Srxnl |
| Psipl | Fdps | Fxrl | Atrx | Anp32e | Rpsl9 | Sms | Chmp4b |
| Cdkn2a | Psmd4 | Sox2 | Magtl | Atp5fl | Rpsl6 | Midi | Procr |
| Lltdl | Acp6 | Noct | Cox7b | Stoml2 | Eif3k | 1810022K09Rik | Dlgap4 |
| Tmem59 | Hadh | PlatrlO | Pgkl | Ctnnall | Spint2 | Ndufcl | Ptpnl |
| Hspbll | Acer2 | Hiatl | Rpl36a | Nasp | Cox6bl | Slc39al | Pmepa1 |
| Uqcrh | Slc2al | Elovl6 | Prpsl | Cdc20 | Rpll3a | Ilf2 | Slco4a1 |
| Ptprf | Gjb5 | Acadm | Fgdl | Ppih | Rpll8 | Larp7 | Pgrmc1 |
| Eif3i | Hdacl | Zfp292 | Prdx4 | Cdca8 | Idh2 | Tet2 | Bgn |
| Atpifl | Hscb | Aqp3 | A830080D0IRik | Zbtb8os | Rps3 | Fubpl | Itm2a |
| Stmnl | Ung | Klf4 | Rbbp7 | Rpa2 | Rpl27a | Anp32b | Fndc3b |
| Enol | Cldn4 | Echdc2 | Zrsr2 | Hmgn2 | Rpsl3 | Smc2 | Sec62 |
| Fgfbpl | Cldn3 | Gjb3 | Ttcl4 | Miip | Rpsl5a | Zfp462 | Postn |
| Shisa3 | Atp6vlf | Fabp3 | Jadel | Apitdl | Uqcrc2 | Puml | Faml98b |
| Scarb2 | Mkrnl | Rps6kal | Vangll | Park7 | Ypel3 | Srrml | SlOOa7a |
| Cops4 | Cct7 | Rsrpl | Ak4 | Tyms | Ifitm3 | Rcc2 | Crctl |
| Gltp | Nful | Tcea3 | Fbliml | Cenpa | Rplp2 | Gm26825 | Ngf |
| Pop5 | Slc2a3 | Usp48 | Zfp600 | Qdpr | Mrpl23 | Tomm7 | Rhoc |
| Pebpl | Fkbp4 | Alpl | Gml3251 | Med28 | Rpsl2 | 4930548H24Rik | Csfl |
| Rpl6 | Ldhb | Gml3154 | 2610305D13Rik | Paics | Rpsl5 | Rfcl | Collla1 |
| Ran | AU018091 | Agtrap | Fbxo6 | G3bp2 | Rpl6l | Grsfl | F3 |
| Mospd3 | Ligl | Insigl | Rbpj | Hnrnpdl | Naca | Hnrnpd | Ostc |
| Hmgbl | Beam | Dnajb6 | Crlf2 | Cit | Rps26 | Golga3 | Cyr61 |
| Ndufa4 | Exosc5 | Yesl | Ppplcc | Rfc5 | Ndufal3 | Mcm7 | Bel 10 |

| Podxl | Gmfg | Lap3 | Arf5 | Chchd2 | Rpll8a | Luc7l2 | Glipr2 |
|-------|------|------|------|--------|--------|--------|--------|
| Akrlb3 | Map4kl | Kit | Stra8 | Rfc2 | Bst2 | Cbx3 | Sec61b |
| Hnrnpa2bl | Ppplrl4a | Rest | Ube2s | Atp5j2 | Cox4il | Immt | Tnc |
| Lsm3 | Tbcb | Sppl | Zfp787 | Lsm5 | Rpll3 | TmsblO | Eval b |
| Trh | Gpil | Mtf2 | Tmeml60 | Tcf7ll | Rpll5 | Dqxl | Errfil |
| Mgstl | Etfb | Pxmp2 | Calm3 | Suclgl | Rps24 | Mcm2 | Ost4 |
| Trappc6a | Ucp2 | Ulkl | Zfp428 | Tpil | Rpl23a-ps3 | Ptms | Ugdh |
| Dmrtc2 | Folrl | Medl3l | Plekha4 | Cdca3 | Rpll3-ps3 | Aebp2 | Apbb2 |
| Fbl | Mrpll7 | Tbx3 | Arrdc4 | Lockd | Rps25 | Fam60a | Igfbp7 |
| Krtdap | Arl6ipl | Sbnol | Eif3f | Peg3 | Fxyd6 | Trim28 | Cxcl5 |
| Prmtl | Aldoa | Cops6 | Septl | Gltscr2 | Rpll0-ps3 | Hnrnpl | Ppbp |
| Bax | Pycard | Slc25al3 | Ctbp2 | Sael | Rpl4 | Polr2i | Cxcl3 |
| Ldha | Bnip3 | Asns | Sycp3 | Lsr | Gsta4 | Sema4b | Cxcll |
| Tm2d3 | Utfl | Trim24 | Nudt4 | Ruvbl2 | Eeflal | Prcl | Cxcl2 |
| I7Rn6 | Ifitm2 | Zc3havl | Sap30 | Bcat2 | Rpl29 | Blm | Ereg |
| Ndufc2 | Cenpw | Ezh2 | Gm2694 | Snrpn | Rpsa | RP23-4H17.3 | U90926 |
| Ndufabl | Ddit4 | Tra2a | Fam25c | Coq7 | Rpll4 | Bclafl | Rsrc2 |
| Tmem219 | Cisdl | Gdf3 | Sapl8 | Plkl | Rps27a | Ptges3 | Denr |
| Vkorcl | Ddt | Dppa3 | Klf5 | Spnsl | Gnb2ll | Arglul | Ubc |
| Mki67 | ChchdlO | Nanog | Khdc3 | Dctppl | Rpl26 | Mcm5 | Serpin el |
| Glrx3 | Pfkl | Lpcat3 | Ooep | Fbxo5 | Rpl23 | Smarca5 | Pcolce |
| Cd81 | Polr2e | Cd9 | Higdla | Sf3b5 | Rpll9 | Cnotl | Kdelr2 |
| Perp | Gpx4 | 2810474019Rik | Mrps24 | Cdkl | Rpl27 | Rps26-psl | Cavl |
| Mif | Cirbp | Apocl | Eif4al | Lsm7 | Dcxr | Aars | Fine |
| Atp5d | 1500009L16Rik | Apoe | Clqbp | Eef2 | Rps23 | Ankrdll | Ptn |
| Ndufs7 | Priml | Pvrl2 | Suzl2 | Mrpl42 | Btf3 | Wapl | Capg |
| Uqcrll | Eif4ebpl | Cox7al | Al662270 | Cct2 | Rps7 | Rpgripl | Rab7 |
| Oazl | Ankrd37 | Tdrdl2 | Dynll2 | Atp5b | Wdr89 | Suptl6 | Fbln2 |
| Slc25a3 | Cope | Tead2 | E130012A19Rik | Ormdl2 | Rpl30 | Zc3hl3 | Sec 13 |
| Ndufal2 | Sin3b | Gtf2hl | Gnal3 | Sarnp | Gml0020 | Uchl3 | Cxcll 2 |
| Cnpy2 | Syce2 | Spty2dl | Snhg20 | Hmgb2 | Rpl8 | Anapcl3 | Tspan9 |

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| Nabp2 | Asnal | Mfge8 | Texl9.1 | Lsm4 | Rpl3 | Gnai2 | Arhgdib |
| Slc25a4 | Mtl | Ticrr | Pfkp | Tecr | Rpl35a | UqcrlO | 1111 |
| Apela | 2700060E02Rik | Zfand6 | Tubb2b | Orc6 | Gm9843 | Actr2 | Ehd2 |
| Isynal | Mrpsl6 | Eed | H2afy | Nudt21 | Sodl | Canx | Pvr |
| Mrpl34 | Tkt | Tmem41b | Cox7c | Cdhl | Psmbl | Alkbh5 | Plaur |
| Ndufb7 | Mphosph8 | Gga2 | Lncencl | Psmb5 | NdufblO | Ncorl | Psmd8 |
| Prdx2 | Esco2 | Nfatc2ip | Nampt | Dhrs4 | RpslO | Pfas | Fxyd5 |
| Pllp | Bnip3l | Mylpf | Ifi27 | Cdca2 | RpllOa | Naa38 | Rcn3 |
| Got2 | Sugtl | Echsl | Tell | Spc24 | Ddah2 | Xafl | Klfl3 |
| PsmblO | Pigyl | Ifitml | Papola | H2afx | Gm26917 | Ywhae | Vimp |
| Rab4a | Psma4 | Taldol | Apobec3 | Slc35f2 | AY036118 | Tafl5 | Lrrc32 |
| Dnajc9 | Cox5a | Fgf4 | Smclb | Pkm | 2410015M2ORik | Npepps | Map6 |
| Itm2b | Morf4ll | Akapl2 | Pim3 | Anp32a | Rpl27-ps3 | Top2a | Adm |
| Atp5l | H2afv | Sgkl | Rpl39l | Snapc5 | Gml0036 | Acly | Mical2 |
| Cadml | Commdl | Tetl | Eif4a2 | Tipin | Prelid2 | Bptf | Tgfbli1 |
| Crabpl | Pttgl | Spic | Adprh | Ccnb2 | Rpsl4 | Fasn | Rnhl |
| 2810417H13Rik | Psmb6 | Csrp2 | Dppa4 | Cox7a2 | Rpll7 | Slcl6a3 | H19 |
| Rps27l | Psmdl2 | Baz2a | Dppa2 | Gpxl | Gm6133 | Dek | Igf2 |
| Gtf2a2 | Atp5h | Ash2l | Cggbpl | Impdh2 | Fau | Rbm25 | Cttn |
| Hmgn3 | Galkl | Zfp42 | Morc3 | Ndufaf3 | Cox8a | Dnajc21 | Rgsl7 |
| Nf2 | Psma2 | Tmeml92 | Brwdl | Uqcrcl | Eeflg | MyolO | Ctgf |
| Ramp3 | Acotl3 | Nr2c2ap | Tmeml81a | Zmat5 | Gm9493 | Rad21 | Sarla |
| Mdhl | Uqcrb | Klf2 | Dynltla | Pold2 | Rpl9-ps6 | Stl3 | Col6a2 |
| Hintl | Cetn3 | AnapclO | Mpcl | Snrnp25 | Gstol | Limal | Pofut2 |
| Aldh3al | Dhfr | Dnase2a | Pgp | Npml | Rpsl2-ps3 | Usp7 | Pttglip |
| Poldip2 | Mycn | Mt2 | Gfer | Hmmr | mt-Co2 | Etv5 | Bsg |
| Krtl9 | Psma6 | Gabarapl2 | Piml | Cdkn2aipnl | | Tfrc | Timp3 |
| Krtl7 | Fkbp3 | Kat6b | Myolf | Tmeml07 | | Gsk3b | Btgl |
| Itgb4 | Atp6vld | Hesxl | Dhxl6 | Cldn7 | | Coxl7 | Atp2bl |
| Secl4ll | Brixl | Zfhx2 | Dazl | Atp5gl | | Gm8186 | Raplb |
| Tkl | Cox6c | Rnaseh2b | Vapa | Cbxl | | Srpkl | Ndufa4 12 |
| Stard3nl | Eif3e | Tdh | Ralbpl | Psmb3 | | Stk38 | Myl6 |

| Histlhlb | Tonsl | Rgcc | Arll4epl | Jup | | Brd4 | Hmoxl |
|---|---|---|---|---|---|---|---|
| Histlhle | Gcat | Zbtb44 | Prrcl | Dcakd | | Gm42418 | Junb |
| Uqcrfsl | Syngrl | Rpp25 | Fbxol5 | Sumo2 | | Uhrfl | Mmp2 |
| Eci2 | Cenpm | Rbpms2 | Gstp2 | Birc5 | | Khsrp | Gm22 |
| Ndufs6 | Ndufa6 | U2surp | D030056L22Rik | Stral3 | | Birc6 | Actal |
| Mrps36 | Atp5g2 | Slc25a36 | | Histlh2ae | | Erdrl | Nrpl |
| Id2 | Paml6 | Amt | | Gmnn | | Matr3 | Vcl |
| Rtnl | Pigx | Arih2 | | Cks2 | | Stipl | Arf4 |
| Sival | Ndufb4 | Slc25a20 | | Higd2a | | Incenp | Selk |
| Ahnak2 | Dynltlf | Tdgfl | | Ccnbl | | Tmem258 | Mustnl |
| Nudtl4 | Thoc6 | Trim71 | | Rrm2 | | Hells | Spcsl |
| Crip2 | Tceb2 | Uppl | | Misl8bpl | | Scd2 | Fermt2 |
| Ptp4a3 | Ccnf | Cct4 | | Mthfdl | | Eif3a | Gjb2 |
| Ly6a | Ndufv3 | Skpla | | Cct5 | | mt-Ndl | Ubl5 |
| Eefld | Ndufa7 | Vdacl | | Cycl | | | Col5a3 |
| Tst | Tubb5 | Gm2a | | Eif3l | | | Cnnl |
| HlfO | Rpp21 | Mpdul | | Tubalb | | | Oaf |
| Pmml | Znrdl | Tmem256 | | Krt8 | | | Thyl |
| Samm50 | Oardl | Scpepl | | Hnrnpal | | | Trappc4 |
| Eif4b | Ndufv2 | Igf2bpl | | Mrpl40 | | | Ncaml |
| 2610318N02Rik | Tgifl | Calcoco2 | | Rfc4 | | | Wdr61 |
| Dgcr6 | Cebpzos | Dnajc7 | | Bbx | | | Cspg4 |
| Fetub | Mta3 | Slc25a39 | | Ezr | | | Sema7a |
| Atp5o | Pfdnl | Grn | | Acat2 | | | Loxll |
| Agpat4 | Impa2 | Ccdc43 | | Cldn6 | | | Mapk6 |
| Nme4 | Smc3 | Ttyh2 | | Ppill | | | Col12a1 |
| Mapkl3 | | Wbp2 | | U2afl | | | Amotl2 |
| Cd320 | | Ubald2 | | Pfdn6 | | | Selm |
| Ly6g6c | | Jarid2 | | Lsm2 | | | Xbpl |
| Ly6g6f | | Ubxn2a | | Polrlc | | | Aebpl |
| Dnphl | | 1110008L16Rik | | Ndufall | | | Ykt6 |
| Cox7a2l | | Esrrb | | Crb3 | | | Tns3 |

| Pigf | | Ckb | | Myl12b | | | Sec61g |
|------|---|------|---|--------|---|---|--------|
| Ecscr | | Atxn10 | | Dpy30 | | | Sertad2 |
| Cyb5a | | Slc25a1 | | Epcam | | | Rtn4 |
| Rnaseh2c | | Morc1 | | Paip2 | | | Adam19 |
| Trmt112 | | Jam2 | | Lmnb1 | | | Sqstm1 |
| Carnmt1 | | Wtap | | Atp5a1 | | | Sparc |
| Avpi1 | | Sod2 | | Ndufs8 | | | Kctd11 |
| Ndufb8 | | Rnf5 | | Rbm4b | | | Gabarap |
| Cuedc2 | | Zfp57 | | Banf1 | | | Cxcl16 |
| Sfr1 | | Cdc5l | | Mrpl49 | | | Tax1bp3 |
| | | Slc29a1 | | Arl2 | | | Pafah1b1 |
| | | Gm7325 | | Fkbp2 | | | Serpinf1 |
| | | Ccnd3 | | | | | Ift20 |
| | | Ppm1b | | | | | Ccl2 |
| | | Msh2 | | | | | Ccl5 |
| | | Msh6 | | | | | Vmp1 |
| | | Cystm1 | | | | | Col1a1 |
| | | Taf7 | | | | | Copz2 |
| | | Dcp2 | | | | | Igfbp4 |
| | | Snx2 | | | | | Eif1 |
| | | Cndp2 | | | | | Timp2 |
| | | Chka | | | | | Klf6 |
| | | Ubxn1 | | | | | Inhba |
| | | Klf9 | | | | | Serpinb6a |
| | | Scd1 | | | | | Card19 |
| | | mt-Co1 | | | | | Pdlim7 |
| | | | | | | | Tmed9 |
| | | | | | | | Smim15 |
| | | | | | | | Plk2 |
| | | | | | | | Rhob |
| | | | | | | | Nfkbia |
| | | | | | | | Arf6 |

| | | | | | | | | Frmd6 |
|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Actn1 |
| | | | | | | | | Ltbp2 |
| | | | | | | | | Dlk1 |
| | | | | | | | | Tnfaip2 |
| | | | | | | | | Crip1 |
| | | | | | | | | Snhg18 |
| | | | | | | | | Cthrc1 |
| | | | | | | | | Ext1 |
| | | | | | | | | Has2 |
| | | | | | | | | Wisp1 |
| | | | | | | | | Myh9 |
| | | | | | | | | Lgals1 |
| | | | | | | | | Kdelr3 |
| | | | | | | | | Atf4 |
| | | | | | | | | Tuba1c |
| | | | | | | | | Itga5 |
| | | | | | | | | Vasn |
| | | | | | | | | Col8a1 |
| | | | | | | | | Ier3 |
| | | | | | | | | Ppp1r11 |
| | | | | | | | | Vegfa |
| | | | | | | | | Ltbp1 |
| | | | | | | | | Crim1 |
| | | | | | | | | Fez2 |
| | | | | | | | | Cdc42ep3 |
| | | | | | | | | Zfp36l2 |
| | | | | | | | | Hbegf |
| | | | | | | | | Yipf5 |
| | | | | | | | | Lox |
| | | | | | | | | Ier3ip1 |
| | | | | | | | | Efemp2 |
| | | | | | | | | Ehbp1l1 |

|  |  |  |  |  |  |  | Ehd1 |
|---|---|---|---|---|---|---|---|
|  |  |  |  |  |  |  | Fads3 |
|  |  |  |  |  |  |  | Ankrd1 |
|  |  |  |  |  |  |  | Dusp5 |

## Table 17 (Cont'd)

| Gene sets related to FIG. 32A | | | | | | |
|---|---|---|---|---|---|---|
| **9** | **10** | **11** | **12** | **13** | **14** | **15** |
| Map4k4 | Snhg6 | Ptp4a1 | Bag2 | Sdhaf4 | Imp4 | Eif5b |
| Bzw1 | Mpzl1 | Actr1b | Mrpl30 | Sumo1 | Tuba4a | Nop58 |
| Raph1 | Creg1 | Hspd1 | Hspe1 | Aamp | Ncl | Rpl37a |
| Arpc2 | Uap1l1 | Bok | Acadl | Eif4e2 | Ssna1 | Myeov2 |
| Tmbim1 | Ptges | Tsn | Stk16 | Timm17a | Surf2 | Sept2 |
| Lrrfip1 | Serf2 | Nucks1 | Adipor1 | Ufc1 | Urm1 | Ddx18 |
| Ube2f | Slc20a1 | Tpr | Phlda3 | Pfdn2 | Ppp2r4 | B930036N10 Rik |
| Hdlbp | Cst3 | Uck2 | Prdx6 | Hspa14 | Dpm2 | Nmt2 |
| Nifk | Gss | Hnrnpu | Mpc2 | Edf1 | Arpc5l | Sptan1 |
| Actr3 | Sdc4 | Eprs | Mgst3 | Dnlz | Timm10 | Exosc2 |
| Csrp1 | Adrm1 | Smyd2 | Cnih4 | 1110008P14 Rik | Ssrp1 | Dync1i2 |
| Arpc5 | Lamp2 | Rbm17 | Aida | Tor2a | Snrpb | Psmc3 |
| Qsox1 | Renbp | Agpat2 | St6galnac4 | Psmb7 | Ppid | Usp50 |
| Prrx1 | S100a1 | Fbxw2 | Pdia3 | Dnmt3b | Gpatch4 | Cse1l |
| Tmco1 | S100a13 | Mtx2 | Mrps26 | Tgif2 | Jtb | Atp5e |
| Tagln2 | Cnn3 | Caprin1 | Naa20 | Rpn2 | Nras | Rps21 |
| Wdr26 | Atp6v1g1 | 1500011K16 Rik | Fkbp1a | Tceal8 | Gar1 | Plk4 |
| Degs1 | Tm2d1 | Nop56 | Id1 | Morf4l2 | Cenpe | Naa15 |
| Capn2 | Atp6v0b | Snx5 | Dynlrb1 | Fabp5 | Sep15 | Rps27 |
| Rrp15 | Snhg12 | Raly | Romo1 | Car2 | Ebna1bp2 | Mrpl9 |
| Hacd1 | Sh3bgrl3 | 1110008F13 Rik | Samhd1 | Selt | Svbp | Sars |

| Surf4 | Pdpn | Srsf6 | Topi | Cct3 | Mrpsl5 | Agl |
|---|---|---|---|---|---|---|
| Ptrhl | Smiml4 | Sysl | Pfdn4 | Ssr2 | Thrap3 | Ccne2 |
| Faml29b | Coxl8 | Rael | Gnas | Rbm8a | Ak2 | Otud6b |
| Gsn | Hspb8 | Ddx3x | Ctsz | 1810037117 Rik | Tmem234 | Vcp |
| Rbmsl | Tmeml2 Oa | Vma21 | Slmo2 | Ube2d3 | Zcchcl7 | TexlO |
| Grbl4 | Arpclb | Ccna2 | Fhll | Dnajal | Hnrnpr | Tmem245 |
| Zak | Gpnmb | Tpm3 | G6pdx | Clta | Ddost | Lepr |
| Nfe2l2 | Malsul | Atplal | Xist | Prdxl | Mrto4 | Ccdcl63 |
| Nckapl | Pole4 | Csdel | Sh3bgrl | Psmb2 | Sdhb | Ybxl |
| Zc3hl5 | Chmp2a | Eif4e | Tmem35 | Marcksll | Szrdl | Gml3075 |
| Itgav | Vasp | Ddahl | Ammecrl | Trnaulap | Mrpl20 | Noc2l |
| Cd44 | Rabacl | Rad23b | Eiflax | Nude | Aurkaipl | Faml33b |
| Emc7 | Blvrb | Ndcl | Stmn2 | Sfn | Lrpapl | Abcblb |
| Eif3jl | Capnsl | Ctps | Lhfp | Tmem60 | Mrfapl | Dhxl5 |
| B2m | Dkkll | Pabpc4 | Tm4sfl | Ppplcb | Lyar | Noal |
| Fbnl | Nuprl | Mycbp | Mbnll | Slbp | Dynlll | Atp5k |
| Prnp | Snx3 | Sfpq | Lxn | Plac8 | Cox6al | Pdapl |
| H13 | Psap | Ptp4a2 | Hdgf | Anapc5 | Arl6ip4 | Ndufa5 |
| Pdrgl | Cstb | Ythdf2 | Mex3a | Por | Mrpsl7 | Rbm28 |
| Maprel | Gadd45b | Srm | S100al6 | Ywhag | Eif4h | Pdia4 |
| Eif6 | Aril | Gnbl | SlOOalO | Capza2 | Mdh2 | Serbpl |
| Myl9 | Ddit3 | Nadk | Mrps21 | Gstkl | Fisl | Hk2 |
| Ywhab | Cd63 | Dbf4 | Phgdh | Ruvbll | Znhitl | Paip2b |
| Timpl | Ifi30 | Dnajc2 | Camk2d | Arpc4 | Fscnl | Snrpg |
| Hs6st2 | Hsbpl | Abcf2 | Cisd2 | Hnrnpf | Arpcla | Gmcll |
| Flna | Mapllc3b | Rheb | Fam92a | M6pr | Pomp | Wbpll |
| Msn | Cyba | Ppmlg | Tmem55a | Mlf2 | 2610001J05Rik | Dennd5b |
| Satl | Tomm20 | Iscu | Ggh | Cops7a | Cycs | Ndufa3 |
| Sh3kbpl | Ghitm | Mlec | Tomm5 | Goltlb | Vamp8 | Cnot3 |
| Anxa5 | Psme2 | RnflO | Txnl | Clptml | Faml 36a | U2af2 |
| Ufml | Ctsb | Atp2a2 | Nfib | Psmc4 | Cnbp | Iqgapl |

| Dclkl | Srpr | Gnb2 | Scp2 | Nup62 | Hmces | Ipo7 |
|---|---|---|---|---|---|---|
| Wwtrl | Tbrgl | Eif3b | Ktil2 | Mesdc2 | Chchd4 | Teadl |
| Serpl | Hexa | Fam220a | Akrlal | Ppp4c | Emgl | 1110004F10 Rik |
| Ssr3 | Rablla | Cczl | Macfl | Bccip | Phb2 | Knopl |
| Crabp2 | Spg21 | Bri3 | Utplll | Phlda2 | Mrpl51 | Bola2 |
| Lmna | Ppib | Gtf3a | Wasf2 | Ltvl | Tsen34 | Fus |
| S100a4 | Rhoa | Hsphl | Mtfrll | Zwint | Napa | Hras |
| Sl00all | Pdlim4 | Mat2a | Id3 | Ube2n | Mrpsl2 | Polr2l |
| Vcaml | Cd68 | Mthfd2 | Hspg2 | Myl6b | Nudtl9 | Ap2a2 |
| Snx7 | Ggnbp2 | H2afj | Minosl | Fam32a | EmclO | Amdl |
| Ppp3ca | Nidi | Strap | Acot7 | Ddx39 | Grwdl | Ddx21 |
| Pdlim5 | Ninjl | Bcatl | Atad3a | Ier2 | Snrpal | Cdc34 |
| Lmo4 | Ctsl | Slcla5 | Cdk6 | Calr | Mrpsll | Metap2 |
| Sh3glbl | Gml0116 | Tomm40 | Sri | Cneplrl | Aen | PetlOO |
| Gng5 | Glrx | Eif4g2 | Mrpl33 | MM | Clnsla | Timm44 |
| Wis | Twistnb | Gdel | Grpell | Ciapinl | Tufm | Haus8 |
| Chchd7 | Npc2 | Mettl9 | Limchl | Gcsh | Ino80e | Gfod2 |
| Impadl | Dap | Eif3c | Ociadl | Emc8 | Bckdk | Nip7 |
| Rab2a | Ndrgl | Kcnqlotl | Ociad2 | Chmpla | Bub3 | 2810004N23 Rik |
| Ndufaf4 | Cyb5r3 | Rwddl | Septll | Gnpnatl | Urah | Gnl3 |
| Ube2jl | Tmbim6 | Ppal | Anxa3 | Bmp4 | Napll4 | Nisch |
| Tpm2 | Litaf | Mbd3 | Pdgfa | Dadl | Snrpd3 | Ktnl |
| Tlnl | Hacd2 | Abhdl7a | Racl | Tsc22dl | Sumo3 | Mrpl52 |
| Plin2 | Hcfclrl | Map2k2 | Kpna7 | Aasdhppt | Timml3 | Loxl2 |
| Mtap | Atp6v0e | Aes | Polrld | Rpusd4 | Thopl | Gml0076 |
| Jun | Ostfl | Rtcb | Shfml | Oaz2 | Dohh | Tafld |
| Jakl | Pdliml | Naplll | Lsm8 | Fam96a | Yeats4 | Gm26737 |
| Mast2 | | Cs | 1810058I24Rik | Rsl24dl | Cdk4 | Arppl9 |
| Elovll | | Dlcl | Gngl2 | Rnf7 | Pa2g4 | Rps27rt |
| Txlna | | Abcel | Aupl | Rbpl | Lsml | Limk2 |
| Clic4 | | Dnaja2 | Bola3 | Rrp9 | Fkbp8 | Nudcd3 |

| Cdc42 | | E2f4 | Actg2 | Nme6 | Ccdcl24 | Hnrnpab |
|---|---|---|---|---|---|---|
| Nppb | | Psmd7 | Arl6ip5 | Ewsrl | Ddal | Larpl |
| Pgd | | Dcunld5 | Foxpl | Arfl | Rbmxll | Mybbpla |
| Cgrefl | | Rp9 | Rhnol | Trp53 | Lsm6 | Ap2bl |
| Ywhah | | Ei24 | Magohb | Car4 | D8Ertd738e | Cite |
| Gml673 | | Rdx | Ybx3 | Slc35bl | 2310036022 Rik | Nfe2ll |
| Wdrl | | Imp3 | Epnl | H3f3b | Cmc2 | Pcgf2 |
| Pcdh7 | | Polr2m | Sepwl | Gaa | Aprt | Nmtl |
| Tpst2 | | Cdv3 | Gemin7 | Anapcll | Vdac2 | Ddx5 |
| Corolc | | Map4 | Egln2 | Dusll | Apexl | Rpl38 |
| Tmed2 | | G3bpl | Tmeml47 | Paklipl | Nedd8 | Srsf2 |
| Aplsl | | Srsfl | Pdcd5 | Emb | N6amt2 | Prpf4b |
| Fam20c | | Lrrc59 | Josd2 | Pdia6 | Reep4 | HnrnpaO |
| Actb | | Snf8 | Aktlsl | Ywhaq | Pinl | Nsa2 |
| Cyth3 | | Kpnbl | Igflr | Max | Tmedl | Smnl |
| Slc7al | | Psme3 | Serpinhl | Eif2sl | Ecsit | Rps29 |
| Colla2 | | Lsml2 | Rrml | Srsf5 | Elofl | Slirp |
| Tes | | Faml 04a | Prkcdbp | Ahsal | Hmbs | 2010107E04 Rik |
| Calu | | Prpsapl | Parva | Subl | Manf | Rpl37 |
| Caldl | | Gpsl | Tspan4 | Mcrsl | Tma7 | Wdr70 |
| Mtpn | | Gdi2 | Ccndl | Tarbp2 | Ccdcl2 | Polr2k |
| Zyx | | Rala | Epb41l2 | Copzl | Cld | Rangapl |
| Tex261 | | Ssrl | Marcks | Glyrl | Nhp2 | Hesl |
| Cyp26bl | | B230219D22 Rik | Cd24a | Ube2v2 | Uqcrq | Son |
| Sec61al | | Cxcll4 | Gjal | Ap2ml | Atoxl | Snhg9 |
| Brkl | | Hnrnpk | Arid5b | Dnajbll | Gukl | Hnrnpm |
| Ltbr | | Nsun2 | Plpp2 | Cct8 | Rangrf | Rps28 |
| Gabarapll | | RablO | Snrpf | Tcpl | Eif5a | Abcfl |
| Empl | | Smc6 | Atxn7l3b | Rabllb | Tmem97 | Ptcra |
| Erccl | | Odcl | Shmt2 | Mrpsl8b | Nmel | Sgoll |
| Cd3eap | | Srp54b | Lrpl | Meal | Mrpl27 | Wdr43 |
| Axl | | Glrx5 | Col4a2 | Calm2 | Phb | Cebpz |

| Actn4 | | Eif5 | Ckap2 | Polr2d | Coa3 | Epb41l4aos |
|---|---|---|---|---|---|---|
| 2200002D01 Rik | | Pabpcl | Vps36 | Eifla | Ictl | Ndufa2 |
| Atf5 | | Ly6e | Fgfrl | BC031181 | Hnl | Rbm22 |
| Emp3 | | Pcbp2 | Nrgl | Pgaml | Mrps7 | Tcofl |
| Prss23 | | Rslldl | Uba52 | Xpnpepl | 1810043H04 Rik | Nars |
| Rrp8 | | Gsptl | Pgls | mt-Co3 | Mrpll2 | Ddbl |
| Ilk | | Mapkl | Scoc | mt-Nd4 | Tmeml4c | Nmrkl |
| Rras2 | | Eif4gl | Nfix | | Nopl6 | Usmg5 |
| Pik3c2a | | Ppplr2 | Arl2bp | | Prelidl | Pdcdll |
| Itpripl2 | | 0610012G03 Rik | Gml0073 | | Lman2 | mt-Nd2 |
| Tnrc6a | | Naa50 | Zfhx3 | | Ddx46 | mt-Atp8 |
| Cdipt | | Tomm70a | 2310022B05 Rik | | 2010111I01Rik | mt-Nd3 |
| Abracl | | Srrm2 | Ube2el | | Mrpl36 | mt-Nd4l |
| Col6al | | Kif5b | Dph3 | | Sf3b6 | mt-Nd5 |
| Slcl9al | | Etfl | Anxa8 | | Sptssa | |
| Ube2g2 | | Hspa9 | Cnihl | | Erh | |
| Cnn2 | | Ube2d2a | Lgals3 | | TmedlO | |
| Nfic | | Psatl | Tptl | | Snwl | |
| Ncln | | Npm3 | Mbnl2 | | Zfp706 | |
| Txnrdl | | | Smco4 | | 9130401M01 Rik | |
| Ckap4 | | | Rexo2 | | Chracl | |
| Elk3 | | | Cryab | | Polr2f | |
| Phldal | | | Anxa2 | | Tomm22 | |
| Llph | | | Nedd4 | | Adsl | |
| Hmga2 | | | Cdl09 | | Rbxl | |
| Tmem5 | | | Iraklbpl | | Phf5a | |
| Col4al | | | Syncrip | | Nhp2ll | |
| Tm2d2 | | | Pcolce2 | | Rrp7a | |
| Rwdd4a | | | Mras | | Tubala | |
| Cpe | | | Pcbp4 | | Ranbpl | |
| Tpm4 | | | Ifrd2 | | Hmgnl | |

| | | | | | |
|---|---|---|---|---|---|
| Dnajbl | | | Cmtm7 | | Tmem242 | |
| Piezol | | | Purb | | Mrpll8 | |
| Tcf25 | | | GrblO | | Rnpsl | |
| Itgbl | | | Sptbnl | | Ube2i | |
| Flnb | | | Ccngl | | Stubl | |
| Gchl | | | Chd3 | | Mrpl28 | |
| Pnp | | | Pfnl | | Srsf3 | |
| Mmpl4 | | | Txndcl7 | | Glol | |
| Esd | | | Emc6 | | Mrpll4 | |
| Kctdl2 | | | Nxn | | Srsf7 | |
| Dnajc3 | | | Timm22 | | Snrpdl | |
| Ipo5 | | | Ccl7 | | Hdac3 | |
| Amotll | | | Duspl4 | | Cdk2ap2 | |
| Tagln | | | Nme2 | | Corolb | |
| Pafahlb2 | | | Spop | | Ppplca | |
| Rcn2 | | | FkbplO | | Mrplll | |
| Csk | | | Ptrf | | Sf3b2 | |
| Tpml | | | Becnl | | Eiflad | |
| Bnip2 | | | Vatl | | Cfll | |
| Tmed3 | | | Limd2 | | Ssscal | |
| Plscrl | | | Syngr2 | | Polr2g | |
| Rassfl | | | Faml95b | | Tmeml09 | |
| Prkar2a | | | Histlh2ap | | Prpfl9 | |
| Crtap | | | Faml 20a | | Rcll | |
| Slc35e4 | | | Gadd45g | | Nolcl | |
| Ccm2 | | | Sfxnl | | Zdhhc6 | |
| Anxa6 | | | Cltb | | mt-Cytb | |
| Mprip | | | Serfl | | | |
| Map2k3 | | | Mast4 | | | |
| Pitpna | | | Sdcl | | | |
| Myolc | | | Soxll | | | |
| FamlOlb | | | Bzw2 | | | |
| Tnfaipl | | | Bazla | | | |
| Mmd | | | Faml 77a | | | |
| Ccdcl37 | | | Timm9 | | | |

| | | | | | |
|---|---|---|---|---|---|
| P4hb | | | Synj2bp | | |
| Arhgdia | | | Calml | | |
| Sox4 | | | Meg3 | | |
| Tubb2a | | | Aktl | | |
| Pxdcl | | | Oxctl | | |
| Txndc5 | | | Ywhaz | | |
| Bicd2 | | | Eny2 | | |
| Tgfbi | | | Myc | | |
| Pdcd6 | | | Txn2 | | |
| Vcan | | | Polr3h | | |
| Tmeml67 | | | Zcrbl | | |
| Zcchc9 | | | Dazap2 | | |
| Maplb | | | Prrl3 | | |
| Gpx8 | | | Carhspl | | |
| Fst | | | Emp2 | | |
| Rock2 | | | Faml62a | | |
| FamllOc | | | Fstll | | |
| Ifrdl | | | Chmp2b | | |
| Cfl2 | | | Cdknla | | |
| Mgat2 | | | Clicl | | |
| Flrt2 | | | Mydgf | | |
| Fbln5 | | | Memol | | |
| Ddx24 | | | Srpl9 | | |
| Kiel | | | Reep5 | | |
| Ghr | | | Dpysl3 | | |
| Baspl | | | Ap3sl | | |
| Mtdh | | | Ppic | | |
| Plec | | | Gml6286 | | |
| Rpsl9bpl | | | Txnl4a | | |
| Desil | | | Gstpl | | |
| Tspo | | | Prdx5 | | |
| Slc48al | | | Famllla | | |
| Fkbpll | | | Ak3 | | |
| Comt | | | | | |
| Vps8 | | | | | |
| Lpp | | | | | |

| | | | | | |
|---|---|---|---|---|---|
| Ccdc50 | | | | | |
| Senp5 | | | | | |
| Ccdc80 | | | | | |
| Phldb2 | | | | | |
| Cldndl | | | | | |
| App | | | | | |
| Tnfrsfl2a | | | | | |
| Uqcc2 | | | | | |
| Slc39a7 | | | | | |
| Ppplrl8 | | | | | |
| Myll2a | | | | | |
| Lbh | | | | | |
| Cyplbl | | | | | |
| Mcfd2 | | | | | |
| Slc39a6 | | | | | |
| Binl | | | | | |
| Egrl | | | | | |
| Smim3 | | | | | |
| Tubb6 | | | | | |
| 1810055G02 Rik | | | | | |
| Fosll | | | | | |
| Neatl | | | | | |
| Rps6ka4 | | | | | |
| Ppplrl4b | | | | | |
| Ahnak | | | | | |
| Fthl | | | | | |
| Ccdc86 | | | | | |
| Anxal | | | | | |
| Acta2 | | | | | |
| Myof | | | | | |
| Tm9sf3 | | | | | |

[0272]    In particular, regulatory analysis identified a series of TFs that were upregulated in cells along the trajectory to iPSCs and predictive of the expression of the pluripotency programs (FIG. 26D). The earliest predictive TFs were expressed at day 9 (including Nanog, Sox2, Mybl2,

Elf3, Tgifl, Klf2, Etv5, and Cdc51) and additional predictive TFs were induced at day 10 (including Klf4, Esrrb, Spic, Zfp42, Hesxl, and Msc). Of these 14 TFs, 9 had previously described roles in regulation of pluripotency (Nanog, Sox2, Mybl2, Klf2, Cdc51, Klf4, Esrrb, Zfp42, and Hesxl) (Aaronson et al., 2016; Boheler, 2009; Buganim et al., 2012; Hu et al., 2009; Jeon et al., 2016; Li et al., 2015; Shi et al., 2006). A further wave of predictive TFs was upregulated in the iPSC trajectory between day 12 and 14, including Obox6, Sohlh2, Ddit3, and Bhlhe40. Among these late TFs, Obox6 and Sohlh2 were particularly notable, because they were not induced in the trajectories to any other cell fate. Obox6 and Sohlh2 had not previously been reported to be involved in regulation of pluripotency, but both had been implicated in maintenance and survival of germ cell development (Park et al., 2016; Rajkovic et al., 2002).

[0273]     An important change known to occur in the late stages of successful reprogramming was the reversal of X-chromosome inactivation in female cells. Our trajectory analysis identified the correct order of events as previously reported, but without the need for specialized experiments. Specifically, a study based on microscopy of cells labeled with antibodies to specific pluripotency proteins and RNA FISH for Xist (Pasque et al., 2014) showed that Xist downregulation preceded X-chromosome reactivation and positioned these events relative to the appearance of four pluripotency-associated proteins in Nanog-positive cells. Consistently, in our model, along the trajectory to successful reprogramming (but not elsewhere), cells at day 10 showed strong downregulation of Xist but did not yet display a signature of X-reactivation (FIGs. 26E, 26F, Methods). X-reactivation was complete at day 18, with the signature score having risen from 1.05 at day 10 to -1.95 at day 18, consistent with the expected increase in X-chromosome expression (FIG. 26F) (Pasque et al., 2014).

[0274]     Development of extra-embryonic-like cells during reprogramming

[0275]     Our trajectories showed that another subset of cells emerges from the MET Region, gained a strong epithelial signature by day 9, and went on to express a clear trophoblast signature (FIG. 27A, 27B). The trophoblast signature was detectable by day 10.5 and peaked by day 12.5, when such cells accounted for -20% of all cells in both serum and 2i conditions (FIG. 24G). Trophoblast and pre-implantation programs had previously been observed late in human reprogramming (Cacchiarelli et al., 2015)

[0276]    The cells spanned a spectrum of developmental programs associated with specific trophoblasts subsets. Briefly, in normal development the extraembryonic trophoblast progenitors (TPs) gave rise to the chorion, which formed labyrinthine trophoblasts (LaTBs), and the ectoplacental cone, which gave rise to various types of spongiotrophoblasts (SpTBs) and trophoblast giant cells (TGCs), including spiral artery trophoblast giant cells (SpA-TGCs). We scored our cells with signatures we derived from placental scRNA-seq (Nelson et al., 2016) for TP, SpT, TG and SpA-TGCs (Table 15), as well as three well-characterized markers (Msx2, Gcml and Cebpa) of LaTBs (Simmons et al., 2008; Ueno et al., 2013), for which no data were available to derive signatures (FIG. 33A). A substantial number of cells expressed TP, SpTB or SpATG signatures in serum conditions and TP or SpTB signatures in 2i conditions, at 10% FDR (Figure 5C). We also observed a cluster of -200 trophoblasts cells that expressed the three LaTBs markers (in 2i but not serum), which were largely separate from those expressing signatures of ectoplacental derivatives. In addition to trophoblast-like cells, -125 cells expressed a signature (Lin et al., 2016) for the primitive endoderm (XEN-like cells), the other cell type that contributes to extraembryonic tissue (FIG. 33B, FDR 0.1%). Notably, these cells were seen only in a single replicate at a single time point (day 15.5) in serum conditions only. Two previous studies reported the generation of XEN-like cells during OKSM-induced reprogramming to iPSCs (Parenti et al., 2016, Zhao et al., 2018).

[0277]    Regulatory analysis associated various TFs with the trajectory from the MET Region to the overall set of trophoblasts (FIG. 27B). TFs at day 10.5 that were predictive of subsequent trophoblast fates included several involved in trophoblast self-renewal (Gata3, Elf5, Mycn, Mybl2) (Kidder and Palmer, 2010) and early trophoblast differentiation (Ovol2, Ascl2) (Latos and Hemberger, 2016), as well as others expressed in trophoblasts but without known roles in trophoblast differentiation (Rhox6, Rhox9, Batf3 and Elf3).

[0278]    Trajectory and regulatory analysis also identified TFs that were predictive of specific cell subsets. Ancestors of cells with the TP signature expressed Gata3, Pparg, Rhox9, Mytll, Hnflb, and Prdml 1. Gata3 was involved for trophoblast progenitor differentiation (Ralston et al., 2010) and Pparg was involved for trophoblast proliferation and differentiation of labyrinthine trophoblasts (Parast et al., 2009). The other TFs were known to be expressed in placenta, but their roles in cellular differentiation had not been well characterized. Ancestors of cells with the

SpTB or LaTB signature expressed Gata2, Gcml, Msx2, Hoxdl3, and Nrlh4. Gata2 was known to be involved for regulation of specific trophoblast programs (Ma et al., 1997). Gcml and Msx2 had specific roles in LaTB differentiation, EMT and trophoblast invasion (Liang et al., 2016; Simmons and Cross, 2005), respectively. Nrlh4 was detected in placental tissue, but its role in trophoblast differentiation had not been characterized. Ancestors of cells with the SpA-TGC signature expressed Handl, Bbx, Rhox6, Rhox9, and Gata2. Handl was known to be necessary for trophoblast giant cell differentiation and invasion (Scott et al., 2000). Bbx was a core trophoblast gene known to induced by upstream TFs Gata3 and Cdx2 (Ralston et al., 2010) **(FIGs. 33A-33E)** .

[0279]    Neural-like cells also emerged from the MET Region during reprogramming in serum conditions.

[0280]    Only in serum conditions, a third subset of cells emerged from the MET Region, gained a strong epithelial signature, and went on to develop clear neural signatures (FIGs. 27D-27F). These cells were not seen in 2i conditions, presumably due to the differentiation inhibitors in this condition. Compared to the trophoblast-like cells, the signature for neural identity emerged more slowly, by roughly two days (FIG. 24G). The ancestors of neural like cells diverged from the ancestors of trophoblasts and iPSCs by day 9 (FIG. 26B), and then underwent a rapid transition at day 12.5, losing their epithelial signatures and gaining neural signatures (FIGs. 27D, 27E). The signature was maintained through day 18, when such cells comprised 21.5% of all cells in serum conditions.

[0281]    In normal neural development, neuroepithelial cells lost their epithelial identity and upregulated glial factors, transforming into radial glial cells (Florio and Huttner, 2014; Ming and Song, 201 1). Radial glial cells gave rise to astrocytes and oligodendrocytes, and in the CNS also served as progenitors for many neurons (Ming and Song, 201 1). To probe these identities, we used scRNA-Seq data from mouse brain to derive signatures that distinguished different cell types and differentiation states (Table 15). These included signatures of (i) astrocytes, oligodendrocyte precursor cells (OPCs), and neurons in adult brain from in the Allen Brain Atlas (http://www.brain-map.org), and (ii) three unlabeled clusters of radial glial cells in E18 mouse brain (Han et al., 2018), each distinguished by high expression of a different gene (Id3, GdflO, and Neurog2, respectively).

**[0282]** Cells in the landscape spanned multiple stages of neuronal differentiation. Cells near the base of the "neural spike" in the landscape (day 12.5-18) expressed radial glial and neural stem-cell markers (including Pax6 and Sox2) and cells further out along the spike (day 15-18) expressed markers of neuronal differentiation (including Neurog2 and Map2. About 70% of the neural-like cells had significant expression (at 10% FDR) of at least one of the six signatures (FIG. 27G). Cells with the three radial glial signatures appeared first, concurrent with the loss of epithelial identity and first gained of neural lineage identity by day 12.5 (FIG. 27F). Cells expressing the signatures derived from adult neurons and glia emerged around day 14 in the neural spike and grew in abundance for the duration of the time course. Their ancestors were concentrated in the radial glial populations on day 13.5, with a particular concentration in the GdflO RG subpopulation. While the glial populations overlapped substantially, the neurons form a distinct population with substantial substructure. The subset of cells with signatures of adult neurons included cells with canonical markers for excitatory and inhibitory neurons (Slcl7a6 and Gadl, respectively). Expression signatures that distinguished these two classes of cells showed strong, albeit incomplete, overlapped with respective programs of excitatory and inhibitory neurons in the Allen Brain Atlas (FIG. 27G, Methods).

**[0283]** Regulatory analysis identified TFs predictive of the overall neural-like cell population, with the top TFs all known to have roles in various stages of neurogenesis. These TFs included those known to promote early neurogenesis (Rarb, Foxp2, Emxl, Pou3f2, Nr2fl, Mytll, Neurod4), regulated late neurogenesis (Scrt2, Nhlh2, Pou2f2), regulated differentiation and survival of neural subtypes (Onecutl, Tal2, Barhll, Pitx2), and played roles in neural tube formation (Msxl, Msx3).

**[0284]** <u>The developmental landscape highlighted potential paracrine signals</u>

**[0285]** As the reprogramming landscape included a substantial and under-appreciated diversity of differentiating cell subsets, including stromal, epithelial, neural and trophoblast cells, we asked how they might affect each other as they undergo dynamic processes concurrently. In particular, paracrine signaling played a key role in normal development and had also been shown to affect reprogramming, with secretion of inflammatory cytokines enhancing reprogramming efficiency (Mosteiro et al., 2016). Accordingly, we systematically cataloged the contemporaneous occurrence of ligand-receptor pairs across cell subsets in the developmental

landscape. We defined an interaction score based on the product of (1) fraction of cells of type A expressing ligand X and (2) the fraction of cells of type B expressing the cognate receptor Y, at the same time t (FIGs. 28A, 28B and 34B, Methods). We examined 180 individual cognate ligand-receptor pairs, as well as an aggregate score across all pairs between cell clusters (FIG. 34A) and across those pairs related to the SASP signature.

[0286]　　　The landscape revealed rich potential for paracrine signaling (FIG. 28B, FIG. 34B, Table 18). In particular, we observed high interaction scores for several SASP ligands in stromal cells with receptors expressed in iPSCs, such as Gdf9 with Tdgfl (Polo et al., 2012) and Cxcll2 with Dpp4 (FIGs. 28C, 28F, 34C).

**Table 18** - Potential ligand-receptor pairs between stromal cells and iPSCs, neural-like cells, and trophoblast cells ranked by standardized interaction scores

| Ligand: Stromal cells. Receptor: iPSCs | | | Ligand: Stromal cells. Receptor: Neural-like cells | | | Ligand: Stromal cells. Receptor: Trophoblast cells | | |
|---|---|---|---|---|---|---|---|---|
| Ligand-Receptor Pair | Maximal standardized interaction score | Peak Score Day | Ligand-Receptor Pair | Maximal standardized interaction score | Peak Score Day | Ligand-Receptor Pair | Maximal standardized interaction score | Peak Score Day |
| | | | | | | | | |
| Gdf9.Tdgfl | 55.83015277 | 14 | Crlfl.Cntfr | 76.16064491 | 16.5 | Csfl.Csflr | 111.8151997 | 18 |
| Cxcll2.Dpp4 | 42.40247659 | 12.5 | Fgf2.Vtn | 66.31283077 | 18 | Cxcl5.Cxcr2 | 102.1031447 | 18 |
| Ngf.Ngfr | 26.79815659 | 12 | Clcfl.Cntfr | 52.04021271 | 15.5 | Cxcll.Cxcr2 | 85.46017232 | 18 |
| Cclll.Dpp4 | 23.75254375 | 14 | Vegfa.Vtn | 39.99828338 | 18 | Il6.Il6ra | 70.79780689 | 18 |
| Kitl.Kit | 20.48156022 | 17.5 | Bdnf.Ntrk2 | 38.24132006 | 17 | Cxcl2.Cxcr2 | 68.04261554 | 18 |
| Ccl5.Dpp4 | 20.22465038 | 12.5 | Tgfb2.Vtn | 37.9492686 | 18 | Cxcl3.Cxcr2 | 62.67646817 | 17.5 |
| Inhba.Acvr2b | 18.91224205 | 17 | Tgfbl.Vtn | 37.71506462 | 18 | Il7.Il2rg | 57.89558657 | 17 |
| Fgf7.Fgfr4 | 18.88448993 | 12 | Tgfb3.Tgfbrl | 32.86035119 | 17 | Vegfa.Fltl | 52.30228603 | 18 |
| Nppc.Nprl | 17.71660947 | 16.5 | Bdnf.Sortl | 29.14910223 | 17 | Tg.Lrp2 | 45.35387653 | 9.5 |
| Fgf7.Fgfr2 | 17.2915253 | 9 | Ill6.Grin2a | 27.83837935 | 13.5 | Ccl2.Ackr2 | 44.70456305 | 17 |
| Grn.Cryl | 17.25111965 | 17 | Inhba.Acvr2b | 25.85377693 | 15.5 | Sppl.Itgbl | 44.39437623 | 18 |
| Fgf2.Fgfr3 | 17.18398331 | 15.5 | Apln.Aplnr | 23.46381586 | 14 | Ill5.Il2rg | 43.96702273 | 18 |
| Sppl.F2 | 16.91745599 | 17 | Bmpl.Adrala | 21.99556814 | 17.5 | Ccl7.Ackr2 | 42.35095481 | 17 |
| Tgfb3.Tgfbrl | 15.80306191 | 9 | Ill6.Grin2b | 21.85263644 | 18 | Tnfsf9.Tnfrsf9 | 41.80288631 | 15.5 |
| Bdnf.Ntrk2 | 15.73929703 | 12 | Vegfa.Ephb2 | 21.76727834 | 17 | Cxcll5.Cxcr2 | 41.37975891 | 18 |
| Avp.Avprlb | 15.6652861 | 15 | Tgfbl.Tgfbrl | 21.71078611 | 17 | Vegfb.Fltl | 40.59359924 | 18 |
| Inhbb.Acvr2b | 15.22902239 | 18 | Ngf.Sortl | 21.55867193 | 16.5 | Fgf2.Fgfrl | 40.1892017 | 18 |
| Tnfsf8.Tnfrsf8 | 14.9661866 | 17.5 | Ereg.Erbb4 | 21.23888338 | 17 | Ill5.Il2rb | 37.23349427 | 18 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Ucn2.Crhr2 | 14.66104887 | 14 | Cxcll2.Cxcr4 | 20.66598418 | 16.5 | II2.II2rg | 34.72049417 | 17 |
| Sst.Sstr3 | 14.53946813 | 12.5 | Nov.Notch 1 | 20.64844205 | 17 | Illrn.Illr2 | 34.60876011 | 18 |
| Cxcll2.Cxcr4 | 13.99702972 | 9.5 | Inhbb.Acvr2b | 20.20541981 | 15.5 | Bmp4.Bmpr2 | 33.37381523 | 18 |
| Fgfl.Fgfr4 | 13.23808582 | 14 | Egf.Vtn | 20.11367671 | 14.5 | Ppbp.Cxcr2 | 33.31119733 | 17 |
| Gdf6.Bmprlb | 13.23695383 | 11.5 | Fgf7.Fgfr2 | 19.85021209 | 9 | Flt3I.FIt3 | 31.32026205 | 17 |
| Gdf9.Bmprlb | 12.81536347 | 11.5 | FgflO.Fgfr2 | 19.77063453 | 12 | Inhba.Acvr2b | 31.21420166 | 16.5 |
| Gdf5.Acvr2b | 12.41295756 | 17.5 | Fgf2.Fgfr3 | 19.20901825 | 18 | II2.II2rb | 31.17852066 | 17 |
| Cxcl3.Cxcr2 | 12.28144255 | 9 | Inhba.Igsfl | 19.00415822 | 13.5 | Inhbb.Acvrlb | 31.08869402 | 18 |
| CxcllO.Dpp4 | 12.0118101 | 16.5 | Pomc.Vtn | 18.61879864 | 14 | Inhba.Acvrlb | 30.95069812 | 18 |
| Tnfsfll.Tnfrsflla | 11.98501062 | 18 | Tgfb2.Tgfbrl | 18.40997602 | 17 | Ccl8.Ackr2 | 30.92303758 | 17 |
| Tnfsfll.Med24 | 11.31495458 | 17 | Gdf9.Tdgfl | 18.12847923 | 10.5 | Pgf.Fltl | 28.55965416 | 17 |
| Bdnf.Inpp5k | 11.02760154 | 17 | Gdnf.Gfral | 17.94758176 | 18 | Tgfb3.Tgfbrl | 28.48415966 | 18 |
| Cxcl5.Cxcr2 | 10.76725496 | 9 | Ednl.Ednrb | 17.81157803 | 17 | Inhba.Tgfbr3 | 27.97080183 | 18 |
| Bmp2.Bmprlb | 10.52856679 | 11.5 | Gdfll.Acvr2b | 16.93911315 | 15.5 | Inhbb.Acvr2b | 27.64710304 | 18 |
| Inhba.Acvrlb | 10.45689595 | 15.5 | Gdf5.Bmprlb | 16.87028377 | 17 | Ccl3.Ackr2 | 27.17947452 | 14.5 |
| Fgfl.Fgfr3 | 9.904359216 | 14 | Gdf5.Acvr2b | 16.68587549 | 15.5 | Tgfb3.Sdc4 | 26.70563028 | 18 |
| Tgfb3.Eng | 9.606914311 | 18 | Igfl.Igflr | 16.40043325 | 17.5 | Inhba.Acvrll | 24.8733331 | 16.5 |
| Crlfl.Cntfr | 9.491489628 | 9 | Ngf.Ngfr | 16.1554284 | 9 | Wnt5a.Fzd5 | 24.08669584 | 18 |
| Tg.Lrp2 | 9.311152429 | 9.5 | Cxcl5.Ackrl | 15.81074369 | 17 | Egf.Erbb3 | 22.88090865 | 18 |
| Nppa.Nr5a2 | 9.196846339 | 15.5 | Tg.Lrp2 | 15.56587296 | 9.5 | Gdf5.Acvr2b | 22.79535492 | 16.5 |
| Sppl.Itgbl | 9.094293313 | 9 | Ill6.KcnjlO | 15.40280917 | 15 | Tgfbl.Itgb6 | 22.73325122 | 18 |
| Tgfb3.Sdc4 | 8.962618473 | 18 | Ccl2.Ackrl | 14.80314224 | 17 | Vegfc.Flt4 | 22.64781847 | 18 |
| Avp.Avpr2 | 8.816318411 | 16 | Illrn.Illr2 | 14.70537108 | 17 | Vegfa.Kdr | 21.61880314 | 13 |
| Bmp4.Bmprlb | 8.789458439 | 11.5 | Wnt5a.Fzd2 | 14.59368545 | 16.5 | Ill8.Ill8rap | 21.45320636 | 18 |
| Gdfll.Acvr2b | 8.657009643 | 17.5 | Inhbb.Igsfl | 14.56070266 | 13.5 | Tgfb2.Tgfbr3 | 21.43696896 | 12.5 |
| Ctgf.Egfr | 8.474450513 | 9 | Ccll2.Ackrl | 14.48343455 | 15 | Fgf7.Fgfr2 | 21.27556999 | 9 |
| Nov.Notch 1 | 7.853128492 | 9.5 | Ccl7.Ackrl | 14.45732094 | 17 | Ccll2.Ackr2 | 20.65465765 | 15 |
| Cxcll.Cxcr2 | 7.825570863 | 9 | Fgfl.Fgfr3 | 13.98128161 | 14 | Tgfbl.Tgfbr3 | 19.078023 33 | 18 |
| Pomc.Mc5r | 7.803289928 | 13 | Cort.Sstr2 | 13.83366019 | 14.5 | Cclll.Ackr2 | 19.06812091 | 16.5 |
| Inhba.Acvr2a | 7.697312114 | 10 | Vegfa.Kdr | 13.52841955 | 17 | Ccl28.Ackr2 | 19.0608243 | 16.5 |
| Ill6.Cd4 | 7.691300029 | 16 | Bmp4.Bmprlb | 13.17024743 | 17 | Kitl.Kit | 18.32774459 | 10 |
| Hcrt.Npffr2 | 7.611421106 | 14.5 | Igfl.Igsfl | 13.1615924 | 13.5 | Gdfll.Acvr2b | 17.1611013 | 16.5 |
| Nppa.Nprl | 7.327171012 | 15.5 | Inhba.Acvr2a | 12.86079359 | 15.5 | Bdnf.Inpp5k | 16.94541624 | 18 |
| Fgf2.Fgfrl | 6.935257539 | 18 | Gdnf.Gfra2 | 12.82585678 | 18 | Ccl5.Ackr2 | 16.65970084 | 10.5 |
| Inhbb.Acvrlb | 6.8878958 | 15.5 | Ntf3.Ntrk2 | 12.69375513 | 14 | Ngf.Ngfr | 16.41502139 | 9 |
| Ccll7.Ccr4 | 6.846358767 | 17 | Cxcll.Ackrl | 12.64243264 | 17 | Igfl.Igflr | 16.27850014 | 18 |

| | | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Ill6.Grin2b | 6.789839819 | 14.5 | Fgf2.Fgfrl | 12.31083274 | 18 | Bmp2.Bmpr2 | 15.99972954 | 18 |
| Bdnf.Sortl | 6.67375428 | 9 | Vegfa.Nrp2 | 12.23441434 | 18 | Tgfbl.Acvrll | 15.96504429 | 16.5 |
| Tgfb2.Tgfbrl | 6.519268162 | 9 | Bmp6.Acvr2b | 12.1758211 | 13.5 | Gdf5.Bmpr2 | 15.58998037 | 16.5 |
| Ntf3.Ntrk2 | 6.438685726 | 12 | Hbegf.Erbb4 | 12.00500039 | 14.5 | Tgfb2.Tgfbrl | 15.53065603 | 18 |
| Ccl3.Ccr5 | 6.407610415 | 12.5 | Vegfc.Kdr | 11.97527882 | 18 | Tgfbl.Tgfbrl | 15.49109459 | 18 |
| Ptn.Plxnb2 | 6.364004505 | 9 | Ccll7.Ackrl | 11.93535268 | 16 | Inha.Tgfbr3 | 14.94814105 | 18 |
| Egf.Erbb3 | 6.33209249 | 17 | Cxcl3.Cxcr2 | 11.79741482 | 9 | Ccl27a.Ackr2 | 14.35654443 | 17 |
| Fgf9.Fgfr3 | 6.17049013 | 15.5 | Wnt2.Fzd9 | 11.76547196 | 14.5 | Pf4.Ldlr | 13.49144052 | 17.5 |
| Ntf3.Ntrk3 | 6.071479576 | 12.5 | Tnfsfll.Med24 | 11.58428169 | 17 | Vegfc.Kdr | 13.42241254 | 12.5 |
| Wnt5a.Fzd5 | 6.049412152 | 17.5 | Cxcll5.Ackrl | 11.39063421 | 16 | FgflO.Fgfr2 | 12.93211376 | 12 |
| Ill6.Kcnj4 | 5.956600472 | 9 | Cxcl5.Cxcr2 | 10.81475088 | 9 | Pdgfc.Pdgfra | 12.7181284 | 18 |
| FgflO.Fgfr2 | 5.735961453 | 10 | Sppl.Itgbl | 10.57557893 | 9 | Ccl25.Ackr2 | 12.58225578 | 10.5 |
| Csf3.Csf3r | 5.660332275 | 18 | Ccl8.Ackrl | 10.24654012 | 18 | Crlfl.Cntfr | 12.56270017 | 9 |
| Ngf.Sortl | 5.631416895 | 9 | Gdf5.Acvr2a | 9.947335355 | 16.5 | Inhba.Acvrl | 12.49512116 | 18 |
| Wnt2.Fzd9 | 5.625683619 | 13 | Inhbb.Acvr2a | 9.83065505 | 17.5 | Inhbb.Acvrl | 12.17571989 | 18 |
| Ngf.Ntrkl | 5.482536008 | 18 | Bmp2.Bmprlb | 9.823905055 | 17 | Bmp4.Bmprla | 12.13592365 | 18 |
| Ccl2.CcrlO | 5.204305876 | 9 | Ngf.Ntrkl | 9.765431603 | 15.5 | Hgf.Met | 11.85706092 | 18 |
| Gdf5.Bmprlb | 5.164323069 | 11.5 | Ctgf.Egfr | 9.510948488 | 9 | Avp.Avprlb | 11.8443167 | 12.5 |
| Ccl7.CcrlO | 5.03794601 | 9 | Ill6.Grin2c | 9.210664243 | 16.5 | Wnt5a.Lrp6 | 11.2866016 | 18 |
| Inhba.Igsfl | 4.652799622 | 16.5 | Igf2.Vtn | 9.08515341 | 15.5 | Illrn.Illrl | 11.21386458 | 18 |
| Igfl.Igsfl | 4.623901723 | 16.5 | Fgf9.Fgfr3 | 8.929720296 | 13 | Npff.Npffr2 | 11.12680175 | 12.5 |
| Kitl.Epor | 4.572546653 | 9 | Ucn2.Crhr2 | 8.529535163 | 10 | Gpil.Amfr | 11.09557616 | 18 |
| Bmp6.Bmprlb | 4.21969712 | 11.5 | Gdf9.Bmprlb | 8.458633534 | 12.5 | Ccl2.Ccr5 | 10.87678026 | 17 |
| Ill6.Grin2a | 4.182303182 | 12 | Cxcll.Cxcr2 | 8.317259429 | 9 | Inhba.Acvr2a | 10.71764165 | 18 |
| Tgfbl.Tgfbrl | 4.165309406 | 9 | Pnoc.Oprll | 8.170486417 | 13 | Inhbb.Acvr2a | 10.62573575 | 18 |
| Hmgbl.Pgr | 4.162814163 | 9.5 | Inha.Acvr2a | 8.005902758 | 15.5 | Ccll7.Ccr4 | 10.22222634 | 11.5 |
| Tnfsfl3b.Tnfrsfl7 | 4.077062584 | 16.5 | Inhba.Acvrlb | 7.58971181 | 9.5 | Vegfa.Lyvel | 9.978529316 | 11.5 |
| Ill6.Grin2c | 3.818702923 | 17 | Fgf7.Fgfr4 | 7.313765731 | 16 | Lif.Lifr | 9.836393324 | 16.5 |
| Crh.Crhr2 | 3.804963778 | 14 | Ptn.Plxnb2 | 7.174330257 | 9 | Il25.Ill7rb | 9.820316363 | 16 |
| Tgfbl.Eng | 3.789167413 | 17 | Btc.Erbb4 | 7.130596933 | 14.5 | Ccl8.Ccr5 | 9.277471947 | 16.5 |
| Ccl5.Ccr5 | 3.765684384 | 10.5 | Grn.Cryl | 7.038337946 | 16.5 | Ill6.KcnjlO | 9.099847388 | 14.5 |
| Ccl3.Ackr4 | 3.748657973 | 12.5 | Ill6.Kcnj2 | 7.031491551 | 18 | Bdnf.Ntrk2 | 9.027486627 | 12.5 |
| Ccl2.Ccr5 | 3.746070011 | 12.5 | Ednl.Ednra | 6.737910303 | 17.5 | Ednl.Ednrb | 8.719812556 | 14 |
| Gdf5.Acvr2a | 3.726614996 | 16 | Avp.Oxtr | 6.701328931 | 16.5 | Cxcll2.Cxcr4 | 8.696493411 | 17 |
| Npff.Npffr2 | 3.71584242 | 14.5 | Tgfb3.Sdc4 | 6.648807091 | 9 | Fgf9.Fgfrl | 8.617860569 | 18 |
| Inhbb.Igsfl | 3.660059949 | 16.5 | Ill6.Kcnj4 | 6.296091418 | 9 | Sppl.F2 | 8.219496273 | 13.5 |
| Bmp6.Acvr2b | 3.613241885 | 13.5 | Sppl.F2 | 6.250718711 | 14.5 | Ptn.Plxnb2 | 8.085698538 | 9 |

224

| | | | | | |
|---|---|---|---|---|---|
| Lif.Lifr | 3.59302184 | 12.5 | Adm.Calcrl | 6.127364131 | 18 |
| Inhbb.Acvr2a | 3.573362535 | 16 | Artn.Gfra3 | 6.100580729 | 18 |
| Tgfb2.Eng | 3.493150482 | 18 | Ccl5.Ackrl | 6.08281121 | 16 |
| Tnfsfl3b.Tnfrsfl3b | 3.485242199 | 14 | Tgfb3.Eng | 6.075334099 | 9 |
| Bmp2.Bmprla | 3.421538818 | 9 | Gdf6.Bmprlb | 5.814695498 | 17.5 |
| Bmp2.Eng | 3.277644443 | 12 | Hmgbl.Pgr | 5.524547346 | 9.5 |
| Pf4.Ldlr | 3.252582504 | 11.5 | Wnt5a.Lrp6 | 5.416442742 | 15 |
| Ntf5.Ngfr | 3.228481212 | 12 | Vegfa.Lyvel | 5.365931818 | 16.5 |
| Ccl5.Ccr4 | 3.054614918 | 17 | Ccll7.Ccr4 | 5.313995351 | 9.5 |
| Pgf.Nrp2 | 3.013909017 | 9 | Sst.Sstr2 | 4.993026408 | 12.5 |
| Fgf8.Fgfr4 | 3.01220056 | 14 | Vegfa.Fltl | 4.860449031 | 13.5 |
| Artn.Gfra3 | 3.008145345 | 16 | Bmp6.Bmprlb | 4.604550067 | 16.5 |
| | | | Egf.Erbb3 | 4.487189494 | 10.5 |
| | | | Kitl.Epor | 4.470894246 | 9 |
| | | | Gdf9.Acvr2a | 4.461925767 | 12.5 |
| | | | Ccl2.CcrlO | 4.287535378 | 9 |
| | | | Fgf9.Fgfr2 | 4.104799154 | 11 |
| | | | Ill6.Cd4 | 4.102677906 | 15.5 |
| | | | Ccl2.Ccr5 | 4.06128803 | 18 |
| | | | Ntf3.Ntrkl | 4.045425855 | 15.5 |
| | | | Bmp2.Bmprla | 4.007512362 | 9 |
| | | | Pdgfc.Pdgfra | 4.000578173 | 18 |
| | | | Bmp4.Bmprla | 3.973107083 | 17 |
| | | | Ghrl.Ptger3 | 3.959803347 | 15 |
| | | | Illl.Illlral | 3.931542903 | 16.5 |
| | | | Ccl7.CcrlO | 3.86216627 | 9 |
| | | | Gdf5.Bmprla | 3.812514632 | 16.5 |
| | | | Ntf5.Ntrk2 | 3.800422565 | 15.5 |
| | | | Ntf3.Ntrk3 | 3.791204113 | 13 |
| | | | Ccl8.Ccr5 | 3.6877203 | 18 |
| | | | Vegfb.Fltl | 3.67289066 | 13.5 |
| | | | Ccl5.Ccr4 | 3.652617678 | 9.5 |
| | | | Inhba.Acvrl | 3.386360757 | 18 |
| | | | Inhbb.Acvrl | 3.330148881 | 18 |
| | | | Wntl.Fzd9 | 3.30422519 | 12.5 |
| | | | Npff.Npffrl | 3.243049647 | 16 |

| | | |
|---|---|---|
| Tnfsfll.Med24 | 8.080587047 | 18 |
| Ctgf.Egfr | 8.025815916 | 9 |
| Ghrl.Ptger3 | 7.831218363 | 15 |
| Ctfl.Lifr | 7.478421588 | 18 |
| Pdgfd.Pdgfrb | 7.440471865 | 18 |
| Gdf5.Acvr2a | 7.437486529 | 17.5 |
| Cxcll2.Dpp4 | 7.386223592 | 12.5 |
| Cclll.Ccr5 | 7.344244377 | 16.5 |
| Gdf5.Bmprla | 7.242141121 | 17.5 |
| Artn.Gfra3 | 6.624252893 | 16 |
| Ill8.Illrl2 | 6.470340015 | 18 |
| Inha.Acvr2a | 6.410004454 | 18 |
| Gdf6.Bmpr2 | 6.362677796 | 18 |
| Ntf3.Ntrk2 | 6.34714587 | 12.5 |
| Gdf5.Acvrl | 6.33836936 | 18 |
| Tslp.Prnp | 6.263327318 | 18 |
| Gdf9.Tdgfl | 6.170602382 | 10.5 |
| Bdnf.Sortl | 5.94172272 | 9 |
| Bmp2.Acvrl | 5.90978443 | 18 |
| Bmp6.Acvr2b | 5.871545931 | 13.5 |
| Tnfsfll.Tnfrsflla | 5.868170248 | 15.5 |
| II6.II6st | 5.857031136 | 18 |
| Kitl.Epor | 5.493268145 | 14 |
| Hmgbl.Pgr | 5.439455664 | 9.5 |
| Gdf9.Bmpr2 | 5.301534907 | 17.5 |
| Ngf.Sortl | 5.181692923 | 9 |
| Tnfsfl3b.Tnfrsfl3b | 5.166928123 | 15.5 |
| Ucn2.Crhr2 | 5.15524664 | 9 |
| Fgfl.Fgfrl | 5.090269326 | 18 |
| Pdgfa.Pdgfra | 4.960203778 | 18 |
| Fgf7.Fgfr4 | 4.959156503 | 12 |
| Nov.Notch 1 | 4.944351734 | 9.5 |
| Bmp2.Bmprla | 4.828229043 | 18 |
| Fgf2.Fgfr3 | 4.718080894 | 13.5 |
| Grn.Cryl | 4.629614942 | 9 |
| Tgfb3.Eng | 4.541775835 | 9 |

| | | | | | | | | TnfsflO.Tnfrs flOb | 4.456880919 | 16.5 |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | Hcrt.Hcrtrl | 4.407762506 | 14.5 |
| | | | | | | | | Ccl5.Ccr5 | 4.218364077 | 16 |
| | | | | | | | | H16.Kcnj4 | 4.184296843 | 9 |
| | | | | | | | | Ghrl.Ptgir | 4.00490292 | 15 |
| | | | | | | | | Cxcll6.Cxcr6 | 3.995533009 | 18 |
| | | | | | | | | Ccl3.Ccr5 | 3.825939759 | 12.5 |
| | | | | | | | | Ill6.Grin2c | 3.804620341 | 14 |
| | | | | | | | | Ccl5.Ccr4 | 3.700028296 | 13 |
| | | | | | | | | Ill7b.Ill7rb | 3.43715641 | 10.5 |
| | | | | | | | | Hmgbl.Ar | 3.425935882 | 11 |
| | | | | | | | | Ntf3.Ntrkl | 3.384388196 | 13 |
| | | | | | | | | Ngf.Ntrkl | 3.213785377 | 13 |
| | | | | | | | | Ccll2.Ccr5 | 3.032941015 | 16 |

[0287]       Analysis of the neural-like cells revealed particularly interesting interaction scores involving Cntfr (FIGs. 28D, 28G, 34D), an 116-family co-receptor whose activation played critical roles in neural differentiation and survival (Elson et al., 2000; Nakashima et al., 1999). On day 11.5 in serum conditions, one day before the early neuronal signatures appear, neural ancestors upregulated expression of Cntfr; expression was 4.6-fold higher in epithelial cells that were neural ancestors versus those that were not. Just before, on day 10.5, stromal cells began expressing three activating ligands for Cntfr (Crlfl, Lif, Clcfl). We speculated that these events may help trigger the program of neural differentiation among a subset of epithelial cells in serum conditions. The analysis also revealed a potential interaction involving the ligand-receptor pair Bdnf-Ntrk2, which had been implicated in promoting neuronal development, maturation and survival (Chen et al., 2015; Jukkola et al., 2006; Yun et al., 2008) (FIGs. 28D, 28G, 34D). The same ligand-receptor interactions were seen in 2i conditions, but the MEK inhibitor in 2i medium would be expected to block Cntfr signaling and subsequent neural differentiation.

[0288]       Trophoblast-like cells also showed notable interaction scores, including Csfl and Csflr (FIGs. 28E, 28H). In early placental development, Csfl was expressed in maternal columnar epithelial cells and Csflr was expressed in fetal trophoblasts, suggesting a functional role of this interaction in trophoblast development and differentiation. Many of the other top-

ranked interactions were between a single receptor in trophoblast cells (Cxcr2) and multiple members of the same ligand family (Cxcl5, Cxcll, Cxcl2, Cxcl3, and Cxcll5) (FIGs. 24E, 24H, 34E). Cxcr2 had been shown to be necessary for trophoblast invasion in human trophoblast cells (Vandercappellen et al., 2008; Wu et al., 2016).

[0289]    RNA expression revealed genomic aberrations in stromal and trophoblast-like cells

[0290]    We hypothesized that some cell types might harbor detectable genomic aberrations. In particular, trophoblasts were known to undergo endocycles of replication in vivo (Edgar et al., 2014), resulting in selective amplification of specific genomic regions containing functionally important genes (Hannibal and Baker 2016). Additionally, our stromal cells exhibited signs of stress and cell death which may be associated with genomic aberrations.

[0291]    To identify potential genomic aberrations, we scored the scRNA-Seq data for large regions showing coherent increases or decreases in gene expression, following successful approaches we developed to identify aberrant regions in individual tumor cells in a patient (Patel et al., 2014). We searched copy-number variations at the level of whole chromosomes and subchromosomal regions spanning 25 consecutive housekeeping genes (median size 25 Mb) (STAR Methods). To evaluate the detection of subchromosomal events, we analyzed scRNA-Seq data from oligodendroglioma (Tirosh et al. 2016): the method had high specificity, but sensitivity to detect only about one-third of events.

[0292]    Whole-chromosome aneuploidies were detected in 4.0% of trophoblast cells and 2.1% of stromal cells, compared to only 1.1% of all other cells across the landscape. Most whole-chromosome events were consistent with loss or gain of a single copy of the chromosome (FIG. 281). Subchromosomal events were detected in 6.9%> of trophoblast cells and 3.2% of stromal cells, compared to only 1.2% in most other cells types and 0.4% in neural cells (Figure 6J); the true proportions are likely to be about 3-fold higher, given the estimated sensitivity.

[0293]    Trophoblast-like cells showed recurrent events at a higher frequency than stromal cells. Among trophoblast cells harboring aberrations, 8.6% were detected as carrying a recurrent event involving apparent duplication (50% higher expression) of a region containing 74 genes (FIG. 28K). Among the genes are Wnt7b, which was required for normal placental development (Parr et al., 2001); Prr5, which mediates Pdfgb signaling required for development of labyrinthine cells (Ohlsson et al., 1999; Woo et al., 2007); and several genes identified as 'core

trophoblast genes' (Cyb5r3, Cenpm, Srebf2, and Pmml). The top 15 recurrent events also included the amplification of the prolactin gene cluster on chromosome 13 in 1% of cells. These observations suggested that the trophoblast-associated mechanisms of genomic alteration may be expressed, to some extent, in our trophoblast-like cells.

[0294]    In the stromal cells with evidence of genomic aberration, the most common recurrent events had lower frequency. Notably, however, the most frequently amplified region contained cell cycle inhibitors Cdkn2a, Cdkn2b, and Cdkn2c, while the most frequently lost region contained Cdkl3, which promotes cell cycling, and Mapk9, loss of which promotes apoptosis. These observations suggested that genomic alterations in these regions may contribute to development stromal cells.

[0295]    Forced expression of Obox6 enhanced reprogramming

[0296]    Finally, we explored whether some of the new TFs identified by regulatory analysis along the trajectory to iPSCs might provide ways to increase reprogramming efficiency. In principle, TFs could increase the efficiency of reprogramming in several ways, including increasing the transition frequency to iPSC precursors, boosting the growth rate of iPSC precursors, reducing alternative fates of other epithelial-related fates, or increasing supportive paracrine signaling from non-iPS cells.

[0297]    We focused on Obox6, which our regulatory analysis discovered as the TF most strongly correlated with reprogramming success, among those not previously implicated in the process. Obox6 (oocyte-specific homeobox 6) is a homeobox gene of unknown function that is preferentially expressed in the oocyte, zygote, early embryos and embryonic stem cells (Rajkovic et al., 2002). (Although Obox6 was the only Obox family member detected in our experiment, we note that a better-studied oocyte-specific homeobox Oboxl has been shown to enhance reprogramming efficiency, promote MET, and be able to substitute for Sox2 in reprogramming (Wu et al., 2017)). While Obox6 was expressed only in a small fraction of cells (<1%) before day 12, cells expressing Obox6 during day 5.5 to day 8 are highly biased toward the MET Region, with 94% being in the top 50% of cells with respect to the proportion of descendants in this region (FIG. 29A).

[0298]    We tested whether expressing Obox6 together with OKSM during days 0-8 can boost reprogramming efficiency. We infected our secondary MEFs with a Dox-inducible lentivirus

carrying either Obox6, the known pluripotency factor Zfp42 (Rajkovic et al., 2002; Shi et al., 2006), or no insert as a negative control. Both Obox6 and Zpf42 increased reprogramming efficiency of secondary MEFs by ~2-fold in 2i and even more so in serum, with the result confirmed in multiple independent experiments (FIGs. 29B, 29C, and 36A-36F). Assays in primary MEFs showed similar increases in reprogramming efficiency (FIGs. 26A-36F).

[0299] Together, these computational and experimental results suggested that the role of Obox6 in reprogramming merits further study.

[0300] In addition, we identified GDF9 that can significantly booster reprogramming efficiency. We added GDF9 to the medium from day 8. We observed more Oct4-GFP positive colonies (iPSCs) **(FIG. 37).** We also confirmed that we saw more iPSCs after adding GDF9 by scRNA sequencing.

[0301] **FIG. 38** shows adding GDF9 to the medium resulted in more iPSCs.

[0302] **Discussion**

[0303] Understanding the trajectories of cellular differentiation was important for studying development and for regenerative medicine. Large-scale, single-cell profiling had dramatically advanced progress toward this goal. However, the challenge of turning snapshots from single-cell profiling into accurate movies of cellular differentiation had not yet been fully solved. Here, we described two resources for the scientific community: a new analytical approach to reconstructing trajectories, and a massive dataset of 315,000 cells from time courses of classic reprogramming from fibroblasts to iPSCs under two conditions. By applying the approach to the dataset, we shed new light on this well-studied problem, and provide a template for future studies in other systems.

[0304] An optimal transport framework to model cell differentiation

[0305] Waddington-OT provided an inherently probabilistic approach that described transitions between time points in terms of stochastic couplings, derived from a modified version of the mathematical method of optimal transport. The approach yielded a natural concept of trajectories in terms of ancestor and descendant distributions for any set of cells at a given time point. This allowed us gracefully to recover, for example, branching events (by the emergence of bimodality in the descendant distribution) or shared vs. distinct ancestry between two cell sets (by convergence of the ancestor distributions) (FIGs. 23C-23E). The trajectories can then be

used to study differentiation between classes of cells at different times, including creating regulatory models to infer TFs involved in activating specific gene-expression programs. Our model did not impose strict structural constraints a priori on the nature of these processes, allowing for gradual changes over time rather than sharp discrete transitions. Moreover, OT can be applied to even a single pair of time points (if the transition is expected to be sufficiently smooth) and thus can be helpful even for a small experimental scheme. Indeed, we validated Waddington-OT by testing its ability to accurately infer cellular distributions at held-out intermediate time points and by showing that its results are robust across wide variation in parameters.

[0306]    Waddington-OT differed from previous approaches because it (i) did not attempt to force cells onto a simple branching graph, (ii) made explicit use of temporal information, and (iii) allowed for cell growth and death. We also found that Waddington-OT appeared to perform better than several graph-based methods, at least for studying cellular reprogramming from fibroblasts to iPSCs (FIGs. 35A-35B, Methods). Specifically, the widely and successfully used program Monocle2 (Qiu et al., 2017) generated trajectories that a) were inconsistent with known information about time (day 18 stromal cells give rise to essentially all cells after day 0), and b) placed neural and iPS together as one terminal state. The recently developed program URD (Farrell et al., 2018) could avoid the latter problem by finding trajectories to specific cell sets of interest, but a) it generated trajectories which contradicted the gradual MET/Stromal fate specification we saw in our data (in URD, the stromal branch completely diverges at day 0.5), and b) the binary nature of the URD tree could not capture the multifurcation of neural, iPS, trophoblast and epithelial cells from MET.

[0307]    Tracking cell differentiation trajectories and fates in a diverse reprogramming landscape

[0308]    Although the reprogramming of fibroblasts to iPSCs had been intensively studied since it was discovered by Yamanaka, our study shedded new light on the process - providing insights that could only be obtained from large-scale single-cell profiles across dense time courses matched with appropriate analytical methods.

[0309]    First, single-cell profiling with large numbers of cells along a dense time course revealed remarkable and unappreciated diversity in the reprogramming landscape, with large

classes of cells having distinct biological programs, related to distinct states and tissues (pluripotency, trophoblasts, neural tissue, epithelium and stroma). In earlier studies based on bulk RNA analysis, we and others had detected expression of individual genes characteristic of various lineages during reprogramming. (Mikkelsen et al., 2008; O'Malley et al., 2013; Parenti et al., 2016). Studying these classes in greater detail, we found a tremendous richness of cells expressing distinct gene-expression programs associated with specific cell types in vivo. Examples included: (i) within iPSC-like cells, programs associated with 2-, 4-, 8-, 16-, and 32-cell stage embryos; (ii) within extra-embryonic-like cells, programs associated with several distinct types of trophoblasts and programs associated with primitive endoderm (at one time point); (iii) within neural-like cells, programs associated with astrocytes, oligodendrocytes, and neurons, as well as specific subprograms associated with excitatory and inhibitory neurons; and (iv) within stromal-like cells, distinct programs associated with a wider range of stromal cells than simply MEFs. Further work will be needed to determine the extent to which these cell types adopt the full identity of natural cell types that they resemble.

[0310]     This dramatic diversity raised several key questions that Waddington-OT has helped us begin to address, including: (1) What are the differentiation and fate trajectories that span these cell subsets? When do they diverge, from which ancestors, and to which cells do they give rise? (2) What cell intrinsic regulatory mechanisms may drive each fate, especially transcription factors? (3) What might be the role of cells of different types at cross-communicating and supporting across differentiation trajectories and fates in general, and for the iPSC fate in particular?

[0311]     First, our trajectory and regulatory analysis allowed us to build a model that synthesizes a comprehensive view of the differentiation and fate trajectories in the landscape (FIG. 29D). We highlighted several key fate decisions, in a manner that allowed us to understand their gradual and continuous nature. During the initial phase of reprogramming, cells began to diverge in two alternative directions: toward stromal cells or toward an MET state (FIG. 29D, blue and purple). In the MET direction this divergence was not sharp: although some ancestors exhibited biases in cell fate as early as day 1.5, cells continued to 'switch' their fate preference from MET to Stromal up to day 8 (FIGs. 29A-29D, arrows from purple to blue zones). In contrast, the Stromal Region was terminal, and the reverse phenomenon was not seen by our

model. Following withdrawal of dox at day 8, the cells in the MET state gave rise to iPSC-, trophoblast-, neural-, and epithelial-like cells. We found no evidence that particular cells had biases towards any of these fates before this point, whereas our analysis clearly distinguished the biases that arise once dox was withdrawn. The ancestors that would lead to iPSCs were distinguished early after withdrawal (day 9), and they passed through a narrow bottleneck towards iPSC. Conversely, other cells in the MET region first assumed an epithelial-like state, with ancestors leading to trophoblasts vs. neural cells (in serum) becoming distinguished a few days later. Within neural cells (in serum) and trophoblast-like cells (in both conditions), there was substantial additional divergence, which we could at times trace to additional divergence between ancestors at later time point. For example, the radial glial population expressing GdflO RG at day 13.5 was enriched for ancestors of later emerging neuron-like cells.

[0312]     Second, by characterizing events that occurred along the trajectory toward any cell class, we identified TFs that might drive subsequent fates (FIG. 29D). Along the path toward pluripotency, we readily rediscovered known TFs, validating our approach, but also identified several new TFs not previously implicated in the process. We tested one such new TF, Obox6, which was associated with a strong bias toward MET early and toward pluripotency late; we found that forced expression of Obox6 increased reprogramming efficiency. Along paths to other fates, we similarly rediscovered TFs known to play a role in differentiation of the corresponding cells in vivo, as well as identified TFs that were expressed in the target cell type but had not been implicated in differentiation per se.

[0313]     Third, contemporaneous expression of receptor-ligand pairs across cell subsets highlighted potential paracrine interactions between the stromal cells and the iPSC-like, neural-like and trophoblast-like cells, which might play key roles in the initial differentiation and maintenance of these cell types. If many of these potential interactions could be validated by experimental assays, it would suggest that efficient reprogramming requires alternative cell types, or the exogenous replacement of the factors they supply. Additionally, single-cell expression revealed likely regions of genomic aberration; the frequency of such events was significantly higher in our trophoblast and stromal cells, consistent with known biological properties of these cell types.

[0314]     Prospects for models and studies of differentiation and development

**[0315]**   Our method captured several key aspects of cellular differentiation and, importantly, can be extended to capture additional features. First, the framework currently assumed that a cell's trajectory depended only on its current gene-expression levels. As it became possible to perform single-cell profiling simultaneously for gene expression and epigenomic states, one can readily incorporate both types of information. Second, our framework for learning regulatory models assume that trajectories are cell autonomous, but may be extended to incorporate intercellular interactions, such as the potential paracrine signaling postulated here, by using optimal transport for interacting particles (Ambrosio et al., 2008; Santambrogio, 2015) (STAR Methods). Third, various methods are being developed for obtaining lineage information about cells, based on the introduction of barcodes at discrete time points or even continuously (Frieda et al., 2017; McKenna et al., 2016). Barcodes can be used to recognize cells that descend from a recent common ancestor cell, but do not currently directly reveal the full gene-expression state of the ancestral cell. However, they can be incorporated into our optimal-transport framework to improve the inference of ancestral cell states. Finally, our method can be refined to analyze multiple time points simultaneously, rather than just pairs of consecutive time points; this can be particularly useful for situations where the number of cells at different time points varies significantly.

**[0316]**   In summary, our findings indicated that the process of reprogramming fibroblasts to iPSCs unleashed a much wider range of developmental programs and subprograms than previously characterized.

**[0317]**   **References**

•   Aaronson, Y., Livyatan, L, Gokhman, D., and Meshorer, E. (2016). Systematic identification of gene family regulators in mouse and human embryonic stem cells. Nucleic Acids Research 44, 4080-4089.

•   Daniel et al., (2018). A revised airway epithelial hierarchy includes CFTR-expressing ionocytes. Nature 2018, accepted.

•   Ambrosio, L., Gigli, N., and Savare, G. (2008). Gradient flows: in metric spaces and in the space of probability measures (Springer Science & Business Media).

•   Bastian, M., Heymann, S., Jacomy, M., et al. (2009). Gephi: an open source software for exploring and manipulating networks. Icwsm, 8:361-362.

• Bendall, S.C., Davis, K.L., Amir, E.-a.D., Tadmor, M.D., Simonds, E.F., Chen, T.J., Shenfeld, D.K., Nolan, G.P., and Pe'er, D. (2014). Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. Cell 157, 714-725.

• Beygelzimer, A., Kakadet, S., Langford, J., Arya, S., Mount, D., Li, S., and Li, M. S. (2015). Package FNN.

• Boheler, K.R. (2009). Stem cell pluripotency: a cellular trait that depends on transcription factors, chromatin state and a checkpoint deficient cell cycle. Journal of cellular physiology 221, 10-17.

• Briggs, J.A., Weinreb, C, Wagner, D.E., Megason, S., Peshkin, L., Kirschner, M.W., and Klein, A.M. (2018). The dynamics of gene expression in vertebrate embryogenesis at single-cell resolution. Science.

• Buganim, Y., Faddah, D.A., Cheng, A.W., Itskovich, E., Markoulaki, S., Ganz, K., Klemm, S.L., van Oudenaarden, A., and Jaenisch, R. (2012). Single-cell expression analyses during cellular reprogramming reveal an early stochastic and a late hierarchic phase. Cell 150, 1209-1222.

• Cacchiarelli, D., Trapnell, C, Ziller, M.J., Soumillon, M., Cesana, M., Karnik, R., Donaghey, J., Smith, Z.D., Ratanasirintrawoot, S., Zhang, X., Ho Sui, S.J., Wu, Z., Akopian, V., Gifford, C.A., Doench, J., Rinn, J.L., Daley, G.Q., Meissner, A., Lander, E.S., and Mikkelsen, T. (2015). Integrative Analyses of Human Reprogramming Reveal Dynamic Nature of Induced Pluripotency. Cell 162.

• Cannoodt, R., Saelens, W., Sichien, D., Tavernier, S., Janssens, S., Guilliams, M., Lambrecht, B.N., De Preter, K., and Saeys, Y. (2016). SCORPIUS improves trajectory inference and identifies novel modules in dendritic cell development. bioRxiv.

• Chen, E.Y., Tan, CM., Kou, Y., Duan, Q., Wang, Z., Meirelles, G.V., Clark, NR., and Ma'ayan, A. (2013). Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. BMC Bioinformatics 14, 128.

• Chen, Q., Zhang, M., Li, Y, Xu, D., Wang, Y, Song, A., Zhu, B., Huang, Y., and Zheng, J.C. (2015). CXCR7 Mediates Neural Progenitor Cells Migration to CXCL12 Independent of CXCR4. Stem cells (Dayton, Ohio) 33, 2574-2585.

- Chizat, L., Peyre, G., Schmitzer, B., and Vialard, F.-X. (2017). Scaling algorithms for unbalanced transport problems. arXiv preprint arXiv:160705816v2.

- Coppe, J.-P., Desprez, P.-Y., Krtolica, A., and Campisi, J. (2010). The senescence-associated secretory phenotype: the dark side of tumor suppression. Annual Review of Pathological Mechanical Disease 5, 99-1 18.

- Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transport. Paper presented at: Advances in neural information processing systems.

- Elson, G.C., Lelievre, E., Guillet, C , Chevalier, S., Plun-Favreau, H., Froger, J., Suard, I., de Coignac, A.B., Delneste, Y., and Bonnefoy, J.-Y. (2000). CLF associates with CLC to form a functional heteromeric ligand for the CNTF receptor complex. Nature neuroscience 3, 867.

- Falco, G , Lee, S.L., Stanghellini, I., Bassey, U.C., Hamatani, T., and Ko, M.S. (2007). Zscan4: a novel gene expressed exclusively in late 2-cell embryos and embryonic stem cells. Developmental biology 307, 539-550.

- Farrell, J.A, Wang, Y., Riesenfeld, S.J., Shekhar, K., Regev, A , and Schier, A.F. (2018). Single-cell reconstruction of developmental trajectories during zebrafish embryogenesis. Science.

- Fincher, C.T., Wurtzel, O., de Hoog, T., Kravarik, K.M., and Reddien, P.W. (2018). Cell type transcriptome atlas for the planarian <em>Schmidtea mediterranea</em>. Science.

- Florio, M., and Huttner, W.B. (2014). Neural progenitors, neurogenesis and the evolution of the neocortex. Development 141, 2182-2194.

- Fonseca, E.T.d., Man?anares, A.C.F., Ambr®,Esio, C.E., and Miglino, M.A.I. (2013). Review point on neural stem cells and neurogenic areas of the central nervous system. Open Journal of Animal Sciences Vol.03No.03, 6.

- Frieda, K.L., Linton, J.M., Hormoz, S., Choi, J., Chow, K.-H.K., Singer, Z.S., Budde, M.W., Elowitz, M.B., and Cai, L. (2017). Synthetic recording and in situ readout of lineage information in single cells. Nature 541, 107.

- Froidure, A., Marchal-Duval, E., Ghanem, M., Gerish, L., Jaillet, M., Crestani, B., and Mailleux, A. (2016). Mesenchyme associated transcription factor PRRXl: A key regulator of IPF fibroblast. European Respiratory Journal 48.

• Gegenschatz-Schmid, K., Verkauskas, G., Demougin, P., Bilius, V., Dasevicius, D., Stadler, M.B., and Hadziselimovic, F. (2017). DMRTC2, PAX7, BRACHYURY/T and TERT Are Implicated in Male Germ Cell Development Following Curative Hormone Treatment for Cryptorchidism-Induced Infertility. Genes 8, 267.

• Goolam, M., Scialdone, A., Graham, S.J.L., Macaulay, I.C., Jedrusik, A., Hupalowska, A., Voet, T., Marioni, J.C., and Zernicka-Goetz, M. (2016). Heterogeneity in Oct4 and Sox2 Targets Biases Cell Fate in 4-Cell Mouse Embryos. Cell 165, 61-74.

• Gouti, M., Briscoe, J., and Gavalas, A. (201 1). Anterior Hox genes interact with components of the neural crest specification network to induce neural crest fates. Stem cells (Dayton, Ohio) 29, 858-870.

• Haghverdi, L., Buettner, F., and Theis, F.J. (2015). Diffusion maps for high-dimensional single-cell analysis of differentiation data. Bioinformatics 31, 2989-2998.

• Haghverdi, L., Buettner, M., Wolf, F.A., Buettner, F., and Theis, F.J. (2016). Diffusion pseudonyme robustly reconstructs lineage branching. bioRxiv, 041384.

• Han, X., Wang, R., Zhou, Y., Fei, L., Sun, H., Lai, S., Saadatpour, A., Zhou, Z., Chen, H., Ye, F., et al. (2018). Mapping the Mouse Cell Atlas by Microwell-Seq. Cell 172, 1091-1107.el017.

• Hayashi, Y., Hsiao, E.C., Sami, S., Lancero, M., Schlieve, C.R., Nguyen, T., Yano, K., Nagahashi, A., Ikeya, M., Matsumoto, Y., et al. (2016). BMP-SMAD-ID promotes reprogramming to pluripotency by inhibiting pl6/INK4A-dependent senescence. Proceedings of the National Academy of Sciences of the United States of America 113, 13057-13062.

• Hou, P., Li, Y., Zhang, X., Liu, C , Guan, J., Li, H., Zhao, T., Ye, J., Yang, W., Liu, K., et al. (2013). Pluripotent Stem Cells Induced from Mouse Somatic Cells by Small-Molecule Compounds. Science 341, 651-654.

• Hu, G , Kim, J., Xu, Q., Leng, Y., Orkin, S.H., and Elledge, S.J. (2009). A genome-wide RNAi screen identifies a new transcriptional module required for self-renewal. Genes & development 23, 837-848.

• Hussein, S.M., Puri, M.C., Tonge, P.D., Benevento, M., Corso, A.J., Clancy, J.L., Mosbergen, R., Li, M., Lee, D.-S., and Cloonan, N. (2014). Genome-wide characterization of the routes to pluripotency. Nature 516, 198.

• Jacomy, M., Venturini, T., Heymann, S., and Bastian, M. (2014). ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. PloS one 9, e98679.

• Jeon, H., Waku, T., Azami, T., Khoa le, T.P., Yanagisawa, J., Takahashi, S., and Ema, M. (2016). Comprehensive Identification of Kruppel-Like Factor Family Members Contributing to the Self-Renewal of Mouse Embryonic Stem Cells and Cellular Reprogramming. PloS one 11, e0150715.

• Jukkola, T., Lahti, L., Naserke, T., Wurst, W., and Partanen, J. (2006). FGF regulated gene-expression and neuronal differentiation in the developing midbrain-hindbrain region. Developmental biology 297, 141-157.

• Kan, L., Israsena, N., Zhang, Z., Hu, M., Zhao, L.R., Jalali, A., Sahni, V., and Kessler, J.A. (2004). Sox1 acts through multiple independent pathways to promote neurogenesis. Developmental biology 269, 580-594.

• Kantorovitch, L. (1958). On the Translocation of Masses. Management Science 5, 1-4.

• Kester, L., and van Oudenaarden, A. (2018). Single-Cell Transcriptomics Meets Lineage Tracing. Cell Stem Cell.

• Kidder, B.L., and Palmer, S. (2010). Examination of transcriptional networks reveals an important role for TCFAP2C, SMARCA4, and EOMES in trophoblast stem cell maintenance. Genome Res 20, 458-472.

• Kim, D.H., Marinov, G.K., Pepke, S., Singer, Z.S., He, P., Williams, B., Schroth, G.P., Elowitz, M.B., and Wold, B.J. (2015). Single-cell transcriptome analysis reveals dynamic changes in IncRNA expression during reprogramming. Cell stem cell 16, 88-101.

• Klein, A.M., Mazutis, L., Akartuna, L, Tallapragada, N., Veres, A., Li, V., Peshkin, L., Weitz, D.A., and Kirschner, M.W. (2015). Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. Cell 161, 1187-1201.

• Kolodziejczyk, Aleksandra A., Kim, Jong K., Tsang, Jason C, Ilicic, T., Henriksson, J., Natarajan, Kedar N., Tuck, Alex C, Gao, X., Buhler, M., Liu, P., et al. (2015). Single Cell RNA-Sequencing of Pluripotent States Unlocks Modular Transcriptional Variation. Cell Stem Cell 17, 471-485.

•      Kumar, R.M., Cahan, P., Shalek, A.K., Satija, R., Jay DaleyKeyser, A., Li, H., Zhang, J., Pardee, K., Gennert, D., Trombetta, J.J., et al. (2014). Deconstructing transcriptional heterogeneity in pluripotent stem cells. Nature 516, 56.

•      Latos, P.A., and Hemberger, M. (2016). From the stem of the placental tree: trophoblast stem cells and their progeny. Development 143, 3650-3660.

•      Lattin, J.E., Schroder, K., Su, A.I., Walker, J.R, Zhang, J., Wiltshire, T., Saijo, K., Glass, C.K., Hume, D.A., Kellie, S., et al. (2008). Expression analysis of G Protein-Coupled Receptors in mouse macrophages. Immunome research 4, 5.

•      Lazarov, O., Mattson, M.P., Peterson, D.A., Pimplikar, S.W., and van Praag, H. (2010). When neurogenesis encounters aging and disease. Trends in neurosciences 33, 569-579.

•      Le´onard, C. (2014). A survey of the schröndinger problem and some of its connections with optimal transport. Discrete and Continuous Dynamical Systems - Series A (DCDS-A), 34(4):1533-1574.

•      Li, R., Liang, J., Ni, S., Zhou, T., Qing, X., Li, H., He, W., Chen, J., Li, F., Zhuang, Q., et al. (2010). A mesenchymal-to-epithelial transition initiates and is required for the nuclear reprogramming of mouse fibroblasts. Cell Stem Cell 7, 51-63.

•      Li, W.-Z., Wang, Z.-W., Chen, L.-L., Xue, H.-N., Chen, X., Guo, Z.-K., and Zhang, Y. (2015). Hesxl enhances pluripotency by working downstream of multiple pluripotency-associated signaling pathways. Biochemical and Biophysical Research Communications 464, 936-942.

•      Liang, H., Zhang, Q., Lu, J., Yang, G, Tian, N., Wang, X., Tan, Y, and Tan, D. (2016). MSX2 Induces Trophoblast Invasion in Human Placenta. PloS one 11, e0153656.

•      Lim, L.S., Loh, Y.H., Zhang, W., Li, Y., Chen, X, Wang, Y, Bakre, M., Ng, H.H, and Stanton, L.W. (2007). Zic3 is required for maintenance of pluripotency in embryonic stem cells. Molecular biology of the cell 18, 1348-1358.

•      Lin, J., Khan, M., Zapiec, B., and Mombaerts, P. (2016). Efficient derivation of extraembryonic endoderm stem cell lines from mouse postimplantation embryos. Scientific reports 6, 39457.

•     Liu, J., Han, Q., Peng, T., Peng, M , Wei, B., Li, D., Wang, X., Yu, S., Yang, J., Cao, S., et al. (2015). The oncogene c-Jun impedes somatic cell reprogramming. Nature cell biology 17, 856-867.

•     Liu, L.L., Brumbaugh, J., Bar-Nur, O., Smith, Z., Stadtfeld, M., Meissner, A., Hochedlinger, K., and Michor, F. (2016). Probabilistic Modeling of Reprogramming to Induced Pluripotent Stem Cells. Cell reports 17, 3395-3406.

•     Ma, G.T., Roth, M.E., Groskopf, J.C., Tsai, F.Y., Orkin, S.H., Grosveld, F., Engel, J.D., and Linzer, D.I. (1997). GATA-2 and GATA-3 regulate trophoblast-specific gene expression in vivo. Development 124, 907-914.

•     Macfarlan, T.S., Gifford, W.D., Driscoll, S., Lettieri, K., Rowe, H.M., Bonanomi, D., Firth, A., Singer, O., Trono, D., and Pfaff, S.L. (2012). Embryonic stem cell potency fluctuates with endogenous retrovirus activity. Nature 487, 57-63.

•     Macosko, E.Z., Basu, A., Satija, R., Nemesh, J., Shekhar, K., Goldman, M., Tirosh, I., Bialas, A.R., Kamitaki, N , and Martersteck, E.M. (2015). Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. Cell 161, 1202-1214.

•     Marco, E., Karp, R.L., Guo, G , Robson, P., Hart, A.H., Trippa, L., and Yuan, G.C. (2014). Bifurcation analysis of single-cell gene expression data reveals epigenetic landscape. Proceedings of the National Academy of Sciences of the United States of America 111, E5643-5650.

•     Matsumoto, H., and Kiryu, H. (2016). SCOUP: a probabilistic model based on the Ornstein-Uhlenbeck process to analyze single-cell expression data during differentiation. BMC Bioinformatics 17, 232.

•     McKenna, A., Findlay, G.M., Gagnon, J.A., Horwitz, M.S., Schier, A.F., and Shendure, J. (2016). Whole-organism lineage tracing by combinatorial and cumulative genome editing. Science 353, aaf7907.

•     Mertins, P., Przybylski, D., Yosef, N., Qiao, J., Clauser, K., Raychowdhury, R., Eisenhaure, T.M., Maritzen, T., Haucke, V., Satoh, T., et al. (2017). An Integrative Framework Reveals Signaling-to-Transcription Events in Toll-like Receptor Signaling. Cell reports 19, 2853-2866.

• Messina, G., Biressi, S., Monteverde, S., Magli, A., Cassano, M., Perani, L., Roncaglia, E., Tagliafico, E., Starnes, L., Campbell, C.E., et al. (2010). Nfix regulates fetal-specific transcription in developing skeletal muscle. Cell 140, 554-566.

• Mikkelsen, T.S., Hanna, J., Zhang, X., Ku, M., Wernig, M., Schorderet, P., Bernstein, B.E., Jaenisch, R., Lander, E.S., and Meissner, A. (2008). Dissecting direct reprogramming through integrative genomic analysis. Nature 454, 49.

• Ming, G.L., and Song, H. (201 1). Adult neurogenesis in the mammalian brain: significant answers and significant questions. Neuron 70, 687-702.

• Mosteiro, L., Pantoja, C , Alcazar, N., Marion, R.M., Chondronasiou, D., Rovira, M., Fernandez-Marcos, P.J., Munoz-Martin, M., Blanco-Aparicio, C , and Pastor, J. (2016). Tissue damage and senescence provide critical signals for cellular reprogramming in vivo. Science 354, aaf4445.

• Nakashima, K., Wiese, S., Yanagisawa, M., Arakawa, H., Kimura, N , Hisatsune, T., Yoshida, K., Kishimoto, T., Sendtner, M., and Taga, T. (1999). Developmental requirement of gpl30 signaling in neuronal survival and astrocyte differentiation. The Journal of neuroscience : the official journal of the Society for Neuroscience 19, 5429-5434.

• Nelson, A.C., Mould, A.W., Bikoff, E.K., and Robertson, E.J. (2016). Single-cell RNA-seq reveals cell type-specific transcriptional signatures at the maternal-foetal interface during pregnancy. Nat Commun 7, 11414.

• O'Malley, J., Skylaki, S., Iwabuchi, K.A., Chantzoura, E., Ruetz, T., Johnsson, A., Tomlinson, S.R., Linnarsson, S., and Kaji, K. (2013). High resolution analysis with novel cell-surface markers identifies routes to iPS cells. Nature 499, 88.

• Ocana, O.H., Corcoles, R , Fabra, A., Moreno-Bueno, G., Acloque, H., Vega, S., Barrallo-Gimeno, A., Cano, A., and Nieto, M.A. (2012). Metastatic colonization requires the repression of the epithelial-mesenchymal transition inducer Prrxl. Cancer cell 22, 709-724.

• Parast, M.M., Yu, H., Ciric, A , Salata, M.W., Davis, V., and Milstone, D.S. (2009). PPARgamma regulates trophoblast proliferation and promotes labyrinthine trilineage differentiation. PloS one 4, e8055.

• Parenti, A., Halbisen, M.A., Wang, K., Latham, K., and Ralston, A. (2016). OSKM induce extraembryonic endoderm stem cells in parallel to induced pluripotent stem cells. Stem cell reports 6, 447-455.

• Park, M., Lee, Y., Jang, H., Lee, O.H., Park, S.W., Kim, J.H., Hong, K., Song, H, Park, S.P., Park, Y.Y., et al. (2016). SOHLH2 is essential for synaptonemal complex formation during spermatogenesis in early postnatal mouse testes. Scientific reports 6, 20980.

• Pasque, V., Tchieu, J., Karnik, R, Uyeda, M., Dimashkie, A.S., Case, D., Papp, B., Bonora, G., Patel, S., and Ho, R. (2014). X chromosome reactivation dynamics reveal stages of reprogramming to pluripotency. Cell 159, 1681-1697.

• Patel, A.P., Tirosh, L, Trombetta, J.J., Shalek, A.K., Gillespie, S.M., Wakimoto, H., Cahill, D.P., Nahed, B.V., Curry, W.T., Martuza, R.L., et al. (2014). Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. Science (New York, NY) 344, 1396-1401.

• Pei, J., and Grishin, N.V. (2012). Unexpected diversity in Shisa-like proteins suggests the importance of their roles as transmembrane adaptors. Cellular signalling 24, 758-769.

• Plass, M., Solana, J., Wolf, F.A., Ayoub, S., Misios, A., Glazar, P., Obermayer, B., Theis, F.J., Kocks, C, and Rajewsky, N. (2018). Cell type atlas and lineage tree of a whole complex animal by single-cell transcriptomics. Science.

• Polo, J.M., Anderssen, E., Walsh, R.M., Schwarz, B.A., Nefzger, CM., Lim, S.M., Borkent, M., Apostolou, E., Alaei, S., and Cloutier, J. (2012). A molecular roadmap of reprogramming somatic cells into iPS cells. Cell 151, 1617-1632.

• Porpiglia, E., Samusik, N., Van Ho, A. T., Cosgrove, B. D., Mai, T., Davis, K. L., Jager, A., Nolan, G. P., Bendall, S. C, Fantl, W. J., et al. (2017). High-resolution myogenic lineage mapping by single-cell mass cytometry. Nature Cell Biol., 19:558-567.


• Qiu, X, Mao, Q., Tang, Y, Wang, L., Chawla, R., Pliner, H., and Trapnell, C. (2017). Reversed graph embedding resolves complex single-cell developmental trajectories. bioRxiv, 110668.

• Rajkovic, A., Yan, C, Yan, W., Klysik, M., and Matzuk, M.M. (2002). Obox, a Family of Homeobox Genes Preferentially Expressed in Germ Cells. Genomics 79, 711-717.

- Ralston, A., Cox, B.J., Nishioka, N., Sasaki, H., Chea, E., Rugg-Gunn, P., Guo, G., Robson, P., Draper, J.S., and Rossant, J. (2010). Gata3 regulates trophoblast development downstream of Tead4 and in parallel to Cdx2. Development 137, 395-403.

- Ramskold, D., Luo, S., Wang, Y.-C, Li, R., Deng, Q., Faridani, O.R., Daniels, G.A., Khrebtukova, I., Loring, J.F., Laurent, L.C., et al. (2012). Full-Length mRNA-Seq from single cell levels of RNA and individual circulating tumor cells. Nature biotechnology 30, 777-782.

- Rashid, S., Kotton, D.N., and Bar-Joseph, Z. (2017). TASIC: determining branching models from time series single cell data. Bioinformatics 33, 2504-2512.

- Richard Jordan, D. K. and Otto, F. (1998). The variational formulation of the fokker. SIAM J. Math. Anal., 29(1): 1-17.

- Rostom, R., Svensson, V., Teichmann, S., and Kar, G. (2017). Computational approaches for interpreting scRNA-seq data. FEBS letters.

- Sakakibara, S., Nakamura, Y., Satoh, H., and Okano, H. (2001). Rna-binding protein Musashi2: developmentally regulated expression in neural precursor cells and subpopulations of neurons in mammalian CNS. The Journal of neuroscience : the official journal of the Society for Neuroscience 21, 8091-8107.

- Samusik, N., Good, Z., Spitzer, M. H., Davis, K. L., and Nolan, G. P. (2016). Automated mapping of phenotype space with single-cell data. Nature methods, 13:493-496.

- Sansom, S.N., Griffiths, D.S., Faedo, A., Kleinjan, D.J., Ruan, Y., Smith, J., van Heyningen, V., Rubenstein, J.L., and Livesey, F.J. (2009). The level of the transcription factor Pax6 is essential for controlling the balance between neural stem cell self-renewal and neurogenesis. PLoS genetics 5, el00051 1.

- Santambrogio, F. (2015). Optimal transport for applied mathematicians. Birkauser, NY, 99-102.

- Satija, R., Farrell, J.A., Gennert, D., Schier, A.F., and Regev, A. (2015). Spatial reconstruction of single-cell gene expression data. Nature Biotechnology 33, 495.

- Scott, I.C., Anson-Cartwright, L., Riley, P., Reda, D., and Cross, J.C. (2000). The HANDl basic helix-loop-helix transcription factor regulates trophoblast differentiation via multiple mechanisms. Molecular and cellular biology 20, 530-541.

•      Setty, M., Tadmor, M.D., Reich-Zeliger, S., Angel, O., Salame, T.M., Kathail, P., Choi, K., Bendall, S., Friedman, N., and Pe'er, D. (2016). Wishbone identifies bifurcating developmental trajectories from single-cell data. Nature biotechnology 34, 637-645.

•      Shalek, A.K., Satija, R., Adiconis, X., Gertner, R.S., Gaublomme, J.T., Raychowdhury, R., Schwartz, S., Yosef, N., Malboeuf, C., Lu, D., et al. (2013). Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. Nature 498, 236.

•      Shi, W., Wang, H., Pan, G., Geng, Y., Guo, Y., and Pei, D. (2006). Regulation of the pluripotency marker Rex-1 by Nanog and Sox2. J Biol Chem 281, 23319-23325.

•      Shu, I., Wu, C., Wu, Y., Li, Z., Shao, S., Zhao, W., Tang, X., Yang, H., Shen, L., Zuo, X., et al. (2013). Induction of pluripotency in mouse somatic cells with lineage specifiers. Cell 153, 963-975.

•      Simmons, D.G., and Cross, J.C. (2005). Determinants of trophoblast lineage and cell subtype specification in the mouse placenta. Developmental biology 284, 12-24.

•      Simmons, D.G, Natale, D.R., Begay, V., Hughes, M., Leutz, A., and Cross, J.C. (2008). Early patterning of the chorion leads to the trilaminar trophoblast cell structure in the placental labyrinth. Development 135, 2083-2091.

•      Stadtfeld, M., Maherali, N., Borkent, M., and Hochedlinger, K. (2010). A reprogrammable mouse strain from gene-targeted embryonic stem cells. Nature methods 7, 53-55.

•      Street, K., Risso, D., Fletcher, R.B., Das, D., Ngai, J., Yosef, N, Purdom, E., and Dudoit, S. (2017). Slingshot: Cell lineage and pseudotime inference for single-cell transcriptomics. bioRxiv.

•      Takahashi, K., and Yamanaka, S. (2006). Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors. cell 126, 663-676.

•      Takahashi, K., and Yamanaka, S. (2016). A decade of transcription factor-mediated reprogramming to pluripotency. Nature Reviews Molecular Cell Biology 17, 183.

•      Takaishi, M., Tarutani, M., Takeda, J., and Sano, S. (2016). Mesenchymal to Epithelial Transition Induced by Reprogramming Factors Attenuates the Malignancy of Cancer Cells. PloS one 11, e0156904.

- Tanay, A., and Regev, A. (2017). Scaling single-cell genomics from phenomenology to mechanism. Nature 541, 331-338.

- Tang, F., Barbacioru, C, Wang, Y., Nordman, E., Lee, C, Xu, N., Wang, X., Bodeau, J., Tuch, B.B., Siddiqui, A., et al. (2009). mRNA-Seq whole-transcriptome analysis of a single cell. Nature Methods 6, 377.

- Tasic, B., Menon, V., Nguyen, T.N., Kim, T.K., Jarsky, T., Yao, Z., Levi, B., Gray, L.T., Sorensen, S.A., Dolbeare, T., et al. (2016). Adult mouse cortical cell taxonomy revealed by single cell transcriptomics. NatNeurosci 19, 335-346.

- Tirosh, I., Venteicher, A.S., Hebert, C, Escalante, L.E., Patel, A.P., Yizhak, K., Fisher, J.M., Rodman, C, Mount, C, and Filbin, M.G. (2016). Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma. Nature 539, 309-313.

- Tonge, P.D., Corso, A.J., Monetti, C, Hussein, S.M., Puri, M.C., Michael, LP., Li, M., Lee, D.-S., Mar, J.C., and Cloonan, N. (2014). Divergent reprogramming routes lead to alternative stem-cell states. Nature 516, 192-197.

- Trapnell, C, Cacchiarelli, D., Grimsby, J., Pokharel, P., Li, S., Morse, M., Lennon, N.J., Livak, K.J., Mikkelsen, T.S., and Rinn, J.L. (2014). The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nature biotechnology 32, 381-386.

- Ueno, M., Lee, L.K., Chhabra, A., Kim, Y.J., Sasidharan, R., Van Handel, B., Wang, Y., Kamata, M., Kamran, P., Sereti, K.-L, et al. (2013). c-Met-dependent multipotent labyrinth trophoblast progenitors establish placental exchange interface. Developmental cell 27, 373-386.

- Vandercappellen, J., Van Damme, J., and Struyf, S. (2008). The role of CXC chemokines and their receptors in cancer. Cancer letters 267, 226-244.

- Villani, C. (2008). Optimal transport: old and new, Vol 338 (Springer Science & Business Media).

- Waddington, C.H. (1936). How animals develop (New York).

- Waddington, C.H. (1957). The strategy of the genes; a discussion of some aspects of theoretical biology (London, Allen & Unwin [1957]).

- Wagner, A., Regev, A., and Yosef, N. (2016). Revealing the vectors of cellular identity with single-cell genomics. Nat Biotech 34, 1145-1 160.

• Wagner, D.E., Weinreb, C , Collins, Z.M., Briggs, J.A., Megason, S.G., and Klein, A.M. (2018). Single-cell mapping of gene expression landscapes and lineage in the zebrafish embryo. Science.

• Watanabe, Y., Stanchina, L., Lecerf, L., Gacem, N., Conidi, A., Baral, V., Pingault, V., Huylebroeck, D., and Bondurand, N. (2017). Differentiation of Mouse Enteric Nervous System Progenitor Cells Is Controlled by Endothelin 3 and Requires Regulation of Ednrb by SOX10 and ZEB2. Gastroenterology 152, 1139-1 150.el 134.

• Weinreb, C , Wolock, S., and Klein, A. (2016). SPRING: a kinetic interface for visualizing high dimensional single-cell expression data. bioRxiv.

• Weinreb, C , Wolock, S., Tusi, B.K., Socolovsky, M., and Klein, A.M. (2017). Fundamental limits on dynamic inference from single cell snapshots. bioRxiv.

• Welch, J.D., Hartemink, A.J., and Prins, J.F. (2016). SLICER: inferring branched, nonlinear cellular trajectories from single cell RNA-seq data. Genome Biology 17, 106.

• Whiteman, EX., Fan, S., Harder, J.L., Walton, K.D., Liu, C.J., Soofi, A., Fogg, V.C., Hershenson, M.B., Dressier, G.R., Deutsch, G.H., et al. (2014). Crumbs3 is essential for proper epithelial development and viability. Molecular and cellular biology 34, 43-56.

• Wu, D., Hong, H., Huang, X., Huang, L., He, Z., Fang, Q., and Luo, Y. (2016). CXCR2 is decreased in preeclamptic placentas and promotes human trophoblast invasion through the Akt signaling pathway. Placenta 43, 17-25.

• Wu, L., Wu, Y., Peng, B., Hou, Z., Dong, Y., Chen, K., Guo, M., Li, H., Chen, X., Kou, X., et al. (2017). Oocyte-Specific Homeobox 1, Oboxl, Facilitates Reprogramming by Promoting Mesenchymal-to-Epithelial Transition and Mitigating Cell Hyperproliferation. Stem Cell Reports 9, 1692-1705.

• Wu, X., Oatley, J.M., Oatley, M.J., Kaucher, A.V., Avarbock, M.R., and Brinster, R.L. (2010). The POU domain transcription factor POU3F1 is an important intrinsic regulator of GDNF-induced survival and self-renewal of mouse spermatogonial stem cells. Biology of reproduction 82, 1103-1 111.

• Yamamizu, K., Sharov, A.A., Piao, Y., Amano, M., Yu, H , Nishiyama, A., Dudekula, D.B., Schlessinger, D., and Ko, M.S. (2016). Generation and gene expression profiling of 48 transcription-factor-inducible mouse embryonic stem cell lines. Scientific reports 6, 25667.

•      Ying, Q.-L., Wray, J., Nichols, J., Batlle-Morera, L., Doble, B., Woodgett, J., Cohen, P., and Smith, A. (2008). The ground state of embryonic stem cell self-renewal. Nature 453, 519.

•      Yu, J., Vodyanik, M.A., Smuga-Otto, K., Antosiewicz-Bourget, J., Frane, J.L., Tian, S., Nie, J., Jonsdottir, G.A., Ruotti, V., Stewart, R., et al. (2007). Induced pluripotent stem cell lines derived from human somatic cells. Science 318, 1917-1920.

•      Yun, C, Mendelson, J., Blake, T., Mishra, L., and Mishra, B. (2008). TGF-beta signaling in neuronal stem cells. Disease markers 24, 251-255.

•      Zhao, T., Fu, Y., Zhu, J., Liu, Y., Zhang, Q., Yi, Z., Chen, S., Jiao, Z., Xu, X., Xu, J., Duo, S., Bai, Y., Tang, C, Li, C, and Deng, H. (2018). Single-Cell RNA-Seq Reveals Dynamic Early Embryonic-like Programs during Chemical Reprogramming. Cell Stem Cell 23, 1-15.

•      Zunder, E.R., Lujan, E., Goltsev, Y., Wernig, M., and Nolan, G.P. (2015). A continuous molecular roadmap to iPSC reprogramming through progression analysis of single-cell mass cytometry. Cell Stem Cell 16, 323-337.

•      Zwiessele, M., and Lawrence, N.D. (2016). Topslam: Waddington Landscape Recovery for Single Cell Experiments. bioRxiv.

**[0318]**    **Key resources**

**[0319]**    Key resources used in this study are shown below.

| REAGENTS or RESOURCE | SOURCE | IDENTIFIER |
|---|---|---|
| Recombinant DNA | | |
| FUW Tet-On vector | Addgene | #20323 |
| *Zfp42* cDNA | Origene | MG203929 |
| *Obox6* cDNA | Origene | MR215428 |
| Chemicals, Peptides, and Recombinant Proteins | | |
| leukemia inhibitory factor (LIF) | Millipore | ESG1107 |
| PD0325901 | Sigma | PZ0162-25MG |
| CHIR99021 | Sigma | PZ0162-25MG |
| Critical Commercial Kits | | |
| Chromium™ Single Cell 3' Reagent | 10X genomics | PN-120230, PN-120231, PN- |

| | | |
|---|---|---|
| Kits v l | | 120232 |
| Chromium™ Single Cell 3' Reagent Kits v2 | 10X genomics | PN- 120237 |
| Fugene HD reagent | Promega | E2311 |
| Cloning Reagents | | |
| Gibson Assembly | NEB | E2611S |
| Sequence-Based Reagents | | |
| Deposited Data | | |
| Single cell RNA-seq raw data (pilot study) | NCBI Gene Expression Omnibus | GSE106340 |
| Single cell RNA-seq raw data | NCBl Gene Expression Omnibus | GSE115943 |
| Experimental Models: Organisms/Strains | | |
| OKSM secondary MEFs | Konrad Hoched linger lab | OKSM x B6.Cg-$Gt(ROSA)26Sor^{Jml(rtTA*M2)Jae}$/J x B6;$l29S4$-$Pou5f l^{tm2Jae}$/J |
| Primary MEFs | Rudolf Jaenisch lab | B6.Cg-$Gt(ROSA)26Sof^{M1<rtTA*M2)Jae}$lJ x B6;129S4-/$^{D}ou5f7^{e\ 2Ja}$7J |
| Software and Algorithms | | |
| Waddington-OT | This paper | https ://github .com/broadinstitute /wot |
| Scaling algontlim for unbalanced transport | (Chizat et al., 2016) | |
| CeilRanger | 10X genomics | v2.0.0 |
| ForceAtlas2 | Gephi | vO.9.2 |
| Seurat | | v2.1.0 |
| Scanpy | | v0.2.8 |
| Monocle2 | (Qiu et al. 2017) | v2.8.0 |
| URD | (Farrell et al 2018) | vl.O |

**[0320]**    **Method Details**

**[0321]**    **I. Modeling developmental processes with optimal transport**

**[0322]**    We developed a method to model development based on Optimal Transport. Section 1 reviews the concept of gene expression space and introduces our probabilistic framework for time series of expression profiles. Section 2 introduces our key modeling assumption to infer temporal couplings over short time scales. Section 3 shows how we can compute an optimal coupling between adjacent time points by solving a convex optimization problem, and how we can leverage an assumption of Markovity to compose adjacent time points and estimate temporal couplings over longer intervals. Section 4 describes how to interpret transport maps. Specifically, Section 4.1 shows how to compute ancestors and descendants of cells, Section 4.2 describes an interesting physical interpretation of entropy-regularization, and Section 4.3 shows how we learn gene regulatory networks to summarize the trajectories.

**[0323]**      1. Developmental processes in gene expression space

**[0324]**    A collection of mRNA levels for a single cell is called an *expression profile* and is often represented mathematically by a vector in *gene expression space.* This is a vector space that has dimension equal to the number of genes, with the value of the $i$th coordinate of an expression profile vector representing the number of copies of mRNA for the $i$th gene. Note that real cells only occupy an integer lattice in gene expression space (because the number of copies of mRNA is an integer), but we pretended that cells can move continuously through a real-valued $G$ dimensional vector space.

**[0325]**    As an individual cell changes the genes it expresses over time, it moves in gene expression space and describes a trajectory. As a population of cells develops and grows, a *distribution* on gene expression space evolves over time. When a single cell from such a population is measured with single cell RNA sequencing, we obtained a noisy estimate of the number of molecules of mRNA for each gene. We represented the measured expression profile of this single cell as a sample from a probability distribution on gene expression space. This sampling captured both (a) the randomness in the single-cell RNA sequencing measurement process (due to subsampling reads, technical issues, etc.) and (b) the random selection of a cell from the population. We treated this probability distribution as *nonparametric* in the sense that it wsa not specified by any finite list of parameters.

**[0326]** In the remainder of this section we introduced a precise mathematical notion for a *developmental process* as a generalization of a stochastic process. Our primary goal was to infer the ancestors and descendants of subpopulations evolving according to an unknown developmental process. This information was encoded in the *temporal coupling* of the process, which is lost because we kill the cells when we perform scRNA-Seq. We claimed it was possible to recover the temporal coupling over short time scales provided that cells don't change too much. Therefore we could make inferences about which cells go where. We showed in the remainder of this section how to do this with *optimal transport.*

**[0327]** 1.1 A mathematical model of developmental processes

**[0328]** We began by formally defining a precise notion of the developmental trajectory of an individual cell and its descendants. Intuitively, it was a continuous path in gene expression space that bifurcated with every cell division. Formally, we defined it as follows:

**[0329]** Definition 1 (single-cell developmental trajectory). Consider a cell $x(0) \in \mathbb{R}^G$. Let k(t) ≥ 0 specifiy the number of descendants at time t, where k(0) = 1. A single-cell development trajectory is a continuous function

$$x : [0, T) \to \underbrace{\mathbb{R}^G \times \mathbb{R}^G \times \ldots \times \mathbb{R}^G}_{k(t) \text{ times}}.$$

*This means that x(t) is a k(t)-tuple of cells, each represented by a vector in $\mathbb{R}^G$ :*

$$x(i) = (x_1(t), \ldots, x_{k(t)}(t)) .$$

*We referred to the cells x\(t), . . . , xk(t)(t) as the descendants of x(0).*

**[0330]** Note that we could not directly measure the temporal dynamics of an individual cell because scRNA-Seq was a destructive measurement process: scRNA-Seq lysed cells so it was possible to measure the expression profile of a cell at a single point in time. As a result, it was not possible to directly measure the descendants of that cell, and the full trajectory was unobservable. However, one can learn something about the probable trajectories of individual cells by measuring snapshots from an evolving population.

**[0331]** Published methods typically represent the aggregate trajectory of a population of cells by means of a graph structure. While this recapitulates the branching path traveled by the descendants of an individual cell, it may over-simplify the stochastic nature of developmental processes. Individual cells have the potential to travel through different paths, but any given cell

travels one and only one such path. Our goal was to assign a likelihood to the set of possible paths, which in general were not finite and therefore cannot be a represented by a graph.

**[0332]** We defined a developmental process to be a time-varying probability distribution on gene expression space. One simple example of a distribution of cells is that we can represent a set of cells

$x_1, \ldots, x_n$ by the distribution

$$\mathbb{P} = \frac{1}{n} \sum_{i=1}^{n} \delta_{x_i}.$$

**[0333]** Similarly, we could represent a set of single-cell trajectories $x_1(t), \ldots, x_n(t)$ with a distribution over trajectories. This was a special case of a developmental process, which we defined as follows:

**Definition 2** (developmental process). A developmental process Pt is a time-varying distribution (i.e. stochastic process) on gene expression space.

**[0334]** Recall that a stochastic process was determined by its temporal dependence structure. This was specified by the coupling (i.e. joint distribution) between random variables at different time points. Given that a cell had a particular expression profile $y$ at time $t_2$, where did it come from at time $t_1$? This was the information lost by not tracking individual cells overtime.

**[0335]** **Definition 3** (temporal coupling). Let Pt be a developmental process and consider two time points s<t. Let Xt ~ Pt denote the expression profile of a random cell at time t and let $X_s$ denote the expression profile of the cell of origin at times.

**[0336]** The temporal coupling $\gamma_{s,t}$ is defined as the law of the joint distribution:

$$\gamma_{s,t} = \mathcal{L}(X_s, X_t).$$

*Equivalently,*

$$\int_{x \in A} \int_{y \in B} \gamma_{s,t}(x,y)\,dx\,dy = \Pr\{X_s \in A, X_t \in B\}$$

*for any seis $A, B \subset S^G$-*

**[0337]** The temporal coupling $\gamma_s t$ was not technically a coupling of $P_s$ and P? in the standard sense because it does not necessarily have marginals $P_s$ and P?:

$$\int \gamma_{s,t}(x,y)\,dx = \mathbb{P}_t(y), \quad \text{but} \quad \int \gamma_{s,t}(x,y)\,dy \neq \mathbb{P}_s(x).$$

**[0338]**    Biologically, this was the case when cells grow at different rates. Then proliferative cells from the earlier time point were over-represented when we look for the origin of cells at the later time point. In the following definition, we introduced a relative growth rate function to describe the relationship between the expression profile of a cell and the average number of living descendants it gave rise to after certain amount of time.

**[0339]**    **Definition 4.** A relative growth rate function associated with a temporal coupling is a function g(x)

*satisfying*

$$\int \gamma_{s,t}(x,y)dy = \mathbb{P}_s(x)\frac{g(x)^{t-s}}{\int g(x)^{t-s}d\mathbb{P}_s(x)}.$$

**[0340]**    The integral on the left-hand side represented the amount of mass coming out of *x* and going to any *y*. The term P(x) on the right hand side accounted for the abundance of cells with expression profile x, and the function *g(x)* represented the exponential increase in mass per unit time.

**[0341]**    Having defined the notion of developmental processes and temporal couplings, we now turned to estimating these from data.

**[0342]**    2. The optimal transport principle for developmental processes

**[0343]**    Single-cell RNA-Seq allowed us to sample cells from a developmental process at various time points, but it did not give any information about the coupling between successive time points. Without making any assumptions, it was impossible to recover the temporal coupling even given infinite data in the form of the full distributions $P_s$ and Pi. However, we claimed that it was reasonable to assume that cells don't change expression by large amounts over short time scales. This assumption allowed us to estimate the coupling and infer which cells go where.

**[0344]**    We began with a simple one-dimensional example to build intuition.

**[0345]**    **Example 1.** Let Xo ~ $N(0, \sigma^2)$ and Xi ~ $N(\mu, \sigma^2)$ be one dimensional Gaussian variables representing the location of a particle at time 0 and at time 1. One simple heuristic to estimate $\hat{\gamma}$ is to minimize the squared distance that the particle moves from time 0 to time 1:

$$\hat{\gamma} \leftarrow \arg\min_{\pi} \mathbb{E}_{\pi}\|X_0 - X_1\|^2.$$

**[0346]**      We minimized over all couplings $\pi$ with marginals $(0, \sigma^2)$ and $(\mu, \sigma^2)$. One can check that the optimal joint distribution is a two dimensional Gaussian with the following dependence structure:

$X_1 = X_0 + \mu$.

**[0347]**      This heuristic to couple marginals was called *optimal transport* (OT). If $c(x, y)$ denoted the cost of transporting a unit mass from $x$ to $y$, and the amount we transferred from $x$ to $y$ is $\pi(\chi, y)$, then the total cost of transporting mass according to such a transport plan $\pi$ is given by

$$\iint c(x, y)\pi(x, y)dxdy.$$

**[0348]**      In this study we focused on the cost defined by the squared-Euclidean distance

$$c(x, y) = \|x - y\|^2,$$

**[0349]**      on an appropriate input space. We made this choice to focus on Wasserstein-2 transport because of the many attractive theoretical properties it enjoyed over Wasserstein-1 transport (Villani, 2008).

**[0350]**      The optimal transport plan minimized the expected cost subject to marginal constraints:

$$\pi(\mathbb{P}, \mathbb{Q}) = \underset{\pi}{\text{minimize}} \quad \iint c(x, y)\pi(x, y)dxdy$$
$$\text{subject to} \quad \int \pi(x, \cdot)dx = \mathbb{Q} \tag{1}$$
$$\int \pi(\cdot, y)dy = \mathbb{P}.$$

**[0351]**      Note that this was a linear program in the variable $\pi$ because the objective and constraints were both linear in $\pi$. The optimal objective value defined the *transport distance* between P and Q (it was also called the Earthmover's distance or Wasserstein distance). Unlike many other ways to compare distributions (such as KL-divergence or total variation), optimal transport took the geometry of the underlying space into account. For example, the KL-Divergence was infinite for any two distributions with disjoint support, but the transport distance depended on the separation of the support. For a comprehensive treatment of the rich mathematical theory of optimal transport, we refer the reader to (Villani, 2008).

**[0352]**     2.1 The optimal transport principle for developmental processes.

**[0353]**     We proposed to use optimal transport to estimate the temporal coupling of a developmental process. We made two modifications to classical optimal transport to adapt it to our biological setting.

**[0354]**     1. Classical optimal transport had conservation of mass built into the constraints (1). We accounted for growth by rescaling the distribution **Pi** before applying OT.

**[0355]**     2. The coupling identified by classical optimal transport was purely deterministic in the sense that each point was transported to a single point. However, for cells whose fates were not completely determined, the true coupling should have a degree of entropy to it. We therefore added a term to the objective to promote entropy in the transport coupling.

**[0356]**     Injecting a small amount of entropy also made sense even for a population of cells with truly deterministic descendant distribution. When we sampled finitely many cells at time $t_2$, the true descendants of any given $t\backslash$ cell were not captured. Therefore entropy in the transport map could be used to represent our statistical uncertainty in the inferred descendant distribution.

**[0357]**     In order to state the optimal transport principle, we first introduced some notation. Let **Pi** denote a developmental process with temporal coupling $y_st$ and with relative growth function $g(x)$. Let $Q_s$ denote the distribution obtained by rescaling $\mathbf{P}_s$ by the relative growth rate:

$$\mathbb{Q}_s(x) = \mathbb{P}_s(x) \frac{g^{t-s}(x)}{\int g^{t-s}(z) d\mathbb{P}_s(z)}.$$

**[0358]**     Finally, let $\pi_{s,t}(\epsilon)$ denote the entropy-regularized optimal transport coupling of $Q_s$ and **Pi,** defined as the solution to the following optimization problem

$$\pi_{s,t}(\epsilon) = \underset{\pi}{\text{minimize}} \quad \int\int c(x,y)\pi(\chi,y)dxdy - \epsilon \int \pi(x,y)\log \pi(\mathrm{x},\mathbf{\textit{y}})\mathbf{\textit{dxdy}}$$

$$\text{subject to} \quad \int \pi(x,\cdot)dx = \mathbb{Q}_s \qquad\qquad\qquad (2)$$

$$\int \pi(\cdot,y)dy = \mathbb{P}_t.$$

**[0359]**     We now stated the optimal transport principle for developmental process

$$s \approx t \implies \pi_{s,t}(\epsilon) \approx \gamma_{s,t}.$$

**[0360]** In words, over short time scales, the true coupling was well approximated by the OT coupling. In section 3, we show how to estimate $\pi_{s_i*}(\epsilon)$ **from** data (we occasionally omit the dependence on $\epsilon$ and write $\pi_{s,t}$). This in turn gives us an estimate of $y_{s,t.}$

**[0361]** 3. Inferring temporal couplings from empirical data

**[0362]** In this section we showed how to estimate the temporal couplings of a developmental process from data.

**[0363]** **Definition 5** (developmental time series). A developmental time series was a sequence of samples from a developmental process Pt on R^. This was a sequence of sets Si, . . . , $S_T \subset R^G$ collected at times ti, . . . ,ïτ GR. Each Si is a set of expression profiles in $R^G$ drawn independently from $P_t$.

**[0364]** From this input data, we formed an empirical version of the developmental process. Specifically, at each time point $t$, we formed the empirical probability distribution supported on the data* $S$. We summarize this in the following definition:

**[0365]** **Definition 6** (Empirical developmental process). An empirical developmental process $\hat{P}t$ is a time vary- ing distribution constructed from a developmental time course Si, . . . , $S_T$:

$$\hat{P}_{t_i} = \frac{1}{|S_i|} \sum_{x \sim v_i} \delta_x. \qquad (3)$$

**[0366]** The empirical developmental process was undefined for t $\notin$ {ti, . . . ,ïτ}.

**[0367]** In order to estimate the coupling from time $t_1$ to time $t_2$, we first constructed an initial estimate the growth rate function g(x). In practice, we form an initial estimate ĝ(x) as the expectation of a birth-death process on gene expression space with birth-rate β(χ) and death rate δ(χ) defined in terms of expression levels of genes involved in cell proliferation and apoptosis. We ultimately leveraged techniques from unbalanced transport (Chizat et al., 2017) to refine this initial estimate to learn cellular growth and death rates automatically from data.

**[0368]** We then form the rescaled empirical distribution

$$\hat{Q}_{t_1}(x) = \hat{P}_{t_1}(x) \frac{\hat{g}(x)^{t_1 - t_2}}{\int \hat{g}(z)^{t_1 - t_2} d\hat{P}_{t_1}(z)},$$

and compute the optimal transport map $\hat{\pi}_{t_1,t_2}$ between $\hat{Q}t_1$ and $\hat{P}t_2$

**[0369]** 3.1 Estimating couplings between adjacent time points

**[0370]**    In order to identify an optimal transport plan connecting Qˆt1 and Pˆt2 , we solved an optimization problem with a matrix-valued optimization variable. In the classical zero-entropy setting (2) with « = 0 was a linear program. While the classical optimal transport linear program could be difficult to solve for large numbers of points, fast algorithms have been recently developed (Cuturi, 2013) to solve the entropically regularized version of the transport program. Entropic regularization speeded up the computations because it made the optimization problem strongly convex, and gradient ascent on the dual could be realized by successive diagonal matrix scalings called Sinkhorn iterations (Cuturi, 2013). These were very fast operations.

**[0371]**    The scaling algorithm for entropically regularized transport had also been extended to work in the setting of unbalanced transport (Chizat et al., 2017), where the equality constraints were relaxed to bounds on the marginals of the transport plan (in terms of KL-divergence or total variation or a general f-divergence). In our application this was very attractive from a modeling perspective for the following reasons:

**[0372]**    1. We may have specified the growth rate function gˆ(x). Unbalanced transport adjusted the input growth rate in order to reduce the transport cost. This allowed us to automatically learn growth rates from scratch.

**[0373]**    2. Even if the growth rates were completely uniform, the random sampling could introduce what looked like growth. For example, suppose there was a rare subpopulation of cells consisting of 5% of the total. If at one time point, we randomly sampled fewer of these cells so that they comprised 4% of the total, and at the next time point we sample 6%, then it would look like this population had increased by 50%. Unbalanced transport could automatically adjust for this apparent growth.

**[0374]**    We used both entropic regularization and unbalanced transport. To compute the transport map between the empirical distributions of expression profiles observed at time $t$, and $t_{i+i}$, we solved the following optimization problem

$$\hat{\pi}_{t_i,t_{i+1}} = \underset{\pi}{\arg\min} \sum_{x \in S_i} \sum_{y \in S_{i+1}} c(x, y)\pi(x, y) - \epsilon \int \nu(x, y)\log \pi(x, y)dxdy$$

$$\text{subject to} \quad KL\left[\sum_{x \in S_i}\pi(x, y)\middle\| d\hat{\mathbb{P}}_{t_{i+1}}(y)\right] \le \frac{1}{\lambda_1} \qquad (4)$$

$$KL\left[\sum_{y \in S_{i+1}}\pi(x, y)\middle\| d\hat{\mathbb{Q}}_{t_i}(x)\right] \le \frac{1}{\lambda_2}$$

**[0375]**    where $\epsilon$, $\lambda_1$ and $\lambda_2$ are regularization parameters.

**[0376]**    This is a convex optimization problem in the matrix variable $\pi \in \mathbb{R}^{N_i \times N_{i+1}}$, where $N_i = |s_i|$ is the number of cells sequenced at time ti. It takes about 5 seconds to solve this unbalanced transport problem using the scaling algorithm of (Chizat et al., 2017) on a standard laptop with $N_i \approx 5000$.

**[0377]**    Note that by default the densities (on the discrete set Si) of the empirical distributions specified in equation (3) are simply $d\hat{\mathbb{P}}_{t_i}(x) = \frac{1}{N_i}$. However, in principle one  could use nonuniform  empirical distributions (e.g., if one wanted to include information about cell quality).

**[0378]**    To summarize: given a sequence of expression profiles $S_i, \ldots, S_T$, we solved the optimization  problem (4) for each successive  pair of time points $S_i, S_{i+1}$. For the pair of time-points $(t_i, t_{i+1})$, this gave us a transport map $\hat{\pi}_{t_i,t_{i+1}}$. With  enough  data, this  may  be  a  good estimate  of $\pi_{t_i,t_{i+1}}$ -because  it is well  known  that transport  maps  are consistent  in the sense that

$$\lim_{N_i,N_{i+1} \to \infty} \hat{\pi}_{t_i,t_{i+1}} = \pi_{t_i,t_{i+1}}.$$

**[0379]**    Taken together with the optimal transport principle:

$$\pi_{t_i,t_{i+1}} \approx \gamma_{t_i,t_{i+1}},$$

**[0380]**    We therefore  could  estimate $\gamma_{t_i,t_{i+1}}$ from $\hat{\pi}_{t_i,t_{i+1}}$ when Ni is large enough.

**[0381]**    3.2 Estimating long-range couplings

**[0382]**    We relied on an assumption of Markovity (or memorylessness) in order to estimate couplings over longer time intervals. Recall that a stochastic process was Markov if the future was independent of the past, given the present. Equivalently, it was fully specified by the

couplings between pairs of time points. We defined Markov developmental processes in a similar spirit:

**[0383]**     **Definition 7** (Markov developmental process). A Markov developmental process $P_t$ is a time-varying distribution on $R^{\wedge}$ that is completely specified by couplings between pairs of time points in the following sense. For any three time points $s < t < \tau$, the long-range coupling $\gamma_{s,\tau}$ was equal to the composition of short-range couplings:

$$\gamma_{t,\tau} \circ \gamma_{s,t} = \gamma_{s,\tau}.$$

**[0384]**     Note that the optimal transport maps $\hat{\pi}_{s,t}$ did not have this compositional property. Composing the OT coupling from time $s$ to $t$ and then from $t$ to $\tau$ was not the same as optimally transporting from $s$ directly to $\tau$. In general, we do not recommend computing OT maps directly between non-adjacent time points. We leveraged the Markovity assumption to estimate couplings over long time intervals by composing estimates over shorter intervals. Formally, for any pair of time points $t_i$, $t_{i+k}$, we estimate the coupling $\hat{\gamma}_{t_i,t_{i+k}}$ by composing as follows:

$$\hat{\gamma}_{t_i,t_{i+k}} = \hat{\pi}_{t_i,t_{i+1}} \circ \hat{\pi}_{t_{i+1},t_{i+2}} \circ \ldots \circ \hat{\pi}_{t_{i+k-1},t_{i+k}}.$$

**[0385]**     These compositions were computed via ordinary matrix multiplication.

**[0386]**     It is an interesting question to what extent developmental processes are Markov. On gene expression space, they were likely not strictly Markov because, for example, the history of gene expression could influence chromatin modifications, which may not themselves be fully reflected in the observed expression profile but could still influence the subsequent evolution of the process. However, it was possible that developmental processes could be considered Markov on some augmented space. Note that our core technique for estimating a single temporal coupling over a short time interval does **not** rely on any Markov assumption.

**[0387]**     4. Interpreting transport maps

**[0388]**     In the previous section we introduced the principle of optimal transport for time series of gene expression profiles. Given a time series of expression profiles $S_1, \ldots, S_T$, we used this principle to compute a sequence of transport maps between subsequent time slices. In this section we define the *ancestors* and *descendants* of any subset of cells from this sequence of transport maps in section 4.1. Then, in section 4.2 we explain an intuitive physical interpretation of

entropy-regularization. Finally, in section 4.3 we describe a connection between optimal transport, gradient flows, and Waddington's landscape.

**[0389]**     4.1 Defining ancestors, descendants and trajectories

**[0390]**     We defined the descendants and ancestors of subgroups of cells evolving according to a Markov (i.e. memoryless) developmental process.

**[0391]**     Our definition of ancestors and descendants relies on a notion of *pushing* sets of cells through a trans- port map. Before defining ancestors and descendants, we introduce this terminology. As a distribution on the product space $\mathrm{R}^G \times \mathrm{R}^G$, a coupling *y* assigns a number *y(A, B)* to any pair of sets *A, B* $\subset R^G$

$$\gamma(A, B) = \int_{x \in A} \int_{y \in B} \gamma(x, y) dx dy.$$

**[0392]**     This number $\pi(A, B)$ represented the amount of mass coming from *A* and going to *B*. When we did not specify a particular destination, the quantity *y{A, )* specified the full distribution of mass coming from *A*. We referred to this action as *pushing A* through the transport plan *y*. More generally, we could also push a *distribution* $\mu$ forward through the transport plan *y* via integration

$$\mu \mapsto \int \gamma(x, \cdot) d\mu(x).$$

**[0393]**     We refer to the reverse operation as pulling a set *B* back through *y*. The resulting distribution $\gamma(\cdot, B)$ encodes the mass ending up at *B*. We can also pull distributions $\mu$ back through *y* in a similar way:

$$\mu \mapsto \int \gamma(\cdot, y) d\mu(y).$$

**[0394]**     We sometimes refer to this as *back-propagating* the distribution $\mu$ (and to pushing $\mu$ forward as *forward propagation)*.

**[0395]**     Equipped with this terminology, we define ancestors and descendants as follows:

**[0396]**     **Definition 8** (descendants in a Markov developmental process). Consider a set of cells $C \subset \mathbb{R}^G$ which lived at time ti were part of a population of cells evolving according to a Markov developmental process Pt. Let $\gamma_{t_1, t_2}$ denote the coupling from time ti to time $t_2$. The descendants of C at time $t_2$ are obtained by pushing C through $\gamma$.

**[0397]**   **Definition 9** (ancestors in a Markov developmental process). Consider a set of cells $C \subset \mathbb{R}^G$, which lived at time $t_2$ and were part of a population of cells evolving according to a Markov developmental process $P_t$. Let $\pi$ denote the transport map for $P_t$ from time $t_2$ to time ti. The ancestors of C at time ti were obtained by pulling C back through $\gamma$.

**[0398]**   **Trajectories:** We defined to the *ancestor trajectory* to a set $C$ as the sequence of ancestor distributions at earlier time points. Similarly, we refer to the *descendant trajectory* from a set $C$ as the sequence of descendant distributions at later time points.

**[0399]**   4.2 A physical interpretation of entropy regularized optimal transport

**[0400]**   In this section we explain an interesting physical interpretation of entropy-regularized optimal transport. Consider a collection of $N$ indistinguishable particles undergoing Brownian motion with diffusion coefficient $\epsilon$. Suppose we observe the $N$ particle positions at time 0 and at time 1. If $N=l,$ the distribution on paths connecting the starting and ending point is called a Brownian bridge. For $N > 1$, the distribution over paths involves two components:

**[0401]**   1. A coupling of the particles specifying which particle goes where (because the particles are indistinguishable, this is not uniquely specified by the observations).

**[0402]**   2. Given a matching, the distribution on paths for each matched pair is a Brownian bridge.

**[0403]**   The coupling was a random permutation that matched points at time 0 to points at time 1. The distribution of this random permutation depends on the variance of the Brownian motion. It turned out that the expected (i.e. average) coupling could be computed by maximum entropy optimal transport. These ideas could be traced back to Schrodinger's 1932 work in statistical electrodynamics (Schrodinger, 1932), but the connection to optimal transport was not made explicit until recently (Le´onard, 2014). We summarize this in the following theorem:

**[0404]**   **Theorem 1.** Entropy regularized optimal transport gives the expectation of the distribution over cou- plings induced by Brownian motion (when the diffusion coefficient of the Brownian motion is equal to the entropy regularization parameter).

**[0405]**   4.3 Gradient flow and Waddington's landscape

**[0406]**   In this section we show how optimal transport can be interpreted as a gradient flow in gene expression space (capturing cell-autonomous processes) or in the space of distributions

(capturing cell-nonautonomous processes). For a full treatment of the rich OT theory of gradient flows, we refer the reader to (Ambrosio et al., 2005; Santambrogio, 2015).

**[0407]** We began by considering the simple setting described by Waddington's landscape, which described a gradient flow in gene expression space and is a special case of what we could capture with optimal transport. Mathematically, Waddington's landscape defined a potential function $\Phi$ assigning potential energy $\Phi(\chi)$ to a cell with expression profile $x$. The cells roll eddownhill according to the gradient of $\Phi$ to describe a trajectory $x(t)$ satisfying the differential equation

$$\frac{dx}{dt} = -\nabla\Phi(x). \tag{5}$$

**[0408]** This equation governing the trajectory of individual cells induced a flow in the distribution of the population of cells:

$$\frac{d\mathbb{P}_t}{dt} = \mathrm{div}\,[\nabla\Phi(x)F_t]. \tag{6}$$

**[0409]** Intuitively, this equation stated that the change in mass for each small volume of space (on the left-hand side) was equal to the flux of mass in and out (given by the divergence on the right hand side).

**[0410]** Optimal transport can capture this type of potential driven dynamics: the true coupling specified by (5) is close to the optimal transport coupling over short time scales. To motivate this, we appeal to a classical theorem establishing a dynamical formulation of optimal transport.

**[0411]** **Theorem 2** (Benamou and Brenier, 2001). The optimal objective value of the transport problem (1) is equal to the optimal objective value of the following optimization problem

$$\begin{aligned}
\underset{\rho,v}{\text{minimize}} \quad & \int_0^1\int_{\mathbb{R}^d} \|v(t,x)\|^2 \rho(t,x)\,dt\,dx \\
\text{subject to} \quad & \rho(0,\cdot) = \mathbb{P}, \quad \rho(1,\cdot) = \mathbb{Q} \\
& \nabla\cdot(\rho v) = \frac{\partial\rho}{\partial t}
\end{aligned} \tag{7}$$

**[0412]** In this theorem, v was a vector-valued velocity field that advected the distribution $p$ from P to Q, and the objective value to be minimized was the kinetic energy of the flow (mass x squared velocity). In our setting, the two distributions were snapshots $P_s$ and $P_t$ of a developmental process at two time points, and the theorem showed that the transport map $\pi_{s,t}$

could be seen as a point-to-point summary of a least-action continuous time flow, according to an unknown velocity field. In the special case when the velocity field was the gradient of a potential Φ (i.e. Waddington landscape), the theorem implied that the coupling (5) achieved the optimal transport cost. In other words, OT could capture potential driven dynamics. In addition, optimal transport could also describe much more general settings. This velocity field could change over time and also depended on the entire distribution of cells, so optimal transport could describe very general developmental processes including those with cell-cell interactions, as described below.

[0413]    We showed that the evolution (6) was a special case of a *Wasserstein gradient flow* to minimize  the linear energy functional

$$E(\mathbb{P}) = \int \Phi(x) d\mathbb{P}(x).$$

[0414]    We then described non-linear gradient flows, which can capture cell-cell interactions. To understand gradient flows, we started with the familiar notion of gradient descent:

$$x_{k+1} = -\eta \nabla E(x_k) + x_k.$$

[0415]    This was rewritten as *a proximal procedure,* where one seeks to minimize $E$ over all $x$ in the proximity of ¾

$$\text{¾}_{+1} = \quad \arg\min_{x} \quad E(x) + \frac{1}{2\eta} \|s - x_k\|^2. \tag{8}$$

[0416]    We performed a similar proximal procedure in the space of distributions, replacing the Euclidean norm $\|\cdot\|^2$ with the Wasseerstein distance:

$$\mathbb{P}_{k+1} = \quad \arg\min_{p} E(\rho) + \frac{1}{2\eta} W_2^2(\rho, \mathbb{P}_k). \tag{9}$$

[0417]    This produced a sequence of iterates $\mathbb{P}_0$, $\mathbb{P}_i$, . . . , $\mathbb{P}_k$. The gradient flow was the limit obtained as we shrink the step-size $\eta \downarrow 0$. In (Richard Jordan and Otto, 1998), it's proven that for the linear energy functional

$$E(\mathbb{P}) = \int \Phi(x) d\mathbb{P}(x),$$

[0418]    the limiting gradient flow converges to a solution of (6).

[0419]    Going beyond the linear energy functional associated with Waddington's landscape, one could describe cell-cell interactions with an interaction energy of the form

$$E(\mathbb{P}) = \iint I(x,y)d\mathbb{P}(x)d\mathbb{P}(y).$$

**[0420]** Gradient flows for interaction potentials are discussed in chapter 7 of (Santambrogio, 2015).

**[0421]** **Learning models of gene regulation** Motivated by this interpretation of optimal transport as a gradient flow according to an unknown vector field, we described a strategy to estimate such a vector field from data in **Waddington-OT: Concepts and Implementation .** We interpreted the vector field as a model of gene regulation - it predicted gene expression at later time points as a function of transcription factor expression at current time points. We assumed that the vector field did not change over time, and described a cell-autonomous flow, but we do not assume that it comes from a potential function.

**[0422]** **II. WADDINGTON-OT : Concepts and Implementation**

**[0423]** Building on the theoretical foundations developed in Modeling developmental processes with optimal transport, we developed WADDINGTON-OT: our method for computing ancestor and descendant trajectories, interpolating developmental processes, inferring gene regulatory models, and visualizing developmental landscapes. We begin with an overview in Section 1, and we then describe the specific details in Sections 2 - 8.

**[0424]** 1. Overview

**[0425]** To apply WADDINGTON-OT to a new dataset. The code is available on GitHub: https://github .com/broadinstitute/wot/

**[0426]** In the sections below we describe our procedures for computing transport maps, computing trajectories to cell sets, fitting local and global regulatory models, visualizing the developmental landscape, interpolating the distribution of cells at held-out time points.

**[0427]** To keep the focus here general-purpose, we deferred all reprogramming-specific details to the subsequent sections Methods.

**[0428]** **Input data:** The input to our suite of methods was a temporal sequence of single cell gene expression matrices, prepared as described in **Preparation of expression matrices .**

**[0429]** **Computing** transport **maps:** Waddington-OT calculated transport maps between consecutive time points and automatically estimated cellular growth and death rates. In Section 2

below we provide guidelines for defining the cost function, selecting regularization parameters and (optionally) providing an initial estimate of growth and death rates.

[0430]    **Ancestors, descendants, and trajectories:** We describe in Section 3 how we computed trajectories plot trends in gene expression. Briefly, the *developmental trajectory* of a subpopulation of cells refers to the sequence of ancestors coming before it and descendants coming after it. Using the transport maps, we calculated the forward or backward transport probabilities between any two classes of cells at any time points. For example, we took successfully reprogrammed cells at day 18 and use back-propagation to infer the distribution over their precursors at day 17.5. We then propagated this back to day 17, and so on to obtain the ancestor distributions at all previous time points. This was the developmental trajectory to iPS cells. We plotted trends in gene expression over time.

[0431]    **Fitting regulatory models:** We describe our method to fit a regulatory model to the transport maps in Section 4. Transcription factors (TFs) that appeared to play important roles along trajectories to key destinations were identified by two approaches. The first approach involved constructing a global regulatory model. Pairs of cells at consecutive time points were sampled according to their transport probabilities; expression levels of TFs in the cell at time t were used to predict expression levels of all non-TFs in the paired cell at time t + 1, under the assumption that the regulatory rules are constant across cells and time points. (TFs were excluded from the predicted set to avoid cases of spurious self-regulation). The second approach involved local enrichment analysis. TFs were identified based on enrichment in cells at an earlier time point with a high probability ($> 80\%$) of transitioning to a given fate vs. those with a low probability ($< 20\%$).

[0432]    Visualizing **the developmental landscape** To visualize the developmental landscape, we first reduced the dimensionality of the data with diffusion components, and then embedded the data in two dimensions with force-directed graph visualization (as described in Section 5). While alternative visualization methods, such as t-distributed Stochastic Neighbor Embedding (t-SNE), were well suited for identifying clusters, they did not preserve global structures relevant to studying trajectories across a time course. FLE better reflected global structures by including repulsive forces between dissimilar points. In particular, these repulsive forces seemed to do a good job of splaying out the spikes present in the diffusionmap embedding.

**[0433]**   **Geodesic interpolation:** To validate the temporal couplings, Waddington-OT could interpolate the distribution of cells at a held-out time point. The method wsa performing well if the interpolated distribution was close to the true held-out distribution (compared to the distance between different batches of the held-out distribution). Otherwise, it was possible that the method requires more data or finer temporal resolution.

**[0434]**   Section 6 describes our method to interpolate the distribution of cells at a held-out time point. Our validation results for IPS reprogramming are presented in the subsequent section on **Validation by geodesic interpolation**. We performed extensive sensitivity analysis to show that our temporal couplings produce valid interpolations over a wide range of parameter settings perturbations to the data (down sampling cells or reads). See **QUANTIFICATION AND STATISTICAL ANALYSIS** for this sensitivity analysis.

**[0435]**   2. Computing transport maps

**[0436]**   Recall that for any pair of time points we computed a transport plan that minimizes the expected cost of re-distributing mass, subject to constraints involving the relative growth rate (see **Modeling developmental processes with optimal transport** for a precise statement of the optimization problem). To compute these transport matrices, we needed to specify a cost function, numerical values for the regularization parameters, and (optionally) an initial estimate for the relative growth rate.

**[0437]**   2.1 Cost function

**[0438]**   To compute the cost of transporting each individual point $x$ from time $t \backslash$ to position^ at time $t_2$, we first performed principal components analysis (PCA) on the data from this pair of time points to reduce to 30 dimensions. This dimensionality reduction was performed separately for each pair of adjacent time points. We defined the cost function to be squared Euclidean distance in this 'local-PCA space'.

**[0439]**   Finally, we normalized the cost matrix by dividing each entry by the median cost for that time interval. Here the cost matrix was the matrix with entries $Cy = c(x_i, y_j)$ for each $x_i$ form time $t_i$ and $y_j$ at time $t_2$. This rescaling of the cost allowed us to refer to specific numerical values of the regularization parameters, without worrying about the global scale of distances.

**[0440]**   2.2 Regularization parameters

**[0441]**   The optimization problem (4) involved three regularization parameters:

**[0442]**      1. The *entropy* parameter $E$ controlled the entropy of the transport map. An extremely large entropy parameter gave a maximally entropic transport map, and an extremely small entropy parameter gave a nearly deterministic transport map. The default value was 0.05.

**[0443]**      2. $\lambda\backslash$ controlled the degree to which transport was unbalanced along the rows. Large values of $\lambda\backslash$ imposed stringent constraints related to relative growth rates. Small values of $\lambda\backslash$ gave the algorithm more flexibility to change the relative growth rates in order to improve the transport objective. The default value was 1. To visually inspect the degree of unbalancedness, we recommend plotting the input row-sums vs the output row-sums of the transport map (See FIGs. 30A-30G).

**[0444]**      3. $\lambda_1$ controlled the degree to which transport is unbalanced along the columns. The default value was $\lambda_2 = 50$. This large value essentially imposed equality constraints for the column marginals. A smaller value of $\lambda_2$ would allow different amounts of mass to transport to some cells at time $t_2$. We recommend keeping a large value for $\lambda_2$ so that the results are balanced along the columns. To visually inspect the degree of unbalancedness, one can plot the input column-sums vs the output column-sums of the transport map.

**[0445]**      As we demonstrate in **QUANTIFICATION AND STATISTICAL ANALYSIS,** our validation results were stable over a wide range of values for $E$ and $\lambda\backslash$.

**[0446]**      2.3 Estimating relative growth rates

**[0447]**      Our method solved the optimization problem (4) several times, using the output row-sums of the optimal transport map $\pi\tilde{\imath}1,\ddot{\imath}2$ as a new estimate for the relative growth rate function $\hat{g}(x)$. By default, we initialize with $g(x) = 1$, so that all cells growed at the same rate. With some prior knowledge of growth rates (e.g. based on gene signatures of proliferation and apoptosis), this could be incorporated in the initial estimate for $\hat{g}(x)$. For our reprogramming data, we showed how we formed an initial estimate for relative growth rates in **Estimating growth and death rates and computing transport maps**.

**[0448]**      3 Ancestors, descendants, and trajectories

**[0449]**      Recall that the transport map $\hat{\pi}_t1, t2$ connecting cells from time $t\backslash$ to cells from time $t_2$ has a row for each cell $x$ at time $t\backslash$ and a column for each cell $y$ at time $t_2$. Each row specifies the *descendant distribution* of a single cell x from time $t_1$. The descendant mass is the sum of all the entries across a row. This row-sum was proportional to the number of descendants that x would

contribute to the next time point. Intuitively, the descendant distribution specified which cells at time $t_2$ were likely to be descendants of $x$ (see section 4.1 of **Modeling developmental processes with optimal transport** for the formal definition of descendants in a developmental process).

[0450]     Similarly, each column specified the ancestor distribution of a cell $y$ from time $t_2$. The ancestor mass was usually the same for each cell $y$. The ancestor distribution told us which cells at time $t_1$ were likely to give rise to the cell $y$.

[0451]     Given a set of cells C, we computed the descendant distribution of the entire set by adding the descendant distributions of each cell in the set. This was computed efficiently via matrix multiplication as follows: Let $S_1$ donote all the cells from time point tl, and let

$$p(x) = \begin{cases} 1 & x \in C \\ 0 & \text{otherwise} \end{cases}$$

[0452]     denote the uniform distribution on $C \subset S$. The descendant distribution of C was given by $\pi\hat{i}_{1},\ddot{i}2\, p$. One could compute ancestor distributions in a similar way

[0453]     After computing the trajectory to or from a cell set $C$ (in the form of a sequence of ancestor and descendant distributions), we computed trends in expression for any gene or gene signature along the trajectory. For each time point, we simply computed the mean expression weighting each cell according to the probability distribution defined by the ancestor or descendant distribution.

[0454]     4. Learning gene regulatory models

[0455]     In this section we describe two strategies to summarize the transport maps by learning models of gene regulation. The first model we describe is a simple local enrichment analysis to identify transcription factors (TFs) enriched in ancestors of a set of cells. The second model is motivated by the dynamical systems formulation of optimal transport, as described above in Section 4.3.

[0456]     4.1 Local model: TF enrichment analysis of top ancestors

[0457]     We performed local enrichment analysis as follows. Given a set of cells $C$ at time $t_2$, we first computed the ancestor distribution of $C$ at an earlier time $t_1$, as described in Section 3 above. We then selected cells contributing the most mass to the ancestor distribution, until a certain amount of mass was accounted for (e.g. 30% of the ancestor mass). We referred to these

as the *top ancestors* at time $t_1$ of the cell set $C$. Finally, we compared the top ancestors to a null set of cells from the same time point. For example, this null cell set could be:

[0458]    all cells except for the top ancestors,

[0459]    the bottom *ancestors* (defined to be all cells except for the top ancestors of a less-strict cut-off),

[0460]    the bottom ancestors restricted to a specialized subset (e.g. all other trophoblasts when $C$ is a specific subset of trophoblasts like spongiotrophoblasts).

[0461]    4.2 Global model: learning a cell-autonomous gradient flow

[0462]    To learn a simple description of the temporal flow, we assumed that a cell's trajectory was cell-autonomous and, in fact, depended only on its own internal gene expression. We knew this was wrong as it ignored paracrine signaling between cells, and we returned to discuss models that include cell-cell communication at the end of this section. However, this assumption is powerful because it exposes the time-dependence of the stochastic process $P_t$ as arising from pushing an initial measure through a differential equation:

$$\dot{x} = f(x)_. \qquad\qquad\qquad (10)$$

[0463]    Here $f$ was a vector field that prescribes the flow of a particle $x$ (see FIG. 4 for a cartoon illustration of a distribution flowing according to a vector field). Our biological motivation for estimating such a function $f$ was that it encoded information about the regulatory networks that created the equations of motion in gene-expression space.

[0464]    We set up a regression to learn a regulatory function $f$ that models the fate of a cell at time $t_{,+1}$ as a function of its expression profile at time $t_{,}$. Our approach involved sampling pairs of points using the couplings from optimal transport:

[0465]    For each pair of time points $t_{,}, t_{,+1}$, we sampled pairs of cells $\left(X_{t_i}, X_{t_{i+1}}\right)$ from the joint distribution specified by the transport map $\hat{\gamma}_{,t_{i+1}}$.

[0466]    Using the training data generated in the first step, we set up the following regression:

$$\min_{f \in \mathcal{F}} \quad \mathbb{E}_{t_i, t_{i+1}} \left\| X_{t_{i+1}} - f(X_{t_i}) \right\|^2,$$

[0467]    where $\wedge$ was a rectified-linear function class defined in terms of a specific generalized logistic function $\ell : \gamma \mapsto \mathbb{R}$:

$$\ell(x; k, b, y_0, XQ) = \frac{k y_0}{y_0 + (k - y_0)e^{-b(x - x_0)}},$$

**[0468]**   where $k, b, y_0, x_0 \in \mathbb{R}$ were parameters of the generalized logistic function $\ell(x)$.

**[0469]**   We define a function class $-\widehat{}$ consisting of functions $f : \mathbb{R}^G \to \mathbb{R}^G$ of the form

$$f(x) = Ui(WTx),$$

**[0470]**   where $i$ was applied entry-wise to the vector $WTx \in \mathbb{R}^M$ to obtain a vector that we multiplied against $U \in \mathbb{R}^{?\times M}$. Here $T \in \mathbb{R}^{G_{TF} \times G}$ denoted a projection operator that selected only the coordinated of $x$ that were transcription factors, and GTF was the number of transcription factors. This gave a set of low-rank, linear functions with sparse factors. Each rank-1 component was interpreted as a regulatory module of transcription factors acting on a module of regulated genes.

**[0471]**   We set up the following optimization over matrices

$$\min_{U,W} \ \mathbb{E}_r \left\| \frac{X_{t_i} - X_{t_{i+1}}}{\Delta_t} - U\ell(WTX_{t_i}) \right\|^2 + \eta_1 \|U\|_1 + \eta_2 \|W\|_1 + \eta_3 \|W\|_2^2 \qquad (11)$$

**[0472]**   s.t.  $U \geq 0.$

**[0473]**   where $(X_{t_i}, Xti+\backslash )$ is a pair of random variables distributed according to the normalized transport map $r$, and $\|U\|_1$ denotes the sparsity-promoting $\ell_1$ norm of $U$, viewed as a vector (that is, the sum of the absolute value of the entries of $U$). Each rank one component (row of $U$ or column of $W$) gives us a group of genes controlled by a set of transcription factors. The regularization parameters $\eta_1$ and $\eta_2$ control the sparsity level (i.e. number of genes in these groups).

**[0474]**   **Implementation:**  We designed a stochastic gradient descent algorithm to solve (11). Over a sequence of epochs, the algorithm sampled batches of points $(Xt_i, Xti+\backslash )$ from the transport maps, computed the gradient of the loss, and updates the optimization variables $U$ and $W$. The batch sizes were determined by the Shannon diversity of the transport maps: for each pair of consecutive time points, we computed the Shannon diversity $S$ of the transport map, then randomly sampled $\max(S 10^{-5}, 10)$ pairs of points to add to the batch. We ran for a total of 10, 000 epochs.

**[0475]**   **Cell non-autonomous processes:**  We concluded our treatment of gene regulatory networks by discussing an approach to cell-cell communication. Note that the gradient flow (10)

only made sense for cell autonomous processes. Otherwise, the rate of change in expression $\dot{x}$ was not just a function of a cell's own expression vector x*(t),* but also of other expression vectors from other cells. We accommodated cell non-autonomous processes by allowing *f* to also depend on the full distribution P?:

$$\frac{dx}{dt} = f(x, \mathbb{P}_t).$$

(12)

**[0476]** Concretely, we could allow *f* to depend on the mean expression levels of specific genes (expressed by any cell) encoding, for example, secreted factors or direct protein measurements of the factors themselves.

**[0477]** 5. Geodesic interpolation

**[0478]** Optimal transport provided an elegant way to interpolate distribution-valued data, analogous to how linear regression can be used to interpolate numerical or vector-valued data. Given two numerical data- points, a simply way to interpolate was to connect them with a line; this was the shortest path connecting the observed data. Given two distributions, we interpolated by finding the shortest path in the space of distributions. To do this we needed a notion of distance between distributions, and for this we use the metric induced by optimal transport. This metric space was called Wasserstein space, and this form of interpolation was called geodesic interpolation (Villani, 2008).

**[0479]** We derived a modified version of geodesic interpolation that took into account cell growth. Ordinarily, an interpolating distribution was computed by first computing a transport map between the distributions, and then connecting each point in the first distribution to points in the second according to the transport map. Finally, an interpolating point cloud was produced by from the midpoints of those line segments. (More generally, instead of taking just midpoints, one could also construct a family of interpolations that sweep from the first distribution to the second). We extended this framework to accommodate growth by changing the mass of the point we placed at the midpoint (to account for the fact that cells would have a different number of descendants at time *t\* than they would at time *ti).*

**[0480]** Specifically, to interpolate at time $s \in (t_1, t_2)$ we first renormalize the rows of the transport map so they sum to roughly $\dfrac{\hat{g}(x)^{s-t_1}}{\int \hat{g}(x)^{s-t_1} d\mathbb{P}_{t_1}}$ instead of $\dfrac{\hat{g}(x)^{t_2-t_1}}{\int \hat{g}(x)^{t_2-t_1} d\mathbb{P}_{t_1}(x)}$ *—This took*

into account the descendant mass each cell would have by time $s$ instead of by time $t_2$. We then sampled points zi, . . . , $z_N$ as follows:

[0481]    1. Sampling a pair of points (x, $y$) from the joint distribution specified by the transport map.

[0482]    2. Identifying the point

$$z = \alpha x + (1 - \alpha)y$$

along the line segment connecting $x$ and $y$. Here $a$ is given by $s = ah + (1 - \alpha)t_2$.

[0483]    By repeating the steps above, we accumulate a point-cloud of points zi, . . . , $z_N$. Finally, we define the interpolating distribution as

$$\hat{P}(s) = \frac{1}{N}\sum_{i=1}^{N}\delta_{z_i}.$$

[0484]    Equipped with this notion of interpolation, we tested the performance of optimal transport by comparing the interpolated distribution to held-out time points. Using the data from time ti and ti+2, we interpolated to estimate the distribution Pti+1 . We then computed the Wasserstein distance between the interpolated distribution and the observed distribution. We compared this distance to a null model generated from the independent coupling where we sample pairs $(x, y)$ independently $x \sim \hat{P}_{t_i}$ and $y \sim \hat{P}_{t_{i+2}}$ in step 1 above. We also compared the interpolated distance to distance between batches of $P_{t_{i+1}}$ . Optimal transport was performing well if the interpolated point cloud was as close to the batches of the held out time point as the batches were to each other, and the null-interpolated point cloud was farther away.

[0485]    **Bibliography**

•    Ambrosio, L., Gigli, N., and Savare, G. (2005). Gradient Flows: In Metric Spaces and in the Space of Probability Measures. Lectures in Mathematics. ETH Zürich. Birkhäuser Basel.

•    Bastian, M., Heymann, S., Jacomy, M., et al. (2009). Gephi: an open source software for exploring and manipulating networks. Icwsm, 8:361-362.

•    Beygelzimer, A., Kakadet, S., Langford, J., Arya, S., Mount, D., Li, S., and Li, M. S. (2015). Package FNN.

•    Chizat, L., Peyre', G, Schmitzer, B., and Vialard, F.-X. (2017).   Scaling algorithms for unbalanced transport problems. Mathematics of Computation.

- Cuturi, M. (2013). Sinkhorn distances: Lightspeed computation of optimal transportation distances. In

- Neural Information Processing Systems (NIPS).

- Jacomy, M., Venturini, T., Heymann, S., and Bastian, M. (2014). Forceatlas2, a continuous graph layout algorithm for handy network visualization designed for the gephi software. PloS one, 9:e98679.

- Lėonard, C. (2014). A survey of the schrödinger problem and some of its connections with optimal transport. Discrete and Continuous Dynamical Systems - Series A (DCDS-A), 34(4):1533-1574.

- Porpiglia, E., Samusik, N., Van Ho, A. T., Cosgrove, B. D., Mai, T., Davis, K. L., Jager, A., Nolan, G. P., Bendall, S. C, Fantl, W. J., et al. (2017). High-resolution myogenic lineage mapping by single-cell mass cytometry. Nature Cell Biol., 19:558-567.

- Richard Jordan, D. K. and Otto, F. (1998). The variational formulation of the fokker. SIAM J. Math. Anal., 29(1): 1-17.

- Samusik, N., Good, Z., Spitzer, M. H., Davis, K. L., and Nolan, G. P. (2016). Automated mapping of phenotype space with single-cell data. Nature methods, 13:493-496.

- Santambrogio, F. (2015). Optimal Transport for Applied Mathematicians: Calculus of Variations, PDEs, and Modeling. Progress in Nonlinear Differential Equations and Their Applications. Springer Inter- national Publishing.

- Schrodinger, E. (1932). Sur la theorie relativiste de l'electron et l'interpretation de la mecanique quan- tique. Ann. Inst. H. Poincare, 2:269-310.

- Villani, C. (2008). Optimal Transport Old and New. Springer.

- Zunder, E. R., Lujan, E., Goltsev, Y., Wernig, M., and Nolan, G. P. (2015). A continuous molecular roadmap to ipse reprogramming through progression analysis of single-cell mass cytometry. Cell Stem Cell, 16:323-337.

**[0486]**     **III Experimental   methods**

**[0487]**     1. Derivation of secondary MEFs

**[0488]**     OKSM secondary Mouse embryonic fibroblasts (MEFs) were derived from E13.5 female embryos with a mixed B6;129 background. The cell line used in this study was homozygous for ROSA26-M2rtTA, homozygous for a polycistronic cassette carrying *Oct4, Klf4,*

*Sox2,* and *Myc* at the *Col1a1* locus and homozygous for an EGFP reporter under the control of the *Oct4* promoter (Stadtfeld et al., 2010). Briefly, MEFs were isolated from E13.5 embryos from timed-matings by removing the head, limbs, and internal organs under a dissecting microscope. The remaining tissue was finely minced using scalpels and dissociated by incubation at 37°C for 10 minutes in trypsin-EDTA (Thermo Fisher Scientific). Dissociated cells were then plated in MEF medium containing DMEM (Thermo Fisher Scientific), supplemented with 10% fetal bovine serum (GE Healthcare Life Sciences), non-essential amino acids (Thermo Fisher Scientific), and GlutaMAX (Thermo Fisher Scientific). MEFs were cultured at 37°C and 4% $CO_2$ and passaged until confluent. All procedures, including maintenance of animals, were performed according to a mouse protocol (2006N000104) approved by the MGH Subcommittee on Research Animal Care.

[0489] 2. Derivation of Primary MEFs

[0490] Primary MEFs were derived from E13.5 embryos with a B6.Cg-Gt(ROSA)$^{26Sortml(rtTA*M2)Ja}$7JxB6;129S4-Pou5fl $^{tm2Jae}$/J background. The cell line was homozygous for ROSA26-M2rtTA, and homozygous for an EGFP reporter under the control of the Oct4 promoter. MEFs were isolated as mentioned above.

[0491] 3. Reprogramming assay

[0492] For the reprogramming assay, 20,000 low passage MEFs (no greater than 3-4 passages from isolation) were seeded in a 6-well plate. These cells were cultured at 37°C and 5% $CO_2$ in reprogramming medium containing KnockOut DMEM (GIBCO), 10% knockout serum replacement (KSR, GIBCO), 10% fetal bovine serum (FBS, GIBCO), 1% GlutaMAX (Invitrogen), 1% nonessential amino acids (NEAA, Invitrogen), 0.055 mM 2-mercaptoethanol (Sigma), 1% penicillin-streptomycin (Invitrogen) and 1,000 U/ml leukemia inhibitory factor (LIF, Millipore). Day 0 medium was supplemented with 2 $\mu$g/mL doxycycline Phase-1(Dox) to induce the polycistronic OKSM expression cassette. Medium was refreshed every other day. At day 8, doxycycline was withdrawn, and cells were transferred to either serum-free 2i medium containing 3 $\mu$M CHIR99021, 1 $\mu$M PD0325901, and LIF (Phase-2(2i)) (Ying et al., 2008) or maintained in reprogramming medium (Phase-2(serum)). Fresh medium was added every other day until the final time point on day 18. Oct4-EGFP positive iPSC colonies should start to appear on day 10, indicative of successful reprogramming of the endogenous Oct4 locus.

**[0493]**     4. Sample collection

**[0494]**     We profiled a total of 315,000 cells from two time-course experiments across 18 days in two different culture conditions: in the first we profiled ~65,000 cells collected over 10 time points separated by ~48 hours; in the second we profiled ~250,000 cells collected over 39 time points separated by ~12 hours across an 18-day time course (and every 6 hours between days 8 and 9). In the larger experiment, duplicate samples were collected at each time point. Cells were also collected from established iPSCs cell lines reprogrammed from the same MEFs, maintained either in Phase-2(2i) conditions or in Phase-2(serum) medium. For all time points, selected wells were trypsinized for 5 mins followed by inactivation of trypsin by addition of MEF medium. Cells were subsequently spun down and washed with IX PBS supplemented with 0.1% bovine serum albumin. The cells were then passed through a 40 micron filter to remove cell debris and large clumps. Cell count was determined using Neubauer chamber hemocytometer to a final concentration of 1000 cells/$\mu$l.

**[0495]**     5. Single-cell RNA-seq

**[0496]**     ScRNA-seq libraries were generated from each time point using the 10X Genomics Chromium Controller Instrument (10X Genomics, Pleasanton, CA) and Chromium™ Single Cell 3' Reagent Kits v1 (~65,000 cells experiment) and v2 (~250,000 experiment) according to manufacturer's instructions. Reverse transcription and sample indexing were performed using the CIOOO Touch Thermal cycler with 96-Deep Well Reaction Module. Briefly, the suspended cells were loaded on a Chromium controller Single-Cell Instrument to first generate single-cell Gel Bead-In-Emulsions (GEMs). After breaking the GEMs, the barcoded cDNA was then purified and amplified. The amplified barcoded cDNA was fragmented, A-tailed and ligated with adaptors. Finally, PCR amplification was performed to enable sample indexing and enrichment of the 3' RNA-Seq libraries. The final libraries were quantified using Thermo Fisher Qubit dsDNA HS Assay kit (Q32851) and the fragment size distribution of the libraries were determined using the Agilent 2100 BioAnalyzer High Sensitivity DNA kit (5067-4626). Pooled libraries were then sequenced using Illumina Sequencing. All samples were sequenced to an average depth of 87 million paired-end reads per sample (see Experimental Methods), with 98 bp on the first read and 10 bp on the second read. In the larger experiment, we profiled 259,155 cells to an average depth of 46,523 reads per cell.

**[0497]**     6. Lentivirus vector construction and particle production

**[0498]**     To test whether transcription factors (TFs) improve late-stage reprogramming efficiency, we generated lentiviral constructs for the top candidates *Zfp42,* and *Obox6.* cDNAs for these factors were ordered from Origene *(ZJp42-*MG203929, and *Obox6-MR2* 15428) and cloned into the FUW Tet-On vector (Addgene, Plasmid #20323) using the Gibson Assembly (NEB, E2611S). Briefly, the cDNA for each TF was amplified and cloned into the backbone generated by removing *Oct4* from the FUW-Teto-Oc *t4* vector. All vectors were verified by Sanger sequencing analysis. For lentivirus production, HEK293T cells were plated at a density of 2.6xlO cells/well in a 10cm dish. The cells were transfected with the lentiviral packaging vector and a TF-expressing vector at 70-80% growth confluency using the Fugene HD reagent (Promega E2311), according to the manufacturer's protocols. At 48 hours after transfection, the viral supernatant was collected, filtered and stored at -80°C for future use.

**[0499]**     7. Reprogramming efficiency of secondary MEFS together with individual TFs

**[0500]**     We sought to determine the ability of the candidate TFs to augment reprogramming efficiency in secondary MEFs; the use of secondary MEFs for reprogramming overcomes limitations associated with random lentiviral integration events at variable genomic locations. Briefly, secondary MEFs were plated at a concentration of 20,000 cells per well of a 6-well plate. Cells were infected with virus containing Zfp42, Obox6, or an empty vector and maintained in reprogramming medium as described above. At day 8 after induction, cells were switched to either Phase-2(2i) or Phase-2(serum). On day 16, reprogramming efficiency was quantified by measuring the levels of the EGFP reporter driven by the endogenous *Oct4* promoter. FACS analyses was performed using the Beckman Coulter CytoFLEX S, and the percentage of Oct4-EGFP⁺ cells was determined. Triplicates were used to determine average and standard deviation.

**[0501]**     8. Reprogramming efficiency of primary MEFS with individual TFs and OKSM

**[0502]**     We also independently tested the performance of TFs in primary MEFs. To this end, lentiviral particles were generated from four distinct FUW-Teto vectors, containing Oct4, Sox2, Klf4, and Myc, previously developed in the Jaenisch lab. MEFs from the background strain B6.Cg-Gt(ROSA)26Sor $_{tmi(rtTA*M2)j\,ae}$/j _ B6;129S4-Pou5fl $^{tm2Jae}$/J were infected with these lentiviral particles, together with a lentivirus expressing tetracycline-inducible Zfp42, Obox6 or

no insert. Infected cells were then induced with 2 $\mu$g/mL doxycycline in ESC reprogramming medium (day 0). At day 8 after induction, cells were switched to either Phase-2(2i) or Phase-2(serum). On day 16, the number of Oct4-EGFP⁺ colonies were counted using a fluorescence microscope. Triplicates for each condition used to determine average values and standard deviation.

**[0503]     IV. Preparation of expression matrices**

**[0504]**     To compute an expression matrix from scRNA-Seq data, we aligned sequenced reads to obtain a matrix $U$ of UMI counts, with a row for each gene and a column for each cell. To reduce variation due to fluctuations in the total number of transcripts per cell, we divide the UMI vector for each cell by the total number of transcripts in that cell. Thus we define the expression matrix $E$ in terms of the UMI matrix $U$ via:

$$E = \frac{U_{ij}}{\sum_{i=1}^{G} U_{ij}} \times 10^4.$$

**[0505]**     In our subsequent analysis, we make use of two variance-stabilizing transforms of the expression matrix $E$. In particular, we define

1. $\tilde{E}$ to be the log-normalized expression matrix. The entries of $\tilde{E}$ are obtained via
$$\tilde{E} = logiEi_j + 1)$$
2. $\bar{E}$ to be the truncated expression matrix. The entries of $\bar{E}$ are obtained by capping the entries of $\bar{E}$ at the 99.5% quantile.

**[0506]**     When we refer to an expression profile, by default we refer to a column of $\tilde{E}$ unless otherwise specified.

**[0507]**     1. Aligning reads

**[0508]**     The 98 bp reads were aligned to the UCSC mmlO transcriptome, and a matrix of UMI counts was obtained using Cellranger from the 10X Genomics pipeline (v2.0.0) with default parameters (https://support. lOxgenomics.com/single-cell-gene-expression/software/pi pelines/latest/installati on). Quality control metrics about barcoding and sequencing such as the estimated number of cells per collection and the median number of genes detected across cells are summarized in Table 14. To estimate expression of exogenous OKSM factors from OKSM cassette, we extracted RBGpA sequence (839 bp) from the OKSM cassette FASTA file, and generated a reference using the mkref function from the Cellranger pipeline.

**[0509]**     2. Downsampling and filtering expression matrix

**[0510]**     The expression matrix was downsampled to 15,000 UMIs per cell. Cells with less than 2000 UMIs per cell in total and all genes that were expressed in less than 50 cells were discarded, leaving 251,203 cells and G= 19,089 genes for further analysis. The elements of expression matrix were normalized by dividing UMI count by the total UMI counts per cell and multiplied by 10,000 i.e. expression level is reported as transcripts per 10,000 reads.

**[0511]**     3. Selecting variable genes

**[0512]**     We used the function MeanVarPlot from the Seurat package (v2.1.0) (Satija et al., 2015) to select 1479 variable genes. First, we divided genes into 20 bins based on their average expression levels across all cells. Second, we computed Fano factor of gene expression in each bin and then z-scored. The Fano factor, defined as the variance divided by the mean, was a measure of dispersion. Finally, by thresholding the z-scored dispersion at 1.0, we obtained a set of 1479 variable genes. After selecting variable genes, we created a variable gene expression matrix by renormalizing as described above.

**[0513]**     **V. Visualization; force-directed layout embedding**

**[0514]**     In this section we introduced our two dimensional visualization technique based on force-directed layout embedding (FLE) (Bastian et al., 2009; Jacomy et al., 2014). FLE was large-scale graph visualization tool which simulated the evolution of a physical system in which connected nodes experience attractive forces, but unconnected nodes experience repulsive forces. It better captured global structures than tSNE. Initial FLE algorithms used simple electrostatic and spring forces, but modern FLE algorithms allowed for more elaborate interactions that could depend on the degree of nodes or included gravity terms that attracted all nodes to the center (this was especially important for disconnected graphs, which would otherwise fly apart). Starting from a random initial position of vertices, the network of nodes evolved in such a manner that at any iteration a new position of vertices was computed from the net forces acting on them.

**[0515]**     We applied FLE to visualize the nearest neighbor graph generated from our data.

**[0516]**     **Implementation:** Our visualization took as input the expression matrix of highly-variable genes, selected as described in the previous section of the STAR Methods. First, we reduced to 100 dimensions by computing a 100 dimensional diffusion component embedding of the dataset using SCANPY (vO.2.8) with default parameters. Second, for each cell we computed

its 20 nearest neighbors in 100-dimensional diffusion component space to produce a nearest neighbor graph. For this step, we used the approximate k-NN algorithm Annoy from the R package RCPPANNOY (vO.0.10). Finally, we computed the force-directed layout on the k-NN graph using the ForceAtlas2 algorithm (Jacomy et al., 2014) from the Gephi Toolkit (vO.9.2) (Bastian et al., 2009).

[0517]      **VI. Creating gene signatures and cell sets**

[0518]      1. Gene signatures

[0519]      We then constructed curated gene signatures from various databases of gene signatures. Given a set of genes, we scored cells based on their gene expression. In particular, for a given cell we computed the z-score for each gene in the set. We then truncated these z-scores at 5 or -5, and defined the signature of the cell to be the mean z-score over all genes in the gene set.

[0520]      The table below summarizes the sources from which we obtained signatures. In two cases (neural identity and epithelial identity), we constructed signatures manually using marker genes. A pluripotency gene signature was determined in this work using the pilot dataset. We performed differential gene expression analysis between two groups of cells: mature iPSCs and cells along the time course D0 to D16 and took the top 100 genes with increased expression in mature iPSCs. A proliferation gene signature was obtained by combining genes expressed at Gl/S and G2/M phases.

[0521]      In several places, we also computed gene signatures based on co-expression with a given gene of interest. For instance, in the stromal region we noticed several genes *(Cxcll2, Ifitml,* and *Matn4)* with expression patterns that were distinct from a signature of long-term cultured MEFs (FIG. 3 ID). For each gene, we computed a co-expression signature by finding the set of genes with expression levels in stromal cells that were $>\backslash 5\%$ correlated with the gene of interest. We found that these gene signatures were significantly overlapping (p-value < 0.01, hypergeometric test) with signatures of stromal cells in neonatal muscle and neonatal skin in the Mouse Cell Atlas. Similarly, in the neural region we derived signatures of genes co-expressed with *Gadl* and with *Slcl7a6* (FIG. 33C). These signatures significantly overlapped signatures of inhibitory and excitatory neurons, respectively, derived from the Allen Brain Atlas.

| Gene Signature | Source |
|---|---|
| M EF identity | (Chen et al., 2013; Han et al., 2018; Lattin et al., 2008) |
| Pluripotency | This work. |
| Proliferation | (Tirosh et al., 2016) |
| ER stress | GO:0034976, Biological Process Ontology |
| Epithelial identity | This work.<br>Marker genes: (Li et al., 2010; Takaishi et al., 2016; Whiteman et al., 2014) |
| ECM rearrangement | GO:0030198, Biological Process Ontology |
| Apoptosis | Hallmark P53 Pathway, MSigDB |
| Senescence | (Coppe et al., 2010) |
| Neural identity | This work.<br>Marker gene sources: (Fonseca et al., 2013; Gouti et al., 2011; Kan et al., 2004; Lazarov et al., 2010; Sakakibara et al., 2001; Sansom et al., 2009; Watanabe et al., 2017) |
| Trophoblast | (Han et al., 2018) |
| X reactivation | chromosome X |
| XEN | (Lin et al., 2016) |
| Trophoblast progenitors | (Han et al., 2018) |
| Spiral Artery Trophpblast Giant Cells | (Han et al., 2018) |
| Oligodendrocyte precursor cells (OPC) | (Tasic et al., 2016) |
| Astrocytes | (Tasic et al., 2016) |
| Cortical Neurons | (Tasic et al., 2016) |
| RadialGlia-Id3 | (Han et al., 2018) |
| RadialGlia-GdflO | (Han et al., 2018) |
| RadialGlia-Neurog2 | (Han et al., 2018) |
| Long-term M EFs | (Han et al., 2018) |

| Embryonic mesenchyme | (Han et al., 2018) |
| --- | --- |
| Cxcl 2 co-expressed | This work. |
| Ifitm I co-expressed | This work. |
| Matn4 co-expressed | This work. |
| 2,4,8, 16,32-cell | (Goolam et al., 2016) |

**[0522]**       2. Cell sets

**[0523]**       Using the gene signatures described above, we created coarse cell sets defining the broad regions of the landscape (iPSC, Trophoblast, Neural, Stromal, Epithelial, and MET), and cell subtype sets defining different cell types within a region (stromal, trophoblast, and neural subtypes, along with 2- through 32-cell stages).

**[0524]**       To define the coarse cell sets, we first computed a rough partitioning of the landscape by clustering cells using the Louvain method of spectral clustering to obtain 65 cell clusters using k=5 nearest neighbors (FIG. 34A). By examining signature score activity levels over clusters, we grouped several clusters to form cell sets for the iPSC, Stromal and Neuronal regions. Because our densely sampled data did not always segregate into distinct clusters, we defined some additional coarse cell sets by signature scores. We defined the trophoblast cell set to include all cells with Trophoblast signature greater than 0.7. We defined the epithelial cell set to include all cells with epithelial identity signature greater than 0.8, minus all cells included in other cell sets (mostly removing the trophoblasts with epithelial signature). Finally, we defined the MET Region as the ancestors of iPS, Trophoblast, Neural and Epithelial cells. In particular, we computed the top ancestors of each major cell set, then merged these cell sets and removed the cells *in* each major cell set.

**[0525]**       Within the Stromal, Trophoblast, Neural and iPSC cell sets, we then conducted more sensitive statistical tests for cell subtype signatures. We did this by calculating empirical p-values for the subtype signature score for each (region-specific) subtype in each cell. In each of 100,000 permutation trials, we randomly and independently shuffled the expression levels of each gene across the cells within a region. In each cell, we then computed signature scores in the permuted data, and generated p-values by determining the frequency at which the permuted score was greater than the original score. While the results shown in figures and discussed in the main text

were based on shuffling genes across cells, we similarly permuted the expression levels within each cell, and found consistent results. Finally, we controlled for multiple hypothesis testing by calculating FDR q-values, and used a threshold FDR of 10% to define cell subtype sets.

**[0526]** **VII. Estimating growth and death rates and computing transport maps**

**[0527]** 1. Initial estimate of growth rates

**[0528]** We formed an initial estimate of the relative growth rate as the expectation of a birth-death process on gene expression space with birth-rate $\beta(\chi)$ and death rate $\delta(\chi)$ defined in terms of expression levels of genes involved in cell proliferation and apoptosis. Multi-state birth-death processes had been used before to model growth, death, and transitions in iPS reprogramming (Liu et al., 2016). A birth-death process was a classical model for how the number of individuals in a population could vary over time. The model was specified in terms of a birth rate $\beta$ and death rate $\delta$: During a time interval $\Delta t$, the probability of a birth was $\beta \Delta \ddot{\imath}$ and the probability of a death was $\delta \Delta \ddot{\imath}$. The doubling time for a birth death process was defined as follows. Starting with N(0) = n, the time $\tau$ it would take to get to an expected population size of $EN(t) \ = \ 2n$ is

$$\tau = \frac{\ln 2}{\beta - \delta}$$

**[0529]** The half-life could be computed in a similar way. We applied a sigmoid function to transform the proliferation score into a birth rate. The sigmoid function smoothly interpolated between maximal and minimal birth rates. We specified the maximal birth rate to be $\beta_{MAX} \ = \ 1.7$. Therefore, the fastest cell doubling time is

$$\frac{\ln 2}{1.7} \approx 0.41 \ days \ \approx \ 9.6 \ hours,$$

by the doubling time equation above. We defined the minimal birth rate as $\beta_{M\ddot{\imath}N} = 0.3$. Therefore the slowest cell doubling time is

$$\frac{\ln 2}{0.3} = 2.3 \ days \ = \ 55 \ hours.$$

**[0530]** Similarly, we transformed the apoptosis signature into an estimate of cellular death rates by applying a sigmoid function to smoothly interpolate between minimal and maximal allowed death rates. We defined the minimal death rate parameter to be $\delta_{MIN} = 0.3$, and the maximal death rate parameter as $\delta_{MAX} = 1.7$. By the calculations above, these correspond to half-lifes of 55 and 9.6 hours respectively.

**[0531]**      2. Learning growth rates and computing transport maps

**[0532]**      Using the growth rates defined in the previous section as an initial estimate, we computed transport maps and automatically improved these growth rates using the Waddington-OT software package (see Section *Computing transport maps).* For the cost function, we used squared Euclidean distance in 30 dimensional local PCA space computed on the variable gene data from the relevant pair of time points. We used the following parameter settings:

$$e = 0.05, \lambda_1 = 1, \lambda_2 = 50, \text{growth\_iters} = 3.$$

**[0533]**      The parameters $\lambda_1$ and $\lambda_2$ control the degree to which the row-sums and column-sums were unbalanced. A larger value of $\lambda_1$ induced a greater correlation between the input and output growth rates. The Waddington-OT package iterated the procedure of computing transport maps based on input growth rates, and then using the output growth rates as new input growth rates to recompute transport maps. We ran this for growth_iters = 3 total iterations.

**[0534]**      This gave us a set of transport maps between each pair of time points, which could be used to estimate the temporal coupling. From this estimate of the temporal coupling, we computed ancestor and descendant distributions to each of the major cell sets defined in the previous section.

**[0535]**      **VIII. Regulatory analysis**

**[0536]**      We performed regulatory analysis to identify modules of transcription factors regulating modules of genes with our global regulatory model from the Waddington-OT software package, described in Section *Learning gene regulatory models*. The optimization began by specifying the number of gene modules, and establishing an initial estimate for each. We used spectral clustering to initialize the modules: genes were clustered into 50 sets, with one module corresponding to each set, and weights set to 0 for genes outside the set, and 1 for genes within the set.

**[0537]**      We then specified a time lag between TF and gene module expression. In order to test for potential regulatory interactions on different time scales, we computed global regulatory models with three time lags: 6hrs, 48hrs, and 96hrs. This allowed us to identify factors that were predictive several days in advance —for instance, Nanog is a very early predictor of pluripotency and was found to be associated with a pluripotency associated gene expression module in the 96 hour model —as well as those predictive on shorter time scales —for instance, we TFs that were

predictive of neural-associated expression modules in the 6 and 48 hour models, but did not find such predictive TFs in the 96 hour model.

**[0538]** Finally, we set regularization and stochastic block size parameters. Default values available in the code online were used in this study. Briefly, regularization parameters were tuned on small training datasets to enforce sparsity (11 penalties) and reduce model complexity (12 penalty) while still achieving a good fit (>60% correlation between predicted and observed expression) in training data. These parameters may be specifically tuned in new datasets. The stochastic block size and number of epochs were set according to available hardware resources.

**[0539]** **IX. Validation by geodesic interpolation**

**[0540]** We validated Waddington-OT by demonstrating that we could accurately interpolate the distribution of cells at held out time points. We applied geodesic interpolation (described in **Waddington-OT; Concepts and Implementation**) to our reprogramming data to predict the distribution of cells at each time point, using only the data from the previous and next time points. In other words, we sought to predict the distribution $P_{t_2}$ at time $t_2$ from the distributions at neighboring time points: $P_{t_1}$ and $P_{t_3}$ (FIGs. 24H, 30D). To determine a baseline for performance, we examined the distance between the two different batches of the held-out distribution (FIGs. 24H, 30D).

**[0541]** To compute the optimal transport coupling from $P_{t_1}$ to $P_{t_3}$, we used the Waddington-OT package with default parameters. For the cost function we computed 30 dimensional local PCA coordinates using only the points from time $t_1$ and $t_3$. We then embedded the data from time $t_2$ into the 30 dimensional local PCA space which was computed using only the data from time $t_1$ and $t_3$. Finally, we used Wasserstein-2 distance to compute distance between point clouds.

**[0542]** **X. Paracrine signaling**

**[0543]** To characterize potential cell-cell interactions between contemporaneous cells during reprogramming, we first collected a list of ligands and receptors found in the GO database. The set of ligands (415 genes) was a union of three gene sets from the following GO terms:

1) *cytokine activity* (GO:0005125),
2) *growth factor activity* (GO:0008083), and
3) *hormone activity* (GO:0005 179).

**[0544]**     The set of receptors (2335 genes) was defined by the GO term *receptor activity* (GO: 0004872). Next, we used a curated database of mouse protein-protein interactions (Mertins et al., 2017) and identified 580 potential ligand-receptor pairs.

**[0545]**     First, we defined an interaction score $I_{A;B;X;Y;t}$ as the product of (1) the fraction of cells ($E_{A;X;t}$) in cell-set A expressing ligand X at time t and (2) the fraction of cells (i¾;Y;t) in cell-set B expressing the cognate receptor Y at time $t$. We define the aggregate interaction score $I_{A;B;t}$ as a sum of the individual interaction scores across all pairs:

$$I_{A;B;t} = \sum_{All\ X\cdot Y\ pairs} I_{A;B;X;Y;t} = \sum_{All\ X\cdot Y\ pairs} F_{A;X;t}\ F_{B;Y;t}$$

**[0546]**     We depicted the aggregate interaction scores for all combinations of cell clusters in FIGs. 28B, 34B.

**[0547]**     Second, we sought to explore individual ligand-receptor pairs at a given day and condition between cell ancestors of interest. For this purpose we defined the interaction score $\wedge_{A;B;X;Y;t}$ as the product of (1) the average expression of the ligand X in ancestors at time t of a cell set A and (2) the average expression of the cognate receptor Y in ancestors at time t of a cell set B. Values of the interaction scores $I_{A;B;X;Y;t}$ are high for ubiquitously expressed ligands and receptors at a given day and may be nonspecific to a pair of cell ancestors of interest. Thus, we used permutations to generate an empirical null distribution of interaction scores. In each of the 10,000 permutations, we randomly shuffled the labels of cells and calculated the interaction score $I^s_{A;B;X;Y;t}$. We then standardized each ligand-receptor interaction score by taking the distance between the interaction score $I_{A;B;X;Y;t}$ and the mean interaction score in units of standard deviations from the permuted data

$$((I_{A;B;X;Y;t} - mean(I^S_{A;B;X;Y;t}))/sd(I^S_{A;B;X;Y;t})).$$

**[0548]**     We depicted examples of standardized interaction scores ranked by their values in FIGs. 28C-28E and 34C-34E. Replacement of the average expression of the ligand with the total expression of the ligand in the calculation of the standardized interaction score did not affect the results.

**[0549]**     **XI. Classification of differential genes along the trajectory to iPSCs**

**[0550]**     To identify differential genes along the successful trajectory to iPSCs we computed the average expression (TPM) of all 19,089 genes in ancestors of iPSCs. The average expression values were log2 transformed and we filtered out genes for which the difference between maximal and minimal expression value between day 0 and day 18 was less than 1, leaving 2311 genes for further analysis. The genes were classified into 15 groups by k-means clustering as implemented in the R package stats. To identify the number of clusters we applied a gap statistic (Tibshirani et al. 2001) using the function clusGap from R package cluster v2.0.6.

**[0551]**     We performed functional enrichment analysis on the identified gene clusters using the findGO.pl program from the HOMER suite (Hypergeometric Optimization of Motif Enrichment, v4.9. 1) (Heinz et al. 2010) with Benjamini and Hochberg FDR correction for multiple hypothesis testing (retaining terms at FDR < 0.05). All genes that passed quality-control filters were used as a background set.

**[0552]**     **XII. Identifying large chromosomal aberrations**

**[0553]**     We have previously developed methods to identify copy number variations (CNVs) in scRNA-Seq data from tumor samples (Patel et al., 2014; Tirosh et al., 2016). That analysis differed from our current study in two key aspects: (1) the data were based on full length scRNA-seq (SMART-Seq2), and sequenced to greater depth in each cell, and (2) there we could rely on the clonal expansion of CNVs to make it easier to identify recurring chromosomal aberrations.

**[0554]**     We performed three types of analysis to detect aberrant expression in large chromosomal regions. First, we searched cells with significant up- or down-regulation at the level of entire chromosomes. Second, we ran a coarse analysis to identify cells with significant net aberrant expression across windows spanning 25 broadly-expressed genes. Focusing on regions that were enriched for cells with significant aberrations found by this coarse filter, we then performed a more sensitive test to compute the significance of aberrations in each window in each cell.

**[0555]**     Empirical p-values and false discovery rates (FDRs) for both analyses were computed by randomly permuting the arrangement of genes in the genome, as described below. Permutations for both types of analysis were done as follows. In each of 100,000 permutations we randomly shuffled the labels of genes in the entire dataset, while preserving the genomic

284

coordinates of genes (with each position having a new label each time) and the expression levels in each cell (so that each cell has the same expression values, but with new labels). We then computed either whole chromosome or subchromosomal aberration scores for each cell.

[0556]   To identify whole-chromosome aberrations scores in each cell, we began by calculating the sum of expression levels in 25Mbp sliding windows along each chromosome, with each window sliding IMbp so that it overlapped the previous window by 24Mbp. For each window in each cell, we then calculated the Z-score of the net expression, relative to the same window in all other cells. We then counted the fraction of windows on each chromosome with an absolute value Z-score $> 2$. This fraction served as the whole-chromosome aberration score for each chromosome in each cell. To assign a p-value to the whole-chromosome score for cell(i) chromosome(j), we calculated the empirical probability that the score for cell(i) chromosome(j) in the randomly permuted data was at least as large as the score in the original data.

[0557]   Subchromosomal aberration scores were computed as follows. We began by identifying the 20% of genes with the most uniform expression across the entire dataset. This was done by calculating the Shannon Diversity $e^{-\frac{3}{4} E_{gc} \ln E_{gc}}$ for each gene $g$ (where $E_{g\,c}$ was the expression matrix as defined above in **Preparation of expression matrices**), and taking the 20% of genes with the largest values. Using these genes, we subset the expression matrix and renormalized by TPM, and then computed in each cell the sum of expression in sliding windows of 25 consecutive genes, with each window sliding by one gene and overlapping the previous window (on the same chromosome) by 24 genes. In each window, we calculated the Z-score relative to all cells at day 0. The net (coarse filter) subchromosomal aberration score for a cell was calculated as the 12-norm of the Z-scores across all windows. To assign a p-value to the subchromosomal aberration score for cell(i), we calculated the empirical probability that the score for cell(i) in the randomly permuted data was at least as large as the score in the original data.

[0558]   Finally, to identify the specific region(s) of genomic aberrations in each cell, we conducted a more sensitive test using just the cells in the stromal and trophoblast regions. Again using 25 housekeeping gene windows, we computed the average z-score of gene expression for genes in each window in each cell. We then compared the scores in all windows in all cells to

similar scores computed for each cell in 100,000 random permutation trials, and then assigned p-values based on the frequency of extremely high (gain) or low (loss) expression values.

**[0559]** For each of the aberration scores and associated p-values described above, we controlled for multiple hypothesis testing by calculating FDR q-values, using a false discovery threshold of 10%.

**[0560]**     **QUANTIFICATION AND STATISTICAL ANALYSIS**

**[0561]**     **I. Analyzing the stability of optimal transport**

**[0562]** To test the stability of our optimal transport analysis to perturbations of the data and parameter settings, we downsampled the number of cells at each time point, downsampled the number of reads in each cell, perturbed our initial estimates for cellular growth and death rates, and perturbed the parameters for entropic regularization and unbalanced transport. We found that our geodesic interpolation results are stable to a wide range of perturbations, summarized in the following table:

| Number of cells per batch | Number of UMIs Per cell | Max Growth $\beta_{MAX}$ | Min Growth $\beta_{MIN}$ | Max Death $\delta_{MAX}$ | Min Death $\delta_{MIN}$ | Entropy regularization $\epsilon$ | Unbalanced transport $\lambda$ |
|---|---|---|---|---|---|---|---|
| Down to: 200 | Down to: 1000 | 33 hrs to 5.5 hrs | None to 9.5 hrs | 33 hrs to 5.5 hrs | None to 9.5hrs | $5 \times 10^{-5}$ to 0.5 | 0.1 to 32 |

**[0563]** To generate this table, we ran geodesic interpolation with all but one of these settings fixed to default values. The default parameter values that we used were:

$$e = 0.05, \lambda_1 = 1, \lambda_2 = 50, \beta_{MAX} = 1.7, \delta_{MAX} = 1.7, \beta_{MIN} = 0.3, \delta_{MIN} = 0.3.$$

**[0564]** Moreover, by default we used all reads per cell and all cells per batch.

**[0565]**     **II. Performance of other methods**

**[0566]**     1. Monocle2

[0567] Monocle2 fitted the data into a graph without using prior information of the number of potential fates (Qiu et al., 2017).

[0568]    We ran Monocle2 (v2.8.0) with default parameters on a subset of our dataset containing 1,000 cells per time point. Running on our full dataset would require more RAM than we had access to.

[0569]    In our data, Monocle2 failed to distinguish iPS, neuronal-like, and trophoblast-like cells as distinct destinations (FIG. 35A-35B). It put together day 18 stromal cells and day 0 MEFs at the root of the tree, and placed iPS, neural-like and trophoblast-like cells on a different branch from cells in the MET Region. Moreover, because the program could incorporate temporal information, it returned a trajectory that was inconsistent with the measured temporal progression. The output of the program implied that day 0 MEF cells gave rise to day 18 stromal cells, which in turn gave rise to everything else.

[0570]    2. URD

[0571]    URD identified trajectories from a user-specified root to a set of user-specified tips by performing random walks according to a Markov diffusion kernel.

[0572]    We ran URD (vl.O) with default parameters on a subset of our dataset containing 1,000 cells per time point. Running on our full dataset would require more RAM than we had access to.

[0573]    In our data, URD predicted that all fates diverge extremely early, with stromal cells diverging from other cells soon after day 0; trophoblast-like cells diverging from neural-like and iPS cells as early as day 1; and neural-like and iPS cells diverging at day 2 (FIGs. 35A-35B). Additionally, URD failed to assign over half (51%) of the cells to any trajectory.

[0574]    Comparing the two branches for iPS and neural (FIGs. 35A-35B - segments 6 and 7) revealed no distinctive pattern between the supposedly divergent trajectories from day 3 - 8. The divergent trajectories appeared to be an artifact of the fact that the method requires a distinct branch point.

[0575]    Moreover, because the method did not incorporate growth rates, the transitions to iPS and Neural come disproportionately from stromal cells.

[0576]    **HI. Pilot study**

[0577]    In our pilot study, we collected 65,000 expression profiles over 16 days at 10 distinct time points (and 9 in serum). We compared results from the larger study to the pilot study in FIGs. 30A-30G, where we showed trends in expression along trajectories to each major cell set:

iPSCs, Neural-like, Trophoblast-like (placenta-like in pilot), and Stromal. We found that the expression trends were reasonably similar. Moreover, by comparing the ancestor divergence plots for the two studies, we found that in both studies the stromal population gradually diverged early in the time course and there was a sharp divergence of iPSC from Neural and Trophoblast just after removal of Dox at day 8.

[0578]    **Data and Software Availability**

[0579]    We have uploaded our data to NCBI Gene Expression Omnibus. The identification numbers are:

| | |
|---|---|
| Single cell RNA-seq raw data (pilot study) | GSE106340 |
| Single cell RNA-seq raw data | GSE1 15943 |

[0580]    Our software package is available on GitHub: https://github.com/broadinstitute/wot

[0581]    S

[0582]    **Reference Cited**

[0583]    1. C. H. Waddington, How animals develop. (New York, 1936).

[0584]    2. C. H. Waddington, The strategy of the genes; a discussion of some aspects of theoretical biology. (London, Allen & Unwin [1957], 1957).

[0585]    3. E. Z. Macosko et al., Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. Cell 161, 1202-1214 (2015).

[0586]    4. A. M. Klein et al., Droplet barcoding for single-cell transcriptomics applied to embryonic stem cells. Cell 161, 1187-1201 (2015).

[0587]    5. G. X. Zheng et al., Massively parallel digital transcriptional profiling of single cells. Nature communications 8, 14049 (2017).

[0588]    6. A. Tanay, A. Regev, Scaling single-cell genomics from phenomenology to mechanism. Nature 541, 331-338 (2017).

[0589]    7. A. Wagner, A. Regev, N. Yosef, Revealing the vectors of cellular identity with single-cell genomics. Nat Biotech 34, 1145-1 160 (2016).

[0590]    8. S. C. Bendall et al., Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. Cell 157, 714-725 (2014).

**[0591]**    9. C. Trapnell et al., The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. Nature biotechnology 32, 381-386 (2014).

**[0592]**    10. M. Setty et al., Wishbone identifies bifurcating developmental trajectories from single-cell data. Nature biotechnology 34, 637-645 (2016).

**[0593]**    11. E. Marco et al., Bifurcation analysis of single-cell gene expression data reveals epigenetic landscape. Proceedings of the National Academy of Sciences of the United States of America 111, E5643-5650 (2014).

**[0594]**    12. J. M. Polo et al., A molecular roadmap of reprogramming somatic cells into iPS cells. Cell 151, 1617-1632 (2012).

**[0595]**    13. Y. Buganim et al., Single-cell expression analyses during cellular reprogramming reveal an early stochastic and a late hierarchic phase. Cell 150, 1209-1222 (2012).

**[0596]**    14. S. M. Hussein et al., Genome-wide characterization of the routes to pluripotency. Nature 516, 198 (2014).

**[0597]**    15. P. D. Tonge et al., Divergent reprogramming routes lead to alternative stem-cell states. Nature 516, 192-197 (2014).

**[0598]**    16. J. O'Malley et al., High resolution analysis with novel cell-surface markers identifies routes to iPS cells. Nature 499, 88 (2013).

**[0599]**    17. X. Qiu et al., Reversed graph embedding resolves complex single-cell developmental trajectories. bioRxiv, 110668 (2017).

**[0600]**    18. S. C. Bendall et al., Single-cell trajectory detection uncovers progression and regulatory coordination in human B cell development. Cell 157, 714-725 (2014).

**[0601]**    19. R. Rostom, V. Svensson, S. Teichmann, G. Kar, Computational approaches for interpreting scRNA-seq data. FEBS letters, (2017).

**[0602]**    20. L. Haghverdi, F. Buettner, F. J. Theis, Diffusion maps for high-dimensional single-cell analysis of differentiation data. Bioinformatics 31, 2989-2998 (2015).

**[0603]**    21. L. Haghverdi, M. Buttner, F. A. Wolf, F. Buettner, F. J. Theis, Diffusion pseudotime robustly reconstructs lineage branching. Nat Meth 13, 845-848 (2016).

**[0604]**    22. K. Campbell, C. Yau, Ouija: Incorporating prior knowledge in single-cell trajectory learning using Bayesian nonlinear factor analysis. bioRxiv, (2016).

**[0605]** 23. R. Cannoodt et al., SCORPIUS improves trajectory inference and identifies novel modules in dendritic cell development. bioRxiv, (2016).

**[0606]** 24. J. D. Welch, A. J. Hartemink, J. F. Prins, SLICER: inferring branched, nonlinear cellular trajectories from single cell RNA-seq data. Genome Biology 17, 106 (2016).

**[0607]** 25. K. Street et al., Slingshot: Cell lineage and pseudotime inference for single-cell transcriptomics. bioRxiv, (2017).

**[0608]** 26. H. Matsumoto, H. Kiryu, SCOUP: a probabilistic model based on the Ornstein-Uhlenbeck process to analyze single-cell expression data during differentiation. BMC Bioinformatics 17, 232 (2016).

**[0609]** 27. S. Rashid, D. N. Kotton, Z. Bar-Joseph, TASIC: determining branching models from time series single cell data. Bioinformatics 33, 2504-2512 (2017).

**[0610]** 28. M. Zwiessele, N. D. Lawrence, Topslam: Waddington Landscape Recovery for Single Cell Experiments. bioRxiv, (2016).

**[0611]** 29. C. Weinreb, S. Wolock, B. K. Tusi, M. Socolovsky, A. M. Klein, Fundamental limits on dynamic inference from single cell snapshots. bioRxiv, (2017).

**[0612]** 30. C. Villani, Optimal transport: old and new. (Springer Science & Business Media, 2008), vol. 338.

**[0613]** 31. M. Cuturi, in Advances in neural information processing systems. (2013), pp. 2292-2300.

**[0614]** 32. L. Chizat, G. Peyre, B. Schmitzer, F.-X. Vialard, Scaling algorithms for unbalanced transport problems. arXiv preprint arXiv: 1607.05816, (2016).

**[0615]** 33. J. H. Levine et al., Data-Driven Phenotypic Dissection of AML Reveals Progenitor-like Cells that Correlate with Prognosis. Cell 162, 184-197 (2015).

**[0616]** 34. K. Shekhar et al., Comprehensive Classification of Retinal Bipolar Neurons by Single-Cell Transcriptomics. Cell 166, 1308-1323.el330 (2016).

**[0617]** 35. R. R. Coifman et al., Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. Proceedings of the National Academy of Sciences of the United States of America 102, 7426-7431 (2005).

**[0618]**      36. M. Jacomy, T. Venturini, S. Heymann, M. Bastian, ForceAtlas2, a continuous graph layout algorithm for handy network visualization designed for the Gephi software. PloS one 9, e98679 (2014).

**[0619]**      37. E. R. Zunder, E. Lujan, Y. Goltsev, M. Wernig, G. P. Nolan, A continuous molecular roadmap to iPSC reprogramming through progression analysis of single-cell mass cytometry. Cell Stem Cell 16, 323-337 (2015).

**[0620]**      38. C. Weinreb, S. Wolock, A. Klein, SPRING: a kinetic interface for visualizing high dimensional single-cell expression data. bioRxiv, (2016).

**[0621]**      39. K. Takahashi, S. Yamanaka, Induction of pluripotent stem cells from mouse embryonic and adult fibroblast cultures by defined factors. cell 126, 663-676 (2006).

**[0622]**      40. J. Yu et al., Induced pluripotent stem cell lines derived from human somatic cells. Science 318, 1917-1920 (2007).

**[0623]**      41. J. Shu et al., Induction of pluripotency in mouse somatic cells with lineage specifiers. Cell 153, 963-975 (2013).

**[0624]**      42. P. Hou et al., Pluripotent Stem Cells Induced from Mouse Somatic Cells by Small-Molecule Compounds. Science 341, 651-654 (2013).

**[0625]**      43. D. H. Kim et al., Single-cell transcriptome analysis reveals dynamic changes in lncRNA expression during reprogramming. Cell stem cell 16, 88-101 (2015).

**[0626]**      44. A. Parenti, M. A. Halbisen, K. Wang, K. Latham, A. Ralston, OSKM induce extraembryonic endoderm stem cells in parallel to induced pluripotent stem cells. Stem cell reports 6, 447-455 (2016).

**[0627]**      45. T. S. Mikkelsen et al., Dissecting direct reprogramming through integrative genomic analysis. Nature 454, 49 (2008).

**[0628]**      46. M. Stadtfeld, N. Maherali, M. Borkent, K. Hochedlinger, A reprogrammable mouse strain from gene-targeted embryonic stem cells. Nature methods 7, 53-55 (2010).

**[0629]**      47. Z. D. Smith, I. Nachman, A. Regev, A. Meissner, Dynamic single-cell imaging of direct reprogramming reveals an early specifying event. Nat Biotechnol 28, 521-526 (2010).

**[0630]**      48. J. Pei, N. V. Grishin, Unexpected diversity in Shisa-like proteins suggests the importance of their roles as transmembrane adaptors. Cellular signalling 24, 758-769 (2012).

**[0631]** 49. M. Meyyappan, H. Wong, C. Hull, K. T. Riabowol, Increased expression of cyclin D2 during multiple states of growth arrest in primary and established cells. Molecular and cellular biology 18, 3163-3172 (1998).

**[0632]** 50. J.-P. Coppe, P.-Y. Desprez, A. Krtolica, J. Campisi, The senescence-associated secretory phenotype: the dark side of tumor suppression. Annual Review of Pathological Mechanical Disease 5, 99-1 18 (2010).

**[0633]** 51. L. Mosteiro et al., Tissue damage and senescence provide critical signals for cellular reprogramming in vivo. Science 354, aaf4445 (2016).

**[0634]** 52. Q.-L. Ying et al., The ground state of embryonic stem cell self-renewal. Nature 453, 519 (2008).

**[0635]** 53. 1. Tirosh et al., Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma. Nature 539, 309-313 (2016).

**[0636]** 54. S. C. Andrews et al., Cdknlc (p57 Kip2) is the major regulator of embryonic growth within its imprinted domain on mouse distal chromosome 7. BMC Developmental Biology 7, 53 (2007).

**[0637]** 55. N. Barker et al., Identification of stem cells in small intestine and colon by marker gene Lgr5. Nature 449, 1003-1007 (2007).

**[0638]** 56. G. C. Elson et al., CLF associates with CLC to form a functional heteromeric ligand for the CNTF receptor complex. Nature neuroscience 3, 867 (2000).

**[0639]** 57. A. Fowden, C. Sibley, W. Reik, M. Constancia, Imprinted genes, placental development and fetal growth. Hormone Research in Paediatrics 65, 50-58 (2006).

**[0640]** 58. A. Ralston et al., Gata3 regulates trophoblast development downstream of Tead4 and in parallel to Cdx2. Development 137, 395-403 (2010).

**[0641]** 59. G. Burton, H.-W. Yung, T. Cindrova-Davies, D. Charnock-Jones, Placental endoplasmic reticulum stress and oxidative stress in the pathophysiology of unexplained intrauterine growth restriction and early onset preeclampsia. Placenta 30, 43-48 (2009).

**[0642]** 60. V. Pasque et al., X chromosome reactivation dynamics reveal stages of reprogramming to pluripotency. Cell 159, 1681-1697 (2014).

**[0643]** 61. K. Tomoda et al., Derivation conditions impact X-inactivation status in female human induced pluripotent stem cells. Cell stem cell 11, 91-99 (2012).

**[0644]** 62. Q. Bai et al., Dissecting the first transcriptional divergence during human embryonic development. Stem Cell Reviews and Reports 8, 150-162 (2012).

**[0645]** 63. A.-H. Monsoro-Burq, E. Wang, R. Harland, Msxl and Pax3 cooperate to mediate FGF8 and WNT signals during Xenopus neural crest induction. Developmental cell 8, 167-178 (2005).

**[0646]** 64. L. Pevny, M. Placzek, SOX genes and neural progenitor identity. Current opinion in neurobiology 15, 7-13 (2005).

**[0647]** 65. V. Y. Wang, H. Y. Zoghbi, Genetic regulation of cerebellar development. Nature reviews. Neuroscience 2, 484 (2001).

**[0648]** 66. Y. Liu, A. W. Helms, J. E. Johnson, Distinct activities of Msxl and Msx3 in dorsal neural tube development. Development 131, 1017-1028 (2004).

**[0649]** 67. M. Bergsland et al., Sequentially acting Sox transcription factors in neural lineage development. Genes Dev 25, 2453-2464 (201 1).

**[0650]** 68. K. Achim et al., The role of Tal2 and Tall in the differentiation of midbrain GABAergic neuron precursors. Biology open 2, 990-997 (2013).

**[0651]** 69. A. Domanskyi, H. Alter, M. A. Vogt, P. Gass, I. A. Vinnikov, Transcription factors Foxal and Foxa2 are required for adult dopamine neurons maintenance. Frontiers in cellular neuroscience 8, 275 (2014).

**[0652]** 70. K. Takebayashi-Suzuki, A. Kitayama, C. Terasaka-Iioka, N. Ueno, A. Suzuki, The forkhead transcription factor FoxBl regulates the dorsal-ventral and anterior-posterior patterning of the ectoderm during early Xenopus embryogenesis. Developmental biology 360, 11-29 (201 1).

**[0653]** 71. G. Hu et al., A genome-wide RNAi screen identifies a new transcriptional module required for self-renewal. Genes & development 23, 837-848 (2009).

**[0654]** 72. W.-Z. Li et al., Hesxl enhances pluripotency by working downstream of multiple pluripotency-associated signaling pathways. Biochemical and Biophysical Research Communications 464, 936-942 (2015).

**[0655]** 73. W. Shi et al., Regulation of the pluripotency marker Rex-1 by Nanog and Sox2. J Biol Chem 281, 23319-23325 (2006).

**[0656]**     74. A. Rajkovic, C. Yan, W. Yan, M. Klysik, M. M. Matzuk, Obox, a Family of Homeobox Genes Preferentially Expressed in Germ Cells. Genomics 79, 711-717 (2002).

**[0657]**     [SI) Villani C. Optimal Transport Old and New. Springer; 2008.

**[0658]**     [S2] Chizat L, Peyre' G, Schmitzer B, Vialard FX. Scaling Algorithms for Unbalanced Transport Problems. Mathematics of Computation. 2017;.

**[0659]**     [S3] Cuturi M. Sinkhorn Distances: Lightspeed Computation of Optimal Transportation Distances. In: Neural Information Processing Systems (NIPS); 2013. .

**[0660]**     [S4] https://support. 10xgenomics.com/single-cell-gene-expression/ software/pipelines/latest/installation.

**[0661]**     [S5] Coifman RR, Lafon S, Lee AB, Maggioni M, Nadler B, Warner F, et al. Geometric diffusions as a tool for harmonic analysis and structure definition of data: Diffusion maps. Proc Natl Acad Sci U S A. 2005;102:7426-7431.

**[0662]**     [S6] Haghverdi L, Buettner F, Theis FJ. Diffusion maps for high-dimensional single-cell analysis of differentiation data. Bioinformatics. 2015;31:2989-2998.

**[0663]**     [S7] Haghverdi L, Buettner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotyme robustly recon- structs lineage branching. bioRxiv. 2016;p. 041384.

**[0664]**     **[S8]** Angerer P, Haghverdi L, Bu˙ttner M, Theis FJ, Marr C, Buettner F. destiny: diffusion maps for large-scale single-cell data in R. Bioinformatics. 2015;32:1241-1243.

**[0665]**     [S9] Moignard V, Woodhouse S, Haghverdi L, Lilly AJ, Tanaka Y, Wilkinson AC, et al. Decoding the regulatory network of early blood development from single-cell gene expression measurements. Nature Biotechn. 2015;33:269-276.

**[0666]**     [S10] SettyM,TadmorMD,Reich-ZeligerS, AngelO, SalameTM, KathailP, et al. Wishbone identifies bifurcating developmental trajectories from single-cell data. Nature Biotechn. 2016;34:637-645.

**[0667]**     [SI 1] Satija R, Farrell JA, Gennert D, Schier AF, Regev A. Spatial reconstruction of single-cell gene expression data. Nature Biotechn. 2015;33:495-502.

**[0668]**     [S12] HeinzS,BennerC,SpannN,BertolinoE,LinYC,LasloP,etal.Simplecombinationso flineage- determining transcription factors prime cis-regulatory elements required for macrophage and B cell identities. Mol cell. 2010;38:576-589.

**[0669]** [S13] Bastian M, Heymann S, Jacomy M, et al. Gephi: an open source software for exploring and manipulating networks. Icwsm. 2009;8:361-362.

**[0670]** [S14] Jacomy M, Venturini T, Heymann S, Bastian M. ForceAtlas2, a continuous graph layout algo- rithm for handy network visualization designed for the Gephi software. PloS one. 2014;9:e98679.

**[0671]** [S15] Beygelzimer A, Kakadet S, Langford J, Arya S, Mount D, Li S, et al.. Package FNN;.

**[0672]** [SI 6] Zunder ER, Lujan E, Goltsev Y, Wernig M, Nolan GP. A continuous molecular roadmap to iPSC reprogramming through progression analysis of single-cell mass cytometry. Cell Stem Cell. 2015;16:323-337.

**[0673]** S17 Porpiglia E, Samusik N, Van Ho AT, Cosgrove BD, Mai T, Davis KL, et al. High-resolution myogenic lineage mapping by single-cell mass cytometry. Nature Cell Biol. 2017;19:558-567.

**[0674]** S18 Samusik N, Good Z, Spitzer MH, Davis KL, Nolan GP. Automated mapping of phenotype space with single-cell data. Nature methods. 2016;13:493-496.

**[0675]** S19 Blondel VD, Guillaume JL, Lambiotte R, Lefebvre E. Fast unfolding of communities in large networks. J Stat Mech Theor Exp. 2008;2008:P10008.

**[0676]** S20 Levine JH, Simonds EF, Bendall SC, Davis KL, El-ad DA, Tadmor MD, et al. Data-driven phenotypic dissection of AML reveals progenitor-like cells that correlate with prognosis. Cell. 2015;162:184-197.

**[0677]** S21 Shekhar K, Lapan SW, Whitney IE, Tran NM, Macosko EZ, Kowalczyk M, et al. Comprehensive classification of retinal bipolar neurons by single-cell transcriptomics. Cell. 2016;166:1308- 1323.

**[0678]** S22 Csardi G, Nepusz T. The igraph software package for complex network research. InterJournal, Complex Systems. 2006;1695:1-9.

**[0679]** S23 Friedman J, Hastie T, Tibshirani R. Sparse inverse covariance estimation with the graphical lasso. Biostatistics. 2008;9:432-441.

**[0680]** S24 Rosvall M, Bergstrom CT. Maps of random walks on complex networks reveal community struc- ture. Proc Natl Acad Sci U S A. 2008;105:1 118-1 123.

**[0681]** S25 Qiu X, Mao Q, Tang Y, Wang L, Chawla R, Pliner H, et al. Reversed graph embedding resolves complex single-cell developmental trajectories. bioRxiv. 2017;p. 110668.

**[0682]** S26 Qiu X, Hill A, Packer J, Lin D, Ma YA, Trapnell C. Single-cell mRNA quantification and differ- ential analysis with Census. Nature methods. 2017;14:309-315.

**[0683]** S27 Mao Q, Wang L, Goodison S, Sun Y. Dimensionality reduction via graph structure learning. In: Proceedings of the 21th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining. ACM; 2015. p. 765-774.

**[0684]** S28 Rashid S, Kotton DN, Bar-Joseph Z. TASIC: determining branching models from time series single cell data. Bioinformatics. 2017;p. btxl73.

**[0685]** S29 Lattin JE, Schroder K, Su AI, Walker JR, Zhang J, Wiltshire T, et al. Expression analysis of G Protein-Coupled Receptors in mouse macrophages. Immunome Res. 2008;4:5.

**[0686]** S30 Chen EY, Tan CM, Kou Y, Duan Q, Wang Z, Meirelles GV, et al. Enrichr: interactive and collaborative HTML5 gene list enrichment analysis tool. BMC Bioinformatics. 2013;14:128.

**[0687]** S31 Tirosh I, Venteicher AS, Hebert C, Escalante LE, Patel AP, Yizhak K, et al. Single-cell RNA-seq supports a developmental hierarchy in human oligodendroglioma. Nature. 2016;539:309-313.

**[0688]** S32 Li R, Liang J, Ni S, Zhou T, Qing X, Li H, et al. A mesenchymal-to-epithelial transition initiates and is required for the nuclear reprogramming of mouse fibroblasts. Cell stem cell. 2010;7:51-63. ]

**[0689]** S33 Whiteman EL, Fan S, Harder JL, Walton KD, Liu CJ, Soofi A, et al. Crumbs3 is essential for proper epithelial development and viability. Mol Cell Biol. 2014;34:43-56.

**[0690]** S34 Takaishi M, Tarutani M, Takeda J, Sano S. Mesenchymal to Epithelial Transition Induced by Re- programming Factors Attenuates the Malignancy of Cancer Cells. PloS one. 2016;l l:e0156904.

**[0691]** S35 Hewitt KJ, Agarwal R, Morin PJ. The claudin gene family: expression in normal and neoplastic tissues. BMC cancer. 2006;6:186.

**[0692]** S36 Coppe´ JP, Desprez PY, Krtolica A, Campisi J. The senescence-associated secretory phenotype: the dark side of tumor suppression. Annu Rev Pathol. 2010;5:99-1 18.

**[0693]** S37 da Fonseca ET, Mane ^nares ACF, Ambro ́sio CE, Miglino MA. Review point on neural stem cells and neurogenic areas of the central nervous system. Open J Anim Sci. 2013;3:242.

**[0694]** S38 Sakakibara Si, Nakamura Y, Satoh H, Okano H. Rna-binding protein Musashi2: developmentally regulated expression in neural precursor cells and subpopulations of neurons in mammalian CNS. J Neurosci. 2001;21:8091-8107.

**[0695]** S39 Gouti M, Briscoe J, Gavalas A. Anterior Hox genes interact with components of the neural crest specification network to induce neural crest fates. Stem cells. 201 1;29:858-870.

**[0696]** S40 Watanabe Y, Stanchina L, Lecerf L, Gacem N, Conidi A, Baral V, et al. Differentiation of Mouse Enteric Nervous System Progenitor Cells Is Controlled by Endothelin 3 and Requires Regulation of Ednrb by SOX10 and ZEB2. Gastroenterology. 2017; 152: 1139— 1150.

**[0697]** S41 Sansom SN, Griffiths DS, Faedo A, Kleinjan DJ, Ruan Y, Smith J, et al. The level of the tran- scription factor Pax6 is essential for controlling the balance between neural stem cell self-renewal and neurogenesis. PLoS Genetics. 2009;5:el00051 1.

**[0698]** S42 SKan L, Israsena N, Zhang Z, Hu M, Zhao LR, Jalali A, et al. Soxl acts through multiple inde- pendent pathways to promote neurogenesis. Dev Biol. 2004;269:580-594.

**[0699]** S43 Lazarov O, Mattson MP, Peterson DA, Pimplikar SW, van Praag H. When neurogenesis encoun- ters aging and disease. Trends Neurosci. 2010;33:569-579.

**[0700]** S44 Tibshirani R, Walther G, Hastie T. Estimating the number of clusters in a data set via the gap statistic. J R Stat Soc Series B Stat Methodol. 2001;63:41 1-423.

**[0701]** S45 Polo JM, Anderssen E, Walsh RM, Schwarz BA, Nefzger CM, Lim SM, et al. A molecular roadmap of reprogramming somatic cells into iPS cells. Cell. 2012;151(7): 1617— 1632.

**[0702]** S46 Mertins P, Przybylski D, Yosef N, Qiao J, Clauser K, Raychowdhury R, et al. An Integrative Framework Reveals Signaling-to-Transcription Events in Toll-like Receptor Signaling. Cell re- ports. 2017;19(13):2853-2866.

**[0703]** S47 ChoiJ, HuebnerAJ, ClementK, WalshRM, SavolA, LinK, etal. Prolonged Mekl/2suppression impairs the developmental potential of embryonic stem cells. Nature. 2017;548:219-223.

**[0704]** S48 Parenti A, Halbisen MA, Wang K, Latham K, Ralston A. OSKM induce extraembryonic endo- derm stem cells in parallel to induced pluripotent stem cells. Stem cell reports. 2016;6(4):447- 455.

**[0705]** [S49] Lin J, Khan M, Zapiec B, Mombaerts P. Efficient derivation of extraembryonic endoderm stem cell lines from mouse postimplantation embryos. Scientific reports. 2016;6.

**[0706]** [S50] Edgar R, Mazor Y, Rinon A, Blumenthal J, Golan Y, Buzhor E, et al. LifeMap Discovery?: the embryonic development, stem cells, and regenerative medicine research portal. PloS one. 2013;8(7):e66629.

\*\*\*

**[0707]** Various modifications and variations of the described methods, pharmaceutical compositions, and kits of the invention will be apparent to those skilled in the art without departing from the scope and spirit of the invention. Although the invention has been described in connection with specific embodiments, it will be understood that it is capable of further modifications and that the invention as claimed should not be unduly limited to such specific embodiments. Indeed, various modifications of the described modes for carrying out the invention that are obvious to those skilled in the art are intended to be within the scope of the invention. This application is intended to cover any variations, uses, or adaptations of the invention following, in general, the principles of the invention and including such departures from the present disclosure come within known customary practice within the art to which the invention pertains and may be applied to the essential features herein before set forth.

# CLAIMS

What is claimed is:

1.      A method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Obox6 into a target cell to produce an induced pluripotent stem cell.

2.      The method of claim 1, further comprising introducing into the target cell at least one nucleic acid encoding a reprogramming factor selected from the group consisting of: Gdf9, Oct3/4, Sox2, Soxl, Sox3, Soxl5, Soxl7, Klf4, Klf2, c-Myc, N-Myc, L-Myc, Nanog, Lin28, Fbxl5, ERas, ECAT15-2, Tell, beta-catenin, Lin28b, Sall l, Sal 14, Esrrb, Nr5a2, Tbx3, and Glisl.

3.      The method of claim 1, further comprising introducing into the target cell at least one nucleic acid encoding a reprogramming factor selected from the group consisting of: Oct4, Klf4, Sox2 and Myc.

4.       The method of claim 1, wherein the nucleic acid encoding Obox6 is provided in a recombinant vector.

5.      The method of claim 4, wherein the vector is a lentivirus vector.

6.      The method of claim 2, where the nucleic acid encoding the reprogramming factor is provided in a recombinant vector.

7.      The method of claim 1, further comprising a step of culturing the cells in reprogramming medium.

8.      The method of claim 1, further comprising a step of culturing the cells in the presence of serum.


9.      The method of claim 1, further comprising a step of culturing the cells in the absence of serum.


10.     The method of claim 1, wherein the induced pluripotent stem cell expresses at least one of a surface marker selected from the group consisting of: Oct4, SOX2, KLf4, c-MYC, LIN28, Nanog, Glisl , TRA-160/TRA-1-81/TRA-2-54,  SSEA1, SSEA4, Sal4, and Esrbbl.


11.     The method of claim 1, wherein the target cell is a mammalian cell.


12.     The method of claim 1, wherein the target cell is a human cell or a murine cell.


13.     The method of claim 1, wherein the target cell is a mouse embryonic fibroblast.


14.     The method of claim 1, wherein the target cell is selected from the group consisting of: fibroblasts, B cells, T cells, dendritic cells, keratinocytes, adipose cells, epithelial cells, epidermal cells, chondrocytes, cumulus cells, neural cells, glial cells, astrocytes, cardiac cells, esophageal cells, muscle cells, melanocytes, hematopoietic cells, pancreatic cells, hepatocytes, macrophages, monocytes, mononuclear cells, and gastric cells, including gastric epithelial cells.


15.     A method of producing an induced pluripotent stem cell comprising introducing at least one of Obox6, Spic, Zfp42, Sox2, Mybl2, Msc, Nanog, Hesxl  and Esrrb into a target cell to produce an induced pluripotent stem cell.

16.     A method of producing an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell to produce an induced pluripotent stem cell.

17.     A method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

18.     A method of increasing the efficiency of production of an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell to produce an induced pluripotent stem cell.

19.     An isolated induced pluripotential stem cell produced by the method of claim 1, 15, or 16.

20.     A method of treating a subject with a disease comprising administering to the subject a cell produced by differentiation of the induced pluripotent stem cell produced by the method of claim 1, 15, or 16.

21.     A composition for producing an induced pluripotent stem cell comprising Obox6 in combination with reprogramming medium.

22.     A composition for producing an induced pluripotent stem cell comprising one or more of the factors identified in or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6 in combination with reprogramming medium.

23.     Use of Obox6 for production of an induced pluripotent stem cell.

24.     Use of a factor identified in or one or more of the factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6 for production of an induced pluripotent stem cell.

25.     A method of increasing the efficiency of reprogramming a cell comprising introducing Obox6 into a target cell to produce an induced pluripotent stem cell.

26.     A method of increasing the efficiency of reprogramming a cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5 and Table 6, into a target cell to produce an induced pluripotent stem cell.

27. A computer-implemented method for mapping developmental trajectories of cells, comprising:

generating, using one or more computing devices, optimal transport maps for a set of cells from single cell sequencing data obtained over a defined time course;

determining, using one or more computing devices, cell regulatory models, and optionally identifying local biomarker enrichment, based on at least the generated optimal transport maps;

defining, using the one or more computing devices, gene modules; and

generating, using the one or more computing devices, a visualization of a developmental landscape of the set of cells.

28. The method of claim 27, wherein determining cell regulatory models comprise sampling pairs of cells at a first time and a second time point according to transport probabilities.

29. The method of claim 28, further comprising using the expression levels of transcription factors at the earlier time point to predict non-transcription factor expression at the second time point.

30. The method of claim 27, wherein identifying local biomarker enrichment comprises identifying transcription factors enriched in cells having a defined percentage of descendants in a target cell population.

31. The method of claim 30, wherein the defined percentage is at least 50% of mass.

32. The method of claim 27, wherein defining gene modules comprises partitioning genes based on correlated gene expression across cells and clusters.

33. The method of claim 32, wherein partitioning comprises partitioning cells based on graph clustering.

34. The method of claim 33, wherein graph clustering further comprises dimensionality reduction using diffusion maps.

35. The method of claim 27, wherein the visualization of the developmental landscape comprises high-dimensional gene expression data in two dimensions.

36. The method of claim 33, wherein the visualization is generated using force-directed layout embedding (FLE).

37. The method of claim 27, wherein the visualization provides one or more cell types, cell ancestors, cell descendants, cell trajectories, gene modules, and cell clusters from the single cell sequencing data.

38. A computer program product, comprising:

a non-transitory computer-executable storage device having computer-readable program instructions embodied thereon that when executed by a computer cause the computer to execute the methods of anyone of claims 27 to 37.

39. A system comprising:

a storage device; and

a processor communicatively coupled to the storage device, wherein the processor executes application code instructions that are stored in the storage device and that cause the system to executed the methods of any one of claims 27 to 37.


40. A method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Gdf9 into a target cell to produce an induced pluripotent stem cell.

1/48

100



FIG. 1

**200**

**205**

Generate optimal transport maps for a set of cells from single cell sequencing data obtained over a defined time course

**210**

Determine cell regulatory models and optionally identifying local biomarker enrichment based on the optimal transport maps

**215**

Define gene modules

**220**

Generate a visualization of a developmental landscape of of the set of cells

# FIG. 2

FIG. 3

**FIG. 4**

FIG. 5A



FIG. 5B



FIG. 5C



FIG. 5D



FIG. 5E



FIG. 5F



FIG. 5G

FIG. 6C



FIG. 6B



FIG. 6A

## FIG. 7A

## FIG. 7B

**2i condition**

Mature iPSCs

pre-iPSCs

Alternative Fates

Placental-like cells

Horn of Transformation

Valley of Stress

Phase-2(2i)

Phase-1(Dox)

Day 4

Day 2

MEFs

D0
D2
D4
D6
D8
D9
D10
D11
D12
D16
iPSCs
BC

## FIG. 7C

**Serum condition**

Mature iPSCs

Placental-like cells

pre-iPSCs

Neural-like cells

Phase-2(serum)

Phase-1(Dox)

## FIG. 7D

**Cell Clusters**

33 32 29
31 30
28
27 26 23 24 25
21 22
20 19 10 11 13 15 16 17
9 12 14
18 8
7
5 3 2
6 4
1

**MEF identity** 1.4 / -0.4

**Pluripotency** 1.8 / -0.7

**OKSM** 5.3 / 0

**Proliferation** 1.5 / -0.6

**ER stress** 0.8 / -0.5

**Epithelial identity** 3.7 / -0.3

**ECM rearrangement** 1.1 / -0.3

**Apoptosis** 0.9 / -0.5

**Senescence (SASP)** 1.3 / -0.4

**Neural identity** 1.9 / -0.2

**Placental identity** 1.8 / -0.3

**X reactivation** 0.6 / -0.4

## FIG. 7E

**Shisa8** 4.7 / 0

**Cdkn2a** 7.1 / 0

**Cdkn1c** 9 / 0

**Cntfr** 3.9 / 0

**H19** 8.3 / 0

**Nanog** 4.8 / 0

**Sox2** 3.8 / 0

**Obox6** 3.6 / 0

**Zfp42** 4.1 / 0

## FIG. 7F

FIG. 8A

FIG. 8B

FIG. 8C

FIG. 8D

FIG. 8E

FIG. 8F

FIG. 9A

10/48



FIG. 9B

FIG. 9C

FIG. 9D

FIG. 10A

FIG. 10B

FIG. 10C

FIG. 11A

FIG. 11B

FIG. 11C

FIG. 11D

Force directed layout (Diffusion Maps)



**FIG. 12A**

tSNE (Principal Components)



**FIG. 12B**

tSNE (Diffusion Maps)



**FIG. 12C**

14/48



**FIG. 13**

FIG. 14B

FIG. 14A

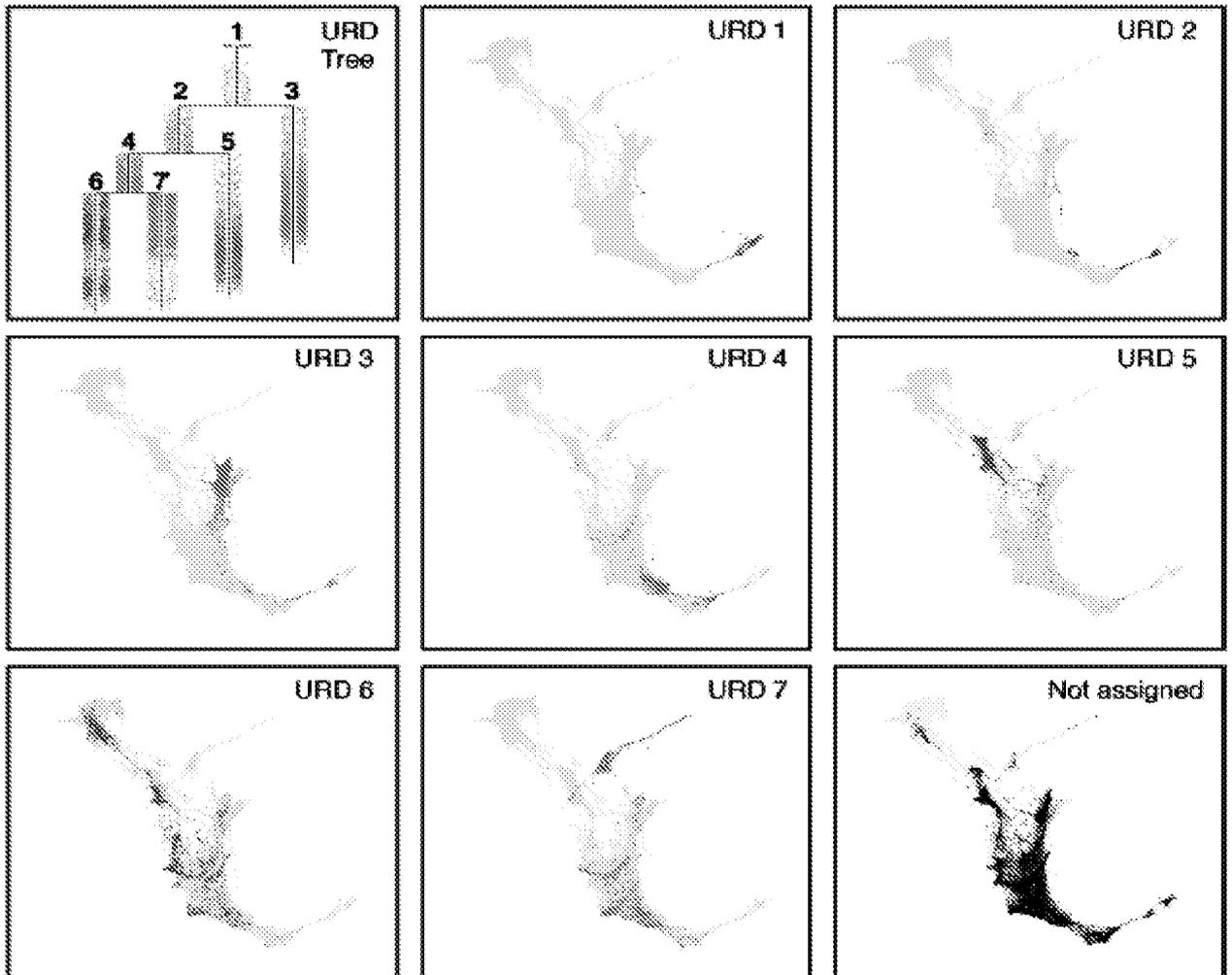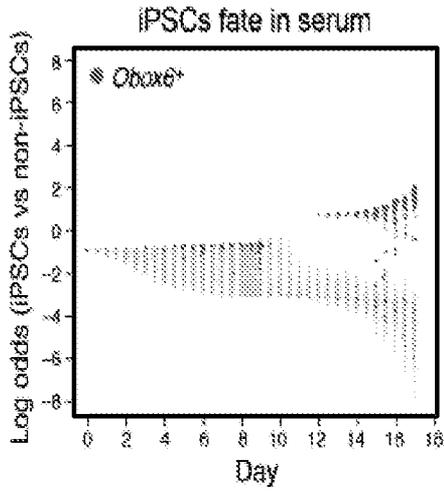FIG. 15

FIG. 16A

FIG. 16B

FIG. 16C

FIG. 16E



FIG. 16D

19/48

Cluster descendants,Serum



**FIG. 17A**

Cluster ancestors, 2i

Cluster ancestors, serum



**FIG. 17B**

**FIG. 17C**

FIG. 18A



FIG. 18B

21/48



FIG. 18F

FIG. 18E

FIG. 18D

FIG. 18C

FIG. 19A

Placenta-like cells, module B

FIG. 19B

FIG. 19C

**FIG. 19D**

Succesful reprogramming, module B

**FIG. 19E**

FIG. 19F

28/48

**2nd MEFs**



FIG. 20A

**Primary MEFs**



FIG. 20B

**Primary MEFs**



FIG. 20C

FIG. 21A

FIG. 21B

FIG. 21C

FIG. 21E



FIG. 21D

FIG. 22A

FIG. 22B



FIG. 22C

**FIG. 23A**



Time course

**FIG. 23B**



Descendants

**FIG. 23C**



Ancestors

**FIG. 23D**



Shared ancestry

**FIG. 23E**



Gene signature



Transcription factor



**FIG. 23F**

**FIG. 24A**



**FIG. 24B**                                          **FIG. 24C**



**FIG. 24D**



**FIG. 24E**              **FIG. 24F**              **FIG. 24G**

**FIG. 24H**

ECM rearrangement        Apoptosis        Senescence (SASP)        *Cdkn2a*



**FIG. 25A**

Stromal ancestors



**FIG. 25B**

Signature trends for stromal ancestors

TF trends for stromal ancestors



**FIG. 25C**

MET



**FIG. 25D**

Signature trends for MET ancestors

TF trends for MET ancestors



**FIG. 25E**



**FIG. 25F**

MET / stromal



**FIG. 25G**



**FIG. 25H**

FIG. 26A

FIG. 26B

FIG. 26C

FIG. 26D

FIG. 26E

FIG. 26F

FIG. 27A

FIG. 27B

FIG. 27C

FIG. 27D

FIG. 27E

FIG. 27F

FIG. 27G

FIG. 28A

FIG. 28B



FIG. 28C

FIG. 28D

FIG. 28E



FIG. 28F

FIG. 28G

FIG. 28H



FIG. 28I

FIG. 28J

FIG. 28K

FIG. 29A

FIG. 29B

FIG. 29C

FIG. 29D

FIG. 30A



FIG. 30B



FIG. 30C



FIG. 30D



FIG. 30E



FIG. 30F



FIG. 30G

FIG. 31A

FIG. 31B

FIG. 31C

FIG. 31D

FIG. 31E

FIG. 31F

FIG. 32A



FIG. 32B

**42/48**



FIG. 32C

FIG. 33A



FIG. 33B



FIG. 33C



FIG. 33D



FIG. 33E

FIG. 34A

FIG. 34B

FIG. 34C

FIG. 34D

FIG. 34E

FIG. 35A



FIG. 35B

FIG. 36A

FIG. 36B

FIG. 36C

FIG. 36D

FIG. 36E

FIG. 36F

FIG. 37

FIG. 38

# INTERNATIONAL SEARCH REPORT

**A. CLASSIFICATION OF SUBJECT MATTER**

**IPC - C 12 N 5/07, 15/63, 15/85 (201 8.01 )**
**CPC -**
**C 12 N 5/0606, 15/63, 15/85**

According to International Patent Classification. (IPC) or to both national classification and IPC

**B. FIELDS SEARCHED**

Minimum documentation searched (classification system followed by classification symbols)
See Search History document

Documentation searched other than minimum documentation to the extent that such documents are included in the fields searched
See Search History document

Electronic data base consulted during the international search (name of data base and, where practicable, search terms used)
See Search History document

**C. DOCUMENTS CONSIDERED TO BE RELEVANT**

| Category* | Citation of document, with indication, where appropriate, of the relevant passages | Relevant to claim No. |
|---|---|---|
| X<br>--<br>Y | (KIM, HM et al.) 'Obox$_4$ regulates the expression of histone family genes and promotes differentiation of mouse embryonic stem cells'; 05 February 2010, FEBS Letters; Volume 584, Issue 3, pages 605-61 1; abstract; page 606, first column, first, second and third paragraphs; page 607, first column, last paragraph; page 607, second column first paragraph; page 608, figure 2 | 1-3, 10-12, 15-18, 19/1 , 19/15-16, 23-24<br>-------------------------------<br>4-9, 13-14, 21-22 |
| Y | US 2014/028751 1 A1 (KO, MSH) 25 September 2014; paragraph [01 13]; claims 1, 5, 7-8 | 4-6 |
| Y | US 2010/0330677 A1 (SMITH, AG) 30 December 2010; paragraphs [0064]-[0066], [0090], [0178] | 7-9, 21-22 |
| Y | US 2013/0295579 A1 (XIE, X et al.) 07 November 2013; paragraphs [0007]-[0008], [0012] | 13-14 |

☐ Further documents are listed in the continuation of Box C.     ☐ See patent family annex.

| | Special categories of cited documents: | "T" | later document published after the international filing date or priority date and not in conflict with the application but cited to understand the principle or theory underlying the invention |
|---|---|---|---|
| "A" | document defining the general state of the art which is not considered to be of particular relevance | | |
| "E" | earlier application or patent but published on or after the international filing date | "X" | document of particular relevance; the claimed invention cannot be considered novel or cannot be considered to involve an inventive step when the document is taken alone |
| "L" | document which may throw doubts on priority claim(s) or which is cited to establish the publication date of another citation or other special reason (as specified) | "Y" | document of particular relevance; the claimed invention cannot be considered to involve an inventive step when the document is combined with one or more other such documents, such combination being obvious to a person skilled in the art |
| "O" | document referring to an oral disclosure, use, exhibition or other means | | |
| "P" | document published prior to the international filing date but later than the priority date claimed | "&" | documen , member of the same patent family |

| Date of the actual completion of the international search | Date of mailing of the international search report |
|---|---|
| 29 October 2018 (29.10.2018) | **1 2 FEB 2019** |

| Name and mailing address of the ISA/ | Authorized officer |
|---|---|
| Mail Stop PCT, Attn: ISA/US, Commissioner for Patents<br>P.O. Box 1450, Alexandria, Virginia 22313-1450<br>Facsimile No. 571-273-8300 | Shane Thomas<br><br>PCT Helpdesk: 571-272-4300<br>PCT OSP: 571-272-7774 |

Form PCT/ISA/2 10 (second sheet) (January 201 5)

**Box No. 1]   Observations where certain claims were found unsearchable (Continuation of item 2 of first sheet)**

This international search report has not been established in respect of certain claims under Article 17(2)(a) for the following reasons:

1.  ▫   Claims Nos.:
        because they relate to subject matter not required to be searched by this Authority,  namely:

2.  ☐   Claims Nos,;
        because they relate to parts of the international application that do not comply with the prescribed requirements to such an
        extent that no meaningful  international  search can be carried out, specifically:

3.  ▫   Claims Nos.:
        because they are dependent claims and are not drafted in accordance with the second and third sentences of Rule 6.4(a).

**Box No. Ill    Observations where unity of invention is lacking (Continuation of item 3 of first sheet)**

This International Searching Authority found multiple inventions in this international application, as follows:

-""-Continued Within the Next Supplemental Box-*" -

1.  ☐   As all required additional search fees were timely paid by the applicant, this international search report covers all searchable
        claims.

2.  ☐   As all searchable claims could be searched without effort justifying additional fees, this Authority did not invite payment of
        additional fees.

3.  ☐   As only some of the required additional search fees were timely paid by the applicant, this international search report covers
        only those claims for which fees were paid, specifically claims Nos.:

4.  ☒   No required additional search fees were timely paid by the applicant.  Consequently,  this international search report is
        restricted to the invention first mentioned in the claims; it is covered by claims Nos.:
        **1-18, 19/1, 19/15-16, 21-24**

**Remark on Protest**   ☐   The additional search fees were accompanied by the applicant's protest and, where applicable, the
                             payment of a protest fee.
                         ☐   The additional search fees were accompanied by the applicant's protest but the applicable protest
                             fee was not paid within the time limit specified in the invitation.
                         ☐   No protest accompanied the payment of additional search fees.

Form PCT/ISA/2 10 (continuation of first sheet (2)) (January 20 IS)

-'''-Continued from Box No. III Observations where unity of invention is lacking -*** -

This application contains the following inventions or groups of inventions which are not so linked as to form a single general inventive concept under PCT Rule 13.1. In order for all inventions to be examined, the appropriate additional examination fees must be paid.
Groups I+, Claims 1-19 (in-part) and 21-24 (in-part) are directed towards a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Obox6 into a target cell (first exemplary transcription factor).
The method and system will be searched to the extent they encompass a transcription factor of Obox6 (first exemplary transcription factor). Applicant is invited to elect additional transcription factor(s), with specified transcription factor(s) for each, to be searched. Additional transcription factor(s) will be searched upon the payment of additional fees. It is believed that claims 1-19 (in-part) and 21-24 (in-part) encompass this first named invention and thus these claims will be searched without fee to the extent that they encompass a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Obox6 into a target cell (first exemplary transcription factor). Applicants must specify the claims that encompass any additionally elected transcription factor(s). Applicants must further indicate, if applicable, the claims which encompass the first named invention, if different than what was indicated above for this group. Failure to clearly identify how any paid additional invention fees are to be applied to the "+" group(s) will result in only the first claimed invention to be searched/examined. An exemplary election would be a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Rhox2a into a target cell (first exemplary elected transcription factor).
Groups I+ share the technical features including a method of producing an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell.
However, these shared technical features are previously disclosed by US 2014/028751 1 A1 (KO).
Ko discloses a method of producing an induced pluripotent stem cell (method of producing induced stem cells; paragraph [0006]) comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell (method comprises introducing Patl2 (transcription factor identified in Table 4) into somatic (target) cells; paragraphs [0006]-[0007]).
Since none of the special technical features of the Groups I+ inventions is found in more than one of the inventions, and since all of the shared technical features are previously disclosed by the Ko reference, unity of invention is lacking.
Group II, Claim 20 is directed towards a method of treating a subject with a disease.
Group III, Claims 25-26 are directed toward a method of increasing the efficiency of reprogramming a cell.
Group IV, Claims 27-39 are directed toward a computer-implemented method for mapping developmental trajectories of cells.
Group V, Claim 40 is directed towards a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Gdf9.
The inventions listed as Groups I+ and II-V do not relate to a single general inventive concept under PCT Rule 13.1 because, under PCT Rule 13.2, they lack the same or corresponding special technical features for the following reasons: the special technical features of Groups I+ include a method of producing an induced pluripotent stem cell comprising introducing at least one of the transcription factors identified in Table 2, Table 3, Table 4, Table 5, and Table 6, into a target cell, which is not present in Groups II-V; the special technical features of Group II include a method of treating a subject with a disease, which is not present in Groups I+ and III-V; the special technical features of Group III include a method of increasing the efficiency of reprogramming a cell, which is not present in Groups I+, II and IV-V; the special technical features of Group IV include a computer-implemented method for mapping developmental trajectories of cells, which is not present in Groups I+, II-III, and V; and the special technical features of Group IV include a method of producing an induced pluripotent stem cell comprising introducing a nucleic acid encoding Gdf9, which is not present in Groups I+ and II-III.
The common technical features of Groups I+, II-III and V are a method of introducing a transcription factor into a target cell.
These common technical features are disclosed by Ko. Ko discloses a method of introducing a transcription factor into a target cell (method comprising introducing Patl2 (transcription factor identified in Table 4) into somatic (target) cells; paragraphs [0006]-[0007]).
Since the common technical features are previously disclosed by Ko, these common features are not special and so Groups I+ and II-IV lack unity.
No technical features are shared between Groups I+ and IV, accordingly, these groups lack unity a priori.
No technical features are shared between Groups II and IV, accordingly, these groups lack unity a priori.
No technical features are shared between Groups III and IV, accordingly, these groups lack unity a priori. No technical features are shared between Groups IV and V, accordingly, these groups lack unity a priori.