Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq

Barbara Treutlein^{1,2*}, Qian Yi Lee^{1,3,4*}, J. Gray Camp⁵, Moritz Mall^{3,4}, Winston Koh¹, Seyed Ali Mohammad Shariati⁶, Sopheak Sim³, Norma F. Neff¹, Jan M. Skotheim⁶, Marius Wernig^{3,4}§ & Stephen R. Quake^{1,7,8}§

Direct lineage reprogramming represents a remarkable conversion of cellular and transcriptome states¹⁻³. However, the intermediate stages through which individual cells progress during reprogramming are largely undefined. Here we use singlecell RNA sequencing⁴⁻⁷ at multiple time points to dissect direct reprogramming from mouse embryonic fibroblasts to induced neuronal cells. By deconstructing heterogeneity at each time point and ordering cells by transcriptome similarity, we find that the molecular reprogramming path is remarkably continuous. Overexpression of the proneural pioneer factor Ascl1 results in a well-defined initialization, causing cells to exit the cell cycle and re-focus gene expression through distinct neural transcription factors. The initial transcriptional response is relatively homogeneous among fibroblasts, suggesting that the early steps are not limiting for productive reprogramming. Instead, the later emergence of a competing myogenic program and variable transgene dynamics over time appear to be the major efficiency limits of direct reprogramming. Moreover, a transcriptional state, distinct from donor and target cell programs, is transiently induced in cells undergoing productive reprogramming. Our data provide a high-resolution approach for understanding transcriptome states during lineage differentiation.

Direct lineage reprogramming bypasses an induced pluripotent stage to directly convert somatic cell types. Using the three transcription factors Ascl1, Brn2 and Myt11 (BAM), mouse embryonic fibroblasts (MEFs) can be directly reprogrammed to induced neuronal (iN) cells within 2 to 3 weeks at an efficiency of up to 20%⁸. Several groups have further developed this conversion using transcription factor combinations that almost always contain Ascl1 (refs 9-12). Recently, one of our groups showed that Ascl1 is an 'on target' pioneer factor initiating the reprogramming process¹³, and inducing conversion of MEFs into functional iN cells alone, albeit at a much lower efficiency compared to BAM¹⁴. These findings raised the question whether and when a heterogeneous cellular response to the reprogramming factors occurs during reprogramming and which mechanisms might cause failure of reprogramming. We hypothesized that single-cell RNA sequencing (RNA-seq) could be used as a high-resolution approach to reconstruct the reprogramming path of MEFs to iN cells and uncover mechanisms limiting reprogramming efficiencies^{4,15,16}.

In order to understand transcriptional states during direct conversion between somatic fates, we measured 405 single-cell transcriptomes (Supplementary Data 1) at multiple time points during iN cell reprogramming (Fig. 1a and Extended Data Fig. 1a). We first explored how individual cells respond to Ascl1 overexpression during the initial phase of reprogramming. We analysed day 0 MEFs and day 2 cells induced with Ascl1 only (hereafter referred to as Ascl1-only cells) using PCA and identified three distinct clusters (A, B, C), which correlated with the level of Ascl1 expression (Fig. 1b-e). Cluster A consisted of all control d0 MEFs and a small fraction of day 2 cells (\sim 12%) which showed no detectable Ascl1 expression, suggesting these day 2 cells were not infected with the Ascl1 virus. This is consistent with typical Ascl1 infection efficiencies of about 80-90%. We found that the day 0 MEFs were surprisingly homogeneous, with much of the variance due to cell cycle (Extended Data Fig. 1b-g, Supplementary Data 3, Supplementary Information). Cluster C was characterized by high expression of Ascl1, Ascl1-target genes (Zfp238, Hes6, Atoh8 and so on) and genes involved in neuron remodelling, as well as the downregulation of genes involved in cell cycle and mitosis (Fig. 1c, e, f and Supplementary Data 2). Cluster B cells represent an intermediate population that expressed Ascl1 at a low level, and were characterized by a weaker upregulation of Ascl1-target genes and less efficient downregulation of cell cycle genes compared to cluster C cells. This suggests that an Ascl1 expression threshold is required to productively initiate the reprogramming process. In addition, we found that forced Ascl1 expression resulted in less intracellular transcriptome variance, a lower number of expressed genes (Fig. 1d) and a lower total number of transcripts per single cell (Extended Data Fig. 2a, b). Notably, the distribution of average expression levels per gene was similar for all experiments independent of Ascl1 overexpression (Extended Data Fig. 2c). We observed that the upregulation of neuronal targets and downregulation of cell cycle genes in response to Ascl1 expression are uniform, indicating that the initial transcriptional response to Ascl1 is relatively homogenous among all cells (Fig. 1e). This suggests that most fibroblasts are initially competent to reprogram and later events must be responsible for the moderate reprogramming efficiency of about 20%.

To explore the effect of transgene copy number variation on the heterogeneity of the early response, we analysed single-cell transcriptomes of an additional 47 cells induced with Ascl1 for two days from secondary MEFs derived via blastocyst injection from a clonal, Ascl1-inducible embryonic stem cell line. As expected, the induction efficiency of *Ascl1* was 100% since the secondary MEFs are genetically identical and all cells carry the transgene in the same genomic location (Fig. 1g). Nevertheless, these clonal MEFs had similar transcriptional responses and heterogeneity as primary infected MEFs at the day 2 time point, as well as comparable reprogramming efficiencies and maturation (Extended Data Fig. 3a). Finally, we compared the early response in our Ascl1-only single-cell RNA-seq data with our previously reported bulk RNA-seq data of Ascl1-only and BAM-mediated reprogramming¹³ (Extended Data Fig. 3b). We found similar downregulation of MEF-related genes and upregulation of pro-neural marker genes in both

*These authors contributed equally to this work.

§These authors jointly supervised this work.

¹Department of Bioengineering, Stanford University, Stanford, California 94305, USA. ²Department of Evolutionary Genetics, Max Planck Institute for Evolutionary Anthropology, Leipzig 04103, Germany. ³Institute for Stem Cell Biology and Regenerative Medicine, Stanford University School of Medicine, Stanford, California 94305, USA. ⁴Department of Pathology, Stanford University School of Medicine, Stanford, California 94305, USA. ⁵Department of Developmental Biology, Stanford University School of Medicine, Stanford, California 94305, USA. ⁶Department of Biology, Stanford University, Stanford University, Stanford, California 94305, USA. ⁷Howard Hughes Medical Institute, Stanford, California 94305, USA. ⁸Department of Applied Physics, Stanford University, Stanford, California 94305, USA. ⁷Howard Hughes Medical Institute, Stanford, California 94305, USA. ⁸Department of Applied Physics, Stanford University, Stanford, California 94305, USA.



Figure 1 | Ascl1 overexpression elicits a homogeneous early response and initiates expression of neuronal genes. a, Mouse embryonic fibroblasts stably integrated with neuronal reporter Tau-eGFP8 were directly transformed to neuronal cells through overexpression of a single (Ascl1), or three factors (Brn2, Ascl1, Myt1l; BAM) as described⁸. Cells were sampled using single-cell RNA-seq at day 0 without infection (d0, 73 cells), day 2 (d2, 81 cells Ascl1-infected and 47 cells clonal), day 5 (d5, 55 cells, eGFP⁺ and eGFP⁻ cells), day 20 (d20, 33 cells, eGFP⁺ cells), and day 22 (d22, 73 cells, eGFP⁺ cells) post-induction with Ascl1. As a comparison, cells reprogrammed using all three BAM factors were analysed at 22 days (d22, 43 cells, eGFP⁺ cells). b, c, PCA of single-cell transcriptomes from day 0 MEFs (circle, 73 cells) and day 2 Ascl1-induced cells (square, 81 cells) shows reduced intercellular variation at day 2. Points are coloured based on hierarchical clustering shown in e (b), or Ascl1 expression (c). d, Left, distribution of transcriptome variance within single cells grouped by cluster assignment of **b** and **e** shows that Ascl1 expression reduces the intracellular transcriptome variance. Right, distribution of total number of genes expressed by single cells grouped by cluster assignment shows that Ascl1 overexpression reduces the range of gene expression. e, Hierarchical clustering of day 0 and day 2 cells (rows) using the top 50 genes (columns) correlating positively (genes I) and negatively (genes II) with PC1. Cells are clustered into three clusters (left sidebar): A (83 cells, MEFs), B (20 cells, intermediates), C (51 cells, day 2 induced cells). f, Top gene ontology enrichments of genes I and II (d) are shown with Bonferroni-corrected P values. BP, biological process; CC, cellular component; reg. exc. memb. pot., regulation of excitatory postsynaptic membrane potential. g, Distribution of PC1 loadings are shown for day 2 cells carrying variable numbers of Ascl1 transgene copies (dark green, Ascl1-infected) or carrying the same Ascl1 copy number and genomic location (yellow, clonal). PC1 effectively separates un-induced MEFs (cluster A) from induced cells highly expressing Ascl1-target genes (cluster C) and both, Ascl1-infected and clonal cells, productively initiate reprogramming. The induction efficiency is higher for clonally induced MEFs, however even in the clonal population Ascl1 induction is variable.

Ascl1- and BAM-mediated reprogramming. These data suggest that the overexpression of Ascl1 focuses the transcriptome and directs the expression of target genes.

We next analysed the transcriptomes of reprogramming cells on day 5. At this time point, the first robust Tau–eGFP signal can be detected in successfully reprogramming cells and we therefore purified

40 Tau-eGFP⁺ and 15 Tau-eGFP⁻ cells for transcriptome analysis by fluorescence-activated cell sorting. We found that Tau-eGFP⁻ cells lacked expression of neuronal Ascl1-target genes (genes B), and maintained expression of fibroblast-associated genes (genes A and C; Fig. 2a, b, Extended Data Fig. 4a, Supplementary Data 4). In addition, we found a positive correlation ($R^2 = 0.49$) between Ascl1 expression and Tau-eGFP intensities (Extended Data Fig. 4b, Fig. 2a, b). Quantitative real-time (qRT)-PCR and western blot analysis of Ascl1 expression on day 5 to day 12 Tau-eGFP-sorted cells validated a significant decrease in Ascl1 expression in Tau-eGFP⁻ cells compared to Tau-eGFP⁺ cells (Fig. 2c, Supplementary Data 5). Thus, Ascl1 expression is correlated to Tau-eGFP levels and expression of neuronal genes at day 5. This raises the hypothesis that Ascl1 is silenced in cells that fail to reprogram. Alternatively, cells with low or no Ascl1 expression at day 5 and day 22 might have never highly expressed Ascl1. To distinguish between these two mechanisms, we used live cell microscopy to track cells over a time course from 3-6 days after Ascl1 induction using an eGFP-Ascl1 fusion construct (Fig. 2d, Extended Data Fig. 5). We immunostained the cells at day 6 using Tuj1 antibodies recognizing the neuronal β 3-tubulin (Tubb3) to identify cells that differentiated towards neuronal fate. We found that transgenic Ascl1 protein levels varied substantially over time and, on average, continued to increase over time in Tuj1⁺ cells, but decreased or plateaued in Tuj1⁻ cells, leading to a significant difference in Ascl1 expression within six days of Ascl1 induction (Fig. 2e, Extended Data Fig. 4c). This time-lapse analysis demonstrated that Ascl1 is silenced in many cells that fail to reprogram.

We next analysed the maturation events occurring during late reprogramming stages. We performed principal component analysis (PCA) on the single-cell transcriptomes of all reprogramming stages analysed, including day 22 cells reprogrammed with Ascl1 alone or with all three BAM factors (Extended Data Fig. 6a). PC1 separated MEFs and early time points (day 2, day 5) from most of the day 22 cells. Surprisingly, PC2 separated most day 22 BAM cells from day 22 Ascl1only cells despite robust Tau-eGFP expression in both groups. We used t-distributed stochastic neighbour embedding (tSNE) to organize all day 22 cells into transcriptionally distinct clusters, and identified differentially expressed genes marking each cluster (Fig. 3a). We identified 3 clusters, which contained cells expressing neuron (Syp), fibroblast (Eln), or myocyte (Tnnc2) marker genes, respectively (Fig. 3b). Consistent with this marker gene expression, cells in each cluster had a maximum correlation with bulk RNA-seq data from purified neurons, embryonic fibroblasts, or myocytes (Fig. 3c). Neuron- and myocyte-like cells expressed a clear signature of each cell type (Fig. 3d). Although we observed cells with complex neuronal morphologies in the Ascl1-only reprogramming experiments as we had reported previously¹⁴ (Fig. 3e), their frequency was too low to be captured in the single-cell RNA-seq experiments. All of the day 22 Ascl1-only cells, and 33% of BAM cells had a highest correlation with myocytes or fibroblasts.

We applied an analytical technique based on quadratic programming to quantify fate conversion and to predict when during reprogramming the alternative muscle program emerges (Extended Data Fig. 6b). This method allowed us to decompose each single cell's transcriptome and express each cell's identity as a linear combination of the transcriptomes from the three different observed fates (neuron, MEF, myocyte; Supplementary Data 6). Using this method, we observed that there is an initial loss of MEF identity concomitant with an increase in neuronal and myocyte identity over the first five days of Ascl1 reprogramming. The neuronal identity is maintained and matures in day 22 cells transduced with BAM (Extended Data Fig. 6c). However, the day 22 Ascl1-only cells failed to mature to neurons and adopted a predominantly myogenic transcriptional program. This divergence was already apparent in some day 5 cells (Extended Data Fig. 6d, e). These findings raised the question whether the additional two reprogramming factors Brn2 and Myt1l suppress the aberrant myogenic program. Compatible with this notion, we observed that Brn2 and Myt11 had low expression



Figure 2 | **Transgenic** *Ascl1* **silencing explains early reprogramming failure. a**, Hierarchical clustering of day 5 cells using genes correlating positively and negatively with PC1 and PC2 from PCA of day 5 Ascl1-only cells. Note that eGFP fluorescence intensity and Ascl1 mRNA expression shown in the left side bar appear correlated. **b**, Violin plots show the distribution of *Ascl1* and neuronal marker *Tubb3* in day 0 MEFs, as well as Tau–eGFP⁺ and Tau–eGFP⁻ day 5 cells. **c**, qRT–PCR for exogenous Ascl1 expression (top, *n* = 4, biological replicates) and western blot of Ascl1 protein levels (bottom, Supplementary Data 5) for unsorted control MEFs and day 2 cells (NA, not applicable), as well as day 5, day 7, day 10 and day 12 cells FAC-sorted using Tau–eGFP as a neuronal marker. Both

in the five day 22 BAM cells that expressed a myogenic program. To directly address this question, we infected MEFs with Ascl1 alone or in combination with Brn2 and/or Myt11 and assessed myogenic



Figure 3 | iN cell maturation competes with an alternative myogenic cell fate that is repressed by Brn2 and Myt1l. a, tSNE reveals alternative cell fates that emerge during direct reprogramming. Shapes and colours indicate the day 20/22 Ascl1-only (dark green) or day 22 BAM-induced (blue) cells. Note that all cells are Tau-eGFP⁺. b, c, tSNE plot from a with cells coloured based on expression level of marker genes (b), or correlation with bulk RNA-seq data from different purified cell types (neurons²⁴, myocytes²⁵, fibroblasts¹³; c). d, Heat map showing expression of genes marking the two alternative fates in day 20/22 Ascl1-only (upper sidebar, dark green) and day 22 BAM (upper sidebar, blue) Tau-eGFP⁺ cells. Genes (rows) have the highest positive and negative correlation with the first principal component in a PCA analysis on all day 20/22 cells and all genes. Columns represent 121 single cells, ordered based on their correlation coefficient with the first principal component. Lower sidebars, Ascl1 transcript level and Tau-eGFP fluorescence for each cell. e, Immunofluorescent detection of Tau-eGFP (green), DAPI (blue), Myh3 (red) and Tubb3 (cyan) for day 22 cells infected with Ascl1 alone, or with all BAM factors. Images are representative of four biological replicates. Right, mean fractions of eGFP⁺ cells that express either Tubb3 or Myh3. Only Tubb3⁺ cells with a neuronal morphology were counted. Six or seven images were analysed for each of four biological replicates. Error bars, s.e.m.

RNA and protein levels of Ascl1 are significantly higher in Tau–eGFP⁺ cells, and gradually decrease in Tau–eGFP⁻ cells (*P < 0.05, **P < 0.01, ***P < 0.001, two-tailed *t*-test; error bars, s.e.m.). **d**, Schematic for live cell imaging experiment. CD1 MEFs were infected with an eGFP–Ascl1 construct at -1 day, induced with doxycycline at day 0, switched to N3 media at day 1 and imaged between 3 and 6 days post doxycycline. Cells were fixed at 6 days and stained for Tubb3 expression. **e**, Average eGFP–Ascl1 intensity (error bars, s.e.m.) was plotted at 45-min intervals for Tuj1⁺ (n = 10) and Tuj1⁻ (n = 12) cells between day 3 and day 6. Tuj1⁺ cells significantly (one-tailed *t*-test, P < 0.05) increased Ascl1 expression through time compared to Tuj1⁻ cells, which appeared to silence *Ascl1*.

and neurogenic fates at day 22 based on immunostaining and qRT–PCR (Fig. 3e, Extended Data Fig. 6f–i). Indeed, myocyte markers (*Myh3*, *Myo18b*, *Tnnc2*) were upregulated in Tau–eGFP-positive versus negative cells and were strongly repressed when Brn2 and/or Myt1l was overexpressed together with Ascl1. Moreover, Brn2 and Myt1l enhanced the expression of the synaptic genes *Gria2*, *Nrxn3*, *Stmn3*, and *Snap25* but not the immature pan-neuronal genes *Tubb3* and *Map2*. As expected, fibroblast markers were repressed in Tau–eGFP⁺ cells.

We next set out to reconstruct the reprogramming path from MEFs to iN cells. By deconstructing heterogeneity at each time point as described above, we removed cells that appeared stalled in reprogramming due to Ascl1 silencing or cells converging on the alternative myogenic fate. We used quadratic programming to order the cells based on fractional similarity to MEF and neuron bulk transcriptomes. This revealed a continuum of intermediate states through the 22-day reprogramming period (Fig. 4a, b). Notably, the total number of transcripts per single cell decreased as a function of fractional neuron identity (Extended Data Fig. 7a). Our ordering of cells based on fractional identities correlated well with pseudotemporal ordering using Monocle¹⁵, an alternative algorithm for delineating differentiation paths (Extended Data Fig. 7b-d). Heat map visualization of genes identified by PCA of all cells on the iN cell lineage revealed two gene regulatory events during reprogramming with many cells at intermediate stages (Fig. 4c, Supplementary Data 7). First, there is an initiation stage where MEFs exit the cell cycle upon Ascl1 induction, and genes involved in mitosis are turned down or off (such as Birc5, Ube2c, Hmga2). Concomitantly, genes associated with cytoskeletal reorganization (Sept3/4, Coro2b, Ank2, Mtap1a, Homer2, Akap9), synaptic transmission (Snca, Stxbp1, Vamp2, Dmpk, Ppp3ca), and neural projections (Cadm1, Dner, Klhl24, Tubb3, Mapt (Tau)) increase in expression. This indicates that Ascl1 induces genes involved in defining neuronal morphology early in the reprogramming process. The initiation phase is followed by a maturation stage whereby MEF extracelluar matrix genes are turned off and genes involved in synaptic maturation are turned on (Syp, Rab3c, Gria2, Syt4, Nrxn3, Snap25, Sv2a). These results are consistent with previous findings that Tuj1⁺ cells with immature neuron-like morphology can be found as early as three days after Ascl1 induction, while functional synapses are only formed 2 to 3 weeks into the reprogramming process⁸. Finally, we constructed a transcription regulator network on the basis of pairwise correlation of transcription regulator expression across all stages of the MEF-to-iN cell reprogramming. This revealed three densely connected sub-networks identifying transcription regulators influencing MEF cell biology, iN cell initiation, and iN cell maturation



Figure 4 | **Reconstructing the direct reprogramming path from MEFs to iN cells. a**, Top, for each cell on the iN cell reprogramming path, the similarity to bulk RNA-seq from either MEFs¹³ or neurons²⁴ was calculated using quadratic programming and plotted as fractional identities (left axis, circle, fractional MEF identity; right axis, triangle, fractional neuron identity). Points are coloured based on the experimental time point. Bottom, Lagrangian residuals of the quadratic programming for each single cell ordered based on their fractional identity as above. Points are coloured based on the experimental time point. **b**, Fractional neuron identities of all cells on the iN cell reprogramming path are shown as a function of the experimental time point. **c**, Ordering of single cells (rows) according to fractional neuron identity. Genes (columns) with the highest positive and negative correlation to PC1 and PC2 are shown. Left sidebars, experimental time point (green/blue) and fractional neuron

(Fig. 4d, Extended Data Fig. 8, Supplementary Data 8, Supplementary Information). Notably, *Ascl1* was found to strongly positively correlate with the transcription regulators in both the initiation and maturation subnetworks and negatively correlate with transcription regulators specific to MEFs. This data corroborates evidence that persistent *Ascl1* expression is required to maintain chromatin states conducive to iN cell maturation¹³.

It has been suggested that direct somatic lineage reprogramming may not involve an intermediate progenitor cell state as seen during induced pluripotent stem cell differentiation¹⁷⁻¹⁹. However, our fractional analysis showed that the identity of intermediate reprogramming cells could not be explained by a simple linear mixture of the differentiated fibroblast and neuron identities, as revealed by an intermediary increase of Lagrangian residuals (Fig. 4a). Therefore, we tested whether a neural precursor cell (NPC) state is transiently induced by adding NPC bulk transcriptome data along with that of MEFs and neurons into the quadratic programming analysis (Fig. 4e). We found that the fractional NPC identity of cells increased specifically for cells at intermediate positions on the MEF-to-iN cell lineage path, and then decreased as a function of iN cell maturation. In addition, several NPC genes (that is, *Gli3*, *Sox9*, *Nestin*, *Fabp7*, *Hes1*) are expressed in intermediates of the iN cell reprogramming path²⁰ (Fig. 4f). However, canonical NPC marker genes such as Sox2 and Pax6 were never induced. This indicates that cells do not go through a canonical NPC stage, yet a unique intermediate transcriptional state is induced transiently that is unrelated to donor and target cell program

identity (yellow/red). Right sidebars, *Ascl1* transcript levels (log₂[FPKM], blue/yellow) and eGFP fluorescence intensities (log₁₀[RFU], black/ white; RFU, relative fluorescence units). **d**, Transcriptional regulator covariance network during iN cell lineage progression. Shown are nodes (transcriptional regulators) with more than three edges, with each edge reflecting a correlation >0.25 between connected transcriptional regulators. **e**, Fractional MEF (left axis) or fractional neural precursor cell (NPC) identities (right axis) are plotted against fractional neuron identity for single cells on the MEF-to-iN cell lineage. Points are shaped based on the experiment. **f**, Expression of selected genes (columns) that mark NPCs, intermediate progenitor cells (IPCs), neurons, or proliferating cells (Prolif.) are shown for cells on the iN cell lineage (rows). Left sidebars, fractional neuron identity (yellow/red) and experimental time point (green/blue).

similar to that which was observed for induced pluripotent stem cell reprogramming $^{21\text{-}23}$

A fundamental question in cell reprogramming is whether there are pre-determined mechanisms that prevent the majority of the fibroblasts from reprogramming or whether all donor cells are competent to reprogram but the reprogramming procedure is inefficient. We did not observe any MEF subpopulations, other than cell cycle variation, that suggested differences in the capacity to initiate reprogramming. Furthermore, we observed that 48 h after infection the majority of the cells induced Ascl1-target genes and silenced MEF-associated genes. This does not preclude the possibility that underlying epigenetic variation in donor cells influences reprogramming outcomes; however, our analysis suggests that it is unlikely that MEF heterogeneity contributes significantly to reprogramming efficiency. We found that divergence from the neuronal differentiation path into an alternative myogenic fate, as well as Ascl1 transgene silencing, were both significant factors contributing to reprogramming efficiency. Though Ascl1 induces lineage conversion, it is inefficient in restricting cells to the neuronal fate. This suggests that intermediate stages of iN cell progression are unstable, perhaps due to epigenetic barriers, and additional factors promote cells to permanently acquire neuron-like identity, rather than revert to MEF-like or diverge towards the alternative myocyte-like fate. In summary, we present a single-cell transcriptomic approach that can be used to dissect direct cellular reprogramming pathways or developmental programs in which cells transform their identity through a series of intermediate states.

Online Content Methods, along with any additional Extended Data display items and Source Data, are available in the online version of the paper; references unique to these sections appear only in the online paper.

Received 10 March 2015; accepted 17 May 2016.

Published online 8 June 2016.

- Xu, J., Du, Y. & Deng, H. Direct lineage reprogramming: strategies, mechanisms, and applications. Cell Stem Cell 16, 119–134 (2015).
- Arlotta, P. & Berninger, B. Brains in metamorphosis: reprogramming cell identity within the central nervous system. *Curr. Opin. Neurobiol.* 27, 208–214 (2014).
- 3. Graf, T. Historical origins of transdifferentiation and reprogramming. *Cell Stem Cell* **9**, 504–516 (2011).
- Treutlein, B. et al. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell RNA-seq. Nature 509, 371–375 (2014).
- Shalek, A. K. et al. Single-cell transcriptomics reveals bimodality in expression and splicing in immune cells. *Nature* 498, 236–240 (2013).
- Zeisel, A. *et al.* Brain structure. Cell types in the mouse cortex and hippocampus revealed by single-cell RNA-seq. *Science* **347**, 1138–1142 (2015).
- Ramsköld, D. et al. Full-length mRNA-seq from single-cell levels of RNA and individual circulating tumor cells. *Nat. Biotechnol.* **30**, 777–782 (2012).
- Vierbuchen, T. et al. Direct conversion of fibroblasts to functional neurons by defined factors. Nature 463, 1035–1041 (2010).
- 9. Pfisterer, U. *et al.* Direct conversion of human fibroblasts to dopaminergic neurons. *Proc. Natl Acad. Sci. USA* **108**, 10343–10348 (2011).
- Yoo, A. S. *et al.* MicroRNA-mediated conversion of human fibroblasts to neurons. *Nature* 476, 228–231 (2011).
- Ambasudhan, R. *et al.* Direct reprogramming of adult human fibroblasts to functional neurons under defined conditions. *Cell Stem Cell* 9, 113–118 (2011).
- 12. Caiazzo, M. et al. Direct generation of functional dopaminergic neurons from mouse and human fibroblasts. Nature **476**, 224–227 (2011).
- 13. Wapinski, O. L. *et al.* Hierarchical mechanisms for direct reprogramming of fibroblasts to neurons. *Cell* **155**, 621–635 (2013).
- 14. Chanda, S. et al. Generation of induced neuronal cells by the single reprogramming factor ASCL1. Stem Cell Rep. **3**, 282–296 (2014).
- Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nat. Biotechnol.* 32, 381–386 (2014).
- Buettner, F. et al. Computational analysis of cell-to-cell heterogeneity in single-cell RNA-sequencing data reveals hidden subpopulations of cells. *Nat. Biotechnol.* 33, 155–160 (2015).
- Merkle, F. T. & Eggan, K. Modeling human disease with pluripotent stem cells: from genome association to function. *Cell Stem Cell* 12, 656–668 (2013).
- Perrier, A. L. *et al.* Derivation of midbrain dopamine neurons from human embryonic stem cells. *Proc. Natl Acad. Sci. USA* **101**, 12543–12548 (2004).
- Li, X. J. et al. Specification of motoneurons from human embryonic stem cells. Nat. Biotechnol. 23, 215–221 (2005).

- Camp, J. G. *et al.* Human cerebral organoids recapitulate gene expression programs of fetal neocortex development. *Proc. Natl Acad. Sci. USA* **112**, 15672–15677 (2015).
- Di Stefano, B. et al. C/EBPα poises B cells for rapid reprogramming into induced pluripotent stem cells. Nature 506, 235–239 (2014).
- 22. Lujan, E. *et al.* Early reprogramming regulators identified by prospective isolation and mass cytometry. *Nature* **521**, 352–356 (2015).
- Takahashi, K. *et al.* Induction of pluripotency in human somatic cells via a transient state resembling primitive streak-like mesendoderm. *Nat. Commun.* 5, 3678 (2014).
- Zhang, Y. *et al.* An RNA-sequencing transcriptome and splicing database of glia, neurons, and vascular cells of the cerebral cortex. *J. Neurosci.* 34, 11929–11947 (2014).
- Trapnell, C. et al. Transcript assembly and quantification by RNA-Seq reveals unannotated transcripts and isoform switching during cell differentiation. *Nat. Biotechnol.* 28, 511–515 (2010).

Supplementary Information is available in the online version of the paper.

Acknowledgements The authors would like to acknowledge B. Passarelli and B. Vernot for discussions regarding bioinformatic pipelines, P. Lovelace for support with FACS and other Quake and Wernig laboratory members for discussions and support. This work was supported by NIH grant RC4NS073015-01 (M.W., S.Q.R., B.T.), the Stinehart-Reed Foundation, the Ellison Medical Foundation, the New York Stem Cell Foundation, CIRM grant RB5-07466 (all to M.W.), a National Science Scholarship from the Agency for Science, Technology and Research (Q.Y.L.), NIH grant GM092925 (S.A.M.S., J.S.), the German Research Foundation (M.M.) and a PhRMA foundation Informatics fellowship (J.G.C.). S.R.Q. is an investigator of the Howard Hughes Medical Institute. M.W. is a New York Stem Cell Foundation (NYSCF) Robertson Investigator and a Tashia and John Morgridge Faculty Scholar at the Child Health Research Institute at Stanford.

Author Contributions B.T., Q.Y.L., M.W. and S.R.Q. conceived the study and designed the experiments. Q.Y.L. performed direct reprogramming. qRT–PCR and western blot experiments; B.T., Q.Y.L., and S.S. performed single-cell RNA-seq experiments; N.F.N. assisted with single-cell RNA-seq experiments and sequenced the libraries; Q.Y.L., S.A.M.S. and M.M. performed time-lapse imaging experiments. B.T., J.G.C. and W.K. analysed single-cell RNA-seq data, Q.Y.L. analysed qRT–PCR and time-lapse imaging data, J.M.S., M.W. and S.R.Q. provided intellectual guidance in the interpretation of the data. B.T., Q.Y.L., J.G.C., M.W., and S.R.Q. wrote the paper.

Author Information The single-cell RNA-seq data were deposited on NCBI GEO with the accession number GSE67310. Reprints and permissions information is available at www.nature.com/reprints. The authors declare competing financial interests: details are available in the online version of the paper. Readers are welcome to comment on the online version of the paper. Correspondence and requests for materials should be addressed to S.R.Q. (quake@stanford.edu) and M.W. (wernig@stanford.edu).

Reviewer Information *Nature* thanks F. Tang and the other anonymous reviewer(s) for their contribution to the peer review of this work.

METHODS

No statistical methods were used to predetermine sample size. The experiments were not randomized and the investigators were not blinded to allocation during experiments and outcome assessment.

Cell derivation, cell culture and iN cell generation. Tau–eGFP reporter MEFs, tested negative for mycoplasma contamination, were isolated, infected with doxy-cycline (dox)-inducible lentiviral constructs and reprogrammed into iN cells as previously described⁸. Day 0 (d0) cells were uninfected MEFs that served as a negative control. Day 2 (d2) cells were infected with Ascl1 and harvested two days after dox-induction. Day 5 (d5) cells were infected with Ascl1, FAC-sorted for Tau–eGFP⁺ and Tau–eGFP⁻ cells five days after dox induction and the two cell populations were mixed again in a 1:1 ratio. Day 20 or 22 (d20/d22) cells were infected either with Ascl1 alone, or combined with Brn2 and Myt1l, plated with glia seven days post dox induction. Each of these groups was then loaded onto separate microfluidic mRNA-seq chips for preparation of pre-amplified cDNA from single cells.

Clonal Ascl1-inducible MEFs were derived as previously described¹³. Twelvewell plates were coated with Matrigel and incubated at 37 °C overnight. 350,000 cells were then plated per well and kept in MEF media. Dox was added a day after plating. For single-cell RNA-seq, cells were harvested two days post dox induction and loaded onto a microfluidic mRNA-seq chip. To evaluate efficiency in reprogramming, MEF + dox media was switched out for N3 + dox media after 48 h, and cells were fixed for immunostaining 12 days post dox.

Capturing of single cells and preparation of cDNA. Single cells were captured on a medium-sized (10-17µm cell diameter) microfluidic RNA-seq chip (Fluidigm) using the Fluidigm C1 system. Cells were loaded onto the chip at a concentration of 350–500 cells µl⁻¹, stained for viability (live/dead cell viability assay, Molecular Probes, Life Technologies) and imaged by phase-contrast and fluorescence microscopy to assess number and viability of cells per capture site. For d5 and d22 experiments, cells were only stained with the dead stain ethidium homodimer (emission ~635 nm, red channel) and Tau-eGFP fluorescence was imaged in the green channel. Only single, live cells were included in the analysis. cDNAs were prepared on chip using the SMARTer Ultra Low RNA kit for Illumina (Clontech). ERCC (External RNA Controls Consortium) RNA spike-in Mix (Ambion, Life Technologies)^{26,27} was added to the lysis reaction and processed in parallel to cellular mRNA. Tau-eGFP fluorescence intensity of each single cell was determined using CellProfiler²⁸ by first identifying the outline of the cell in the image of the respective capture site and then integrating over the signal in the eGFP channel. RNA-seq library construction and cDNA sequencing. Size distribution and concentration of single-cell cDNA was assessed on a capillary electrophoresis based fragment analyser (Advanced Analytical Technologies) and only single cells with high quality cDNA were further processed. Sequencing libraries were constructed in 96-well plates using the Illumina Nextera XT DNA Sample Preparation kit according to the protocol supplied by Fluidigm and as described previously²⁹. Libraries were quantified by Agilent Bioanalyzer using High Sensitivity DNA analysis kit as well as fluorometrically using Qubit dsDNA HS Assay kits and a Qubit 2.0 Fluorometer (Invitrogen, Thermo Fisher Scientific). Up to 110 single-cell libraries were pooled and sequenced 100 bp paired-end on one lane of Illumina HiSeq 2000 or 75 bp paired-end on one lane of Illumina NextSeq 500 to a depth of 1-7 million reads. CASAVA 1.8.2 was used to separate out the data for each single cell using unique barcode combinations from the Nextera XT preparation and to generate *.fastq files. In total, the transcriptome of a total of 405 cells was measured from the following eight independent experiments: d0 (73 cells, 1 experiment), d2 (Ascl1-only in regular MEFs, 81 cells, 1 experiment; Ascl1-only in clonal MEFs, 47 cells, 1 experiment), d5 (Ascl1-only, 55 cells, 1 experiment) and d20 (Ascl1-only, 33 cells, 1 experiment) and d22 (BAM, 43 cells, 1 experiment; Ascl1-only, 34 and 39 cells, 2 independent experiments). See Supplementary Data 1 for the transcriptome data for all 405 cells with annotations (quantification in log₂[FPKM]).

Processing, analysis and graphic display of single cell RNA-seq data. Raw reads were pre-processed with sequence grooming tools FASTQC³⁰, cutadapt³¹, and PRINSEQ³² followed by sequence alignment using the Tuxedo suite (Bowtie³³, Bowtie2³⁴, TopHat³⁵ and SAMtools³⁶) using default settings. Transcript levels were quantified as fragments per kilobase of transcript per million mapped reads (FPKM) generated by TopHat/ Cufflinks²⁵.

After seven days of reprogramming, Tau–eGFP reporter MEFs (with C57BL/6J and 129S4/SvJae background) were co-cultured with glia derived from CD-1 mice. To determine if any feeder cells contaminated the 20–22-day time points, we used the single cell RNA-seq reads to identify positions that differ from the mouse reference genome (mm10, built from strain C57BL/6J mice). We used the mpileup fuction in samtools to generate a multi-sample variant call format file (vcf), and a custom python script to genotype the cells by requiring coverage in all cells for all positions, with a coverage depth of five reads, a phred GT likelihood = 0 for called genotype and \geq 40 for next-best genotype. This resulted in 95 informative sites

distinguishing more than one cell from the reference genome. We clustered cells based on their genotype (homozygous reference, heterozygous, homozygous alternate), and identified cells that were strongly different from the reference genome. These cells expressed either astrocyte (*Gfap*) or microglia marker genes suggesting they were contaminants from the feeder cell culture. We removed these cells from subsequent analyses.

Approximate number of transcripts was calculated from FPKM values by using the correlation between number of transcripts of exogenous spike-in mRNA sequences and their respective measured mean FPKM values (Extended Data Fig. 2). The number of spike-in transcripts per single cell lysis reaction was calculated using the concentration of each spike-in provided by the vendor (Ambion, Life Technologies), the approximate volume of the lysis chamber (10 nl) as well as the dilution of spike-in transcripts in the lysis reaction mix ($40,000 \times$). Transcript levels were converted to the log-space by taking the logarithm to the base 2 (Supplementary Data 1). R studio³⁷ (https://www.rstudio.com/) was used to run custom R³⁸ scripts to perform principal component analysis (PCA, FactoMineR package), hierarchical clustering (stats package), variance analysis and to construct heat maps, correlation plots, box plots, scatter plots, violin plots, dendrograms, bar graphs, and histograms. Generally, ggplot2 and gplots packages were used to generate data graphs.

The Seurat package^{39,40} implemented in R was used to identify distinct cell populations present at d22 of Ascl1-only and BAM reprogramming (Fig. 3a–d). t-distributed stochastic neighbour embedding (tSNE) was performed on all d20/d22 cells using the most significant genes ($P < 1 \times 10^{-3}$, with a maximum of 100 genes per principal component) that define the first three principal components of a PCA analysis on the data set. To further estimate the identity of each cell on the tSNE plot, we colour coded cells based on Pearson correlation of each single cell's expression profile with the expression profile of bulk cortical neurons^{13,24}, myocytes²⁵, and MEFs¹³ (Fig. 3). The Monocle package¹⁵ was used to order cells on a pseudo-time course during MEF to iN cell reprogramming (Extended Data Fig. 7). Covariance network analysis and visualizations were done using igraph implemented in R⁴¹ (http://igraph.sf.net).

To generate PCA plots and heat maps in Figs 1c–e, 2a, 3a and 4c, PCA was performed on cells using all genes expressed in more than two cells and with a variance in transcript level (\log_2 [FPKM]) across all single cells greater than 2. This threshold resulted generally in about 8,000–12,000 genes. Subsequently, genes with the highest PC loadings (highest (top 50–100) positive or negative correlation coefficient with one of the first one to two principal components) were identified and a heat map was plotted with genes ordered based on their correlation coefficient with the respective PC (Figs 1e, 2a, 4c). Cells in rows were ordered based on unsupervised hierarchical clustering using Pearson correlation as distance metric (Figs 1e, 2a) or based on their fractional identity as determined by quadratic programming (Fig. 4c, see below)

Gene ontology enrichment analyses were performed using DAVID Bioinformatics Resources 6.7 of the National Institute of Allergy and Infectious Diseases⁴². Functional annotation clustering was performed and GO terms representative for top enriched annotation clusters are shown in Fig. 1f, Extended Data Figs 1e and 4a with their Bonferroni corrected *P* values. In addition, results of GO enrichment analyses are provided in the Supplementary Data.

To express a single cell transcriptome as a linear combination of primary cell type transcriptomes, we used published bulk RNA-seq data sets for primary murine neurons²⁴, myocytes²⁵, and embryonic fibroblasts¹³ (Extended Data Fig. 6b, c), neurons²⁴ and embryonic fibroblasts¹³ (Fig. 4a) or neurons²⁴, embryonic fibroblasts¹³ and neuronal progenitor cells¹³ (Fig. 4e). In each quadratic programming analysis, we first identified genes that were specifically (log₂ fold change of 3 or higher) expressed in each of the bulk data sets compared to the respective others (Supplementary Data 6). Using these genes, we then calculated the fractional identities of each single cell using quadratic programming (R package 'quadprog'). The resulting fractional neuron identities of cells on the MEF-to-iN cell reprogramming path (265 cells in total, excluding cells that were Tau-eGFP-negative at d5 or myocyte- and fibroblast-like cells at d22) were used to order cells in a pseudotemporal manner (Fig. 4a-c, e, f). We compared this fractional neuron identity based cell ordering with pseudo-temporal ordering of cells based on Monocle (Extended Data Fig. 7b-d), an algorithm that combines differential dimension reduction using independent component analysis with minimal spanning tree construction to link cells along a pseudotemporally ordered path¹⁵. Monocle analysis was performed using genes differentially expressed between neuron²⁴ and embryonic fibroblast¹³ bulk RNA-seq data (same gene set that was used when calculating fractional neuron and fibroblast identities in Fig. 4a, genes listed in Supplementary Data 6).

For the transcription factor network analysis (Fig. 4d), we computed a pairwise correlation matrix (Pearson correlation, visualized in correlogram in Extended Data Fig. 8a) for transcriptional regulators annotated as such in the Animal

Transcription Factor Database (http://www.bioguo.org/AnimalTFDB/)43 and identified those transcriptional regulators (TRs) with a Pearson correlation of greater than 0.35 with at least five other TRs (82 TRs, shown in Extended Data Fig. 8b). We used a permutation approach to determine the probability of finding TRs meeting this threshold by chance. We performed 500 random permutations of the expression matrix of all TRs across cells on the MEF-to-iN cell lineage, and calculated the pairwise correlation matrix for each permutation of the input data frame. All randomized data frames resulted in 0 TRs that met our threshold. This shows that our correlation threshold is strict, and all nodes and connections that we present in the TR network are highly unlikely to be by chance. We used the pairwise correlation matrix for the selected TRs as input into the function graph.adjacency() of igraph implemented in R⁴¹ (http://igraph.sf.net) to generate a weighted network graph, in which the selected TRs are presented as vertices and all pairwise correlations >0.25 are presented as edges linking the respective vertices. The network graph was visualized using the Fruchterman-Reingold layout and the three clear subnetworks (MEF, initiation, maturation) were manually colour coded.

We used Pearson correlation of each single cell expression profile with the expression profile of bulk cortical neurons^{13,24}, myocytes²⁵, and MEFs¹³ to further estimate the identity of each single cell and to estimate when alternative fates emerge (Fig. 3c, Extended Data Fig. 6d, e). For this analysis, we considered the same cell type specific gene sets that were used in the quadratic programming analysis, that is, were genes specifically expressed (log₂ fold change of 3 or higher) in a respective bulk RNA-seq data set compared to the others (Supplementary Data 6).

To estimate intercellular heterogeneity of d0 MEFs, we calculated the variance for each gene across all MEF cells as well as across mouse embryonic stem cells under 2iLIF culture conditions⁴⁴ and across glioblastoma cells⁴⁵. We then plotted the distribution of variances for all genes per cell population as box plots.

Quantitative RT-PCR and immunostaining. Ascl1 infected Tau–eGFP reporter MEFs were FAC-sorted 5, 7, 10, 12 or 22 days post-Ascl1 induction with dox. RNA was then extracted from both Tau–eGFP positive and negative populations from each time point, as well as uninfected control MEFs and unsorted d2 Ascl1-infected MEFs using the TRIzol RNA isolation protocol (Invitrogen, 15596-018). Reverse transcription into cDNA was performed using the SuperScript III First-strand Synthesis System (Invitrogen, 18080-051) and qRT–PCR was performed using Sybr Green (Thermo Fisher Scientific, 4309155). Immunostaining was performed as previously described⁸. Antibodies and qRT–PCR primers are listed in the Supplementary Information.

Time-lapse imaging of Ascl1 expression. MEFs were isolated from E13.5 CD-1 embryos (Charles River) and infected with a dox-inducible, N-terminal-tagged eGFP-Ascl1 fusion construct using the protocol previously described¹. Cells were plated on 35 cm glass bottom dishes (MatTek), coated with polyorthinine (Sigma P3655) and laminin (Invitrogen 23017-015). Imaging experiments were performed between 3 and 6 days post dox induction, in a temperature- and CO₂-controlled chamber. Images were taken for up to 10 positions per dish, for 3 dishes, every 45 min with a Zeiss AxioVert 200M microscope with an automated stage using an EC Plan-Neofluar 5×/0.16 NA Ph1 objective or an A-plan 10×/0.25 NA Ph1 objective. Cells were fixed at 6 days and immunostained using Tuj1 antibodies recognizing neuronal Tubb3 (Covance MRB-435P) to confirm neuronal identity. We used ImageJ to segment individual cells and measure the level of GFP for 7 Tuj1+ cells and 7 Tuj1⁻ cells over time. Average intensity was obtained by normalizing the average intensity of a cell segment by the average background intensity of an adjacent segment of the same size. A *t*-test was performed comparing Tuj1⁺ and Tuj1⁻ cells at each time point to evaluate significance.

Antibodies. Rabbit anti-Ascl1 (Abcam ab74065), chicken anti-GFP (Abcam ab13970), rabbit anti-Tubb3 (Covance MRB-435P), mouse anti-Tubb3 (Covance MMS-435P), mouse anti-Map2 (Sigma M4403), rabbit anti-Myh3 (Santa Cruzsc-20641), goat anti-Dlx3 (Santa Cruz sc-18143), mouse anti- β -Actin (Sigma A5441), rabbit anti-Tcf12 (Bethyl A300-754A).

Primers. *General. Gapdh* (forward: AGGTCGGTGTGAACGGATTTG, reverse: TGTAGACCATGTAGTTGAGGTCA); *Ascl1* (TetO) (forward: CCGAA TTCGCTAGCCACCAT, reverse: AAGAAGCAGGCTGCGGG).

Initiation factors. Atoh8 (forward: GCCAAGAAACGGAAGGAGTGA, reverse: CTGAGAGATGGTACACGGGC); Dlx3 (forward: CGCCGCTCCAA GTTCAAAAA, reverse: GTGGTACCAGGAGTTGGTGG); Hes6 (forward: TACCGAGGTGCAGGCCAA, reverse: AGTTCAGCTGAGACAGTGGC); Sox11 (forward: CCTGTCGCTGGTGGATAAGG, reverse: CTGCGCCTCTC AATACGTGA); Sox9 (forward: CGAGCACTCTGGGCAATCTCA, reverse: ATGACGTCGCTGCTCAGTTC); Tcf4 (forward: CAGTGCGATGT TTTCGCCTC, reverse: ATGTGACCCAAGATCCCTGC); Tcf12 (forward: GTCTCGAATGGAAGACCGCT, reverse: GTTCCGACCATCGAAGCTGA). Maturation factors. Camta1 (forward: CCCCTAAGACAAGACCGCAG, reverse: ACATAGCAGCCGTACAAGCA); Insm1 (forward: GACCCGG CACATCAACAAGT, reverse: GAAGCGAAGCGAAGAGGACA); My111 (forward: ATGTTCCCACAACCACCAC, reverse: TACCGCTTGGCATCG TCATA); St18 (forward: TGCCAAGGAGCTGAGATGAA, reverse: GAAGG

CTGCTTGCGTTGAAT). Neuronal genes. Gria2 (forward: GGGGACAAGGCGTGGAAATA, reverse: GTACCCAATCTTCCGGGGGTC); Map2 (forward: CAGAGAAA CAGCAGAGGAGGT, reverse: TTTGTTCTGAGGCTGGCGAT); Nrxn3 (forward: TGTGAACCAAGTACAGATAAGAGT, reverse: CAGCTCAGGGGAC AAAGAGG); Snap25 (forward: TTCATCCGCAGGGTAACAAA, reverse: GTTGCACGTTGGTTGGCTT); Stmn3 (forward: AGCACCGT ATCTGCCTACAAG, reverse: TGGTAGATGGTGTTCGGGTG); Tubb3 (forward: CAGATAGGGCCAAGTTCTGG, reverse: GTTGTCGGGCCTGAATAGGT). Myocyte genes. Acta1 (forward: CTAGACACCATGTGCGGACGA, reverse: CATACCTACCATGACACCCTGG); Myh3 (forward: AAATGAAGGGGACG CTGGAG, reverse: CAGCTGGAAGGTGACTCTGG); Myo18b (forward: GCCCTCTTCAGGGAAGGTA, reverse: GAGCTTCTCCACTGACACCC); Innc2 (forward: CAACCATGACGGACCAACAG, reverse: GTGTCTGCC CTAGCATCCTC).

Fibroblast genes. Col1a2 (forward: AGTCGATGGCTGCTCCAAAA, reverse: ATTTGAAACAGACGGGGCCA); *Dcn* (forward: GCAAAATCAGT CCAGAGGCA, reverse: CGCCCAGTTCTATGACAAGC).

- Baker, S. C. et al. The External RNA Controls Consortium: a progress report. Nat. Methods 2, 731–734 (2005).
- Jiang, L. et al. Synthetic spike-in standards for RNA-seq experiments. Genome Res. 21, 1543–1551 (2011).
- Carpenter, A. E. et al. CellProfiler: image analysis software for identifying and quantifying cell phenotypes. *Genome Biol.* 7, R100 (2006).
- Wu, A. R. et al. Quantitative assessment of single-cell RNA-sequencing methods. Nat. Methods 11, 41–46 (2014).
- Babraham Institute. Babraham Bioinformatics. FASTQC. http://www. bioinformatics.bbsrc.ac.uk/projects/fastgc
- Martin, M. Cutadapt removes adapter sequences from high-throughput sequencing reads. *EMBnet journal* 17, 10–12 (2011).
- Schmieder, R. & Edwards, R. Quality control and preprocessing of metagenomic datasets. *Bioinformatics* 27, 863–864 (2011).
- Langmead, B., Trapnell, C., Pop, M. & Salzberg, S. L. Ultrafast and memoryefficient alignment of short DNA sequences to the human genome. *Genome Biol.* **10**, R25 (2009).
- Langmead, B. & Salzberg, S. L. Fast gapped-read alignment with Bowtie 2. Nat. Methods 9, 357–359 (2012).
- Trapnell, C., Pachter, L. & Salzberg, S. L. TopHat: discovering splice junctions with RNA-seq. *Bioinformatics* 25, 1105–1111 (2009).
- Li, H. et al. The sequence alignment/map format and SAMtools. Bioinformatics 25, 2078–2079 (2009).
- RStudio. Integrated Development for R. RStudio, Inc., Boston, MA URL http://www.rstudio.com/ (2015).
- R Core Team. R: A language and environment for statistical computing. R Foundation for Statistical Computing. http://www.R-project.org/
- Macosko, E. Z. et al. Highly parallel genome-wide expression profiling of individual cells using nanoliter droplets. Cell 161, 1202–1214 (2015).
- Satija, R., Farrell, J. A., Gennert, D., Schier, A. F. & Regev, A. Spatial reconstruction of single-cell gene expression data. *Nat. Biotechnol.* 33, 495–502 (2015).
- Csardi, G. & Nepusz, T. The igraph software package for complex network research. InterJournal 1695 (2006).
- Huang, W., Sherman, B. T. & Lempicki, R. A. Systematic and integrative analysis of large gene lists using DAVID bioinformatics resources. *Nat. Protocols* 4, 44–57 (2009).
- Zhang, H. M. et al. AnimalTFDB: a comprehensive animal transcription factor database. Nucleic Acids Res. 40, D144–D149 (2012).
- Kumar, R. M. et al. Deconstructing transcriptional heterogeneity in pluripotent stem cells. Nature 516, 56–61 (2014).
- Patel, A. P. et al. Single-cell RNA-seq highlights intratumoral heterogeneity in primary glioblastoma. Science 344, 1396–1401 (2014).



Extended Data Figure 1 | The majority of MEFs are actively undergoing cell cycle, but exit cell cycle upon Ascl1 induction. a, Live cell imaging of Tau-eGFP reporter over the course of BAM-mediated iN cell reprogramming. Tau-eGFP fluorescence normalized to the maximum expression is shown in relation to days post-BAM induction. Tau-eGFP expression began at day 5 and reached a peak at day 8 after induction. Shown are representative images from day 0, day 5 and day 9. b, Box plots of intercellular transcriptome variance showed that MEFs are more heterogeneous than mouse embryonic stem cells under 2iLIF culture conditions⁴⁴ and less heterogeneous than glioblastoma cells⁴⁵. c, PCA of genes with most variance in day 0 MEFs revealed MEF heterogeneity (blue, A). Density plot showing the distribution of number of cells along PC1 loading is shown above the PCA plot. d, Heat map and hierarchical clustering of genes used for the PCA in panel c shows to major MEF subpopulations. Each column represents a single cell, and each row a gene. Subpopulation A is highlighted in blue in the dendrogram. e, GO enrichment for genes in c shows that MEF subpopulation A is distinguished by the low or lack of expression of genes enriched for cell cycle terms. f, g, PCA and heat map of the same genes used in panels c-e, this time including day 0 MEFs (circles, light green) and day 2 cells (squares, dark green), showed that most of the day 2 cells had the same cell cycle signature as MEF subpopulation A. Cells in columns of both heat maps are ordered based on PC1 loading.



Extended Data Figure 2 | **Total number of transcripts per cell decreases during MEF-to-iN cell reprogramming. a**, Average detected transcript levels (mean FPKM, log₂) for 92 ERCC RNA spike-ins as a function of provided number of molecules per lysis reaction for each of the 8 independent single-cell RNA-seq experiments. Linear regression fits through data points are shown. The length of each ERCC RNA spike-in transcript is encoded in the size of the data points. No particular bias towards the detection of shorter versus longer transcripts is observed. The linear regression fit was used to convert FPKM values to approximate number of transcripts. b, Box plots showing the distribution of the total

number of transcripts per single cell for each experiment. Number of transcripts per cell were calculated from the FPKM values of all genes in each cell using the correlation between number of transcripts of exogenous spike-in mRNA sequences and their respective measured mean FPKM values (calibration curves are shown in panel **a**). The total number of transcripts expressed by a single cell and detected by single-cell RNA-seq is highest in MEFs and is more than twofold decreased upon overexpression of Ascl1 or BAM. **c**, Box plots showing the distribution of the median transcript number per gene across all cells of one experiment. The distributions are similar over the course of iN cell reprogramming.



Extended Data Figure 3 | Clonal MEFs reprogram successfully into iN cells, and Ascl1-only and BAM induce similar responses during early iN cell reprogramming. a, Immunostaining of heterogenous Ascl1-infected MEFs and clonal MEFs with homogenous Ascl1 transgene insertions, fixed 12 days after Ascl1 induction, using rabbit anti-Tubb3 (red) and mouse anti-Map2 (cyan) antibodies and DAPI (blue) as a nuclear stain. Reprogramming efficiencies are comparable regardless of variation in Ascl1 copy numbers. Images are representative for one reprogramming experiment. b, Bar plots showing expression of Ascl1-target genes (*Hes6*, *Zfp238*, *Snca*, *Cox8b*, *Bex1*, *Dner*) and MEF marker genes averaged across single cells from day 0 MEFs and day 2 Ascl1-only cells, as well as from bulk RNA-seq data from MEFs, day 2 BAM, and day 2 Ascl1-only cells. This data shows that the initiation of reprogramming at day 2 is similar for Ascl1-alone and BAM-mediated reprogramming.



Extended Data Figure 4 | Failed reprogramming at day 5 correlates with silencing of Ascl1. a, Bonferroni-corrected *P* values for gene ontology enrichments are shown for each group of genes from Fig. 2a, with representative genes listed (Supplementary Data 4). b, Biplot showing Tau–eGFP fluorescence intensity as a function of *Ascl1* transcript level in day 5 cells. Point size is proportional to eGFP transcript levels in log₂[FPKM]. There is a positive correlation (R^2 =0.49) indicating that cells with higher Ascl1 expression are more likely to reprogram. c, Heat map of eGFP–Ascl1 expression in 14 individual cells (columns) during live cell imaging. Rows represent time post Ascl1 induction in 45-min intervals.



Extended Data Figure 5 | Live cell imaging shows diminishing of eGFP-Ascl1 signal in cells that fail to reprogram. a, Immunostaining for Tubb3 and Map2 at day 12 post induction of Ascl1, C-terminal tagged Ascl1–eGFP and N-terminal tagged eGFP-Ascl1 in CD-1 MEFs. eGFP-Ascl1 has comparable reprogramming efficiency with untagged Ascl1 while Ascl1–eGFP has a much reduced reprogramming efficiency, so eGFP-Ascl1 was chosen for live cell imaging. Images are representative for one reprogramming experiment per condition. **b**, Representative images from live cell imaging showing an example of diminishing of eGFP signal in a cell that failed to reprogram (that is, cell was Tuj1-negative at day 6). **c**, Live cell imaging of eGFP signal of eGFP–Ascl1 infected MEFs between 3–6 days post dox induction. **d**, eGFP imaging of live cells 6 days post induction of Ascl1 and corresponding immunostaining for Tubb3 after fixation.



Extended Data Figure 6 | Brn2 and Myt1l repress alternative fates that compete with the iN cell fate during advanced Ascl1 reprogramming. a, Scatter plot showing PC1 and PC2 loadings from principal component analysis (PCA) of single cells from all time points with experimental time point and reprogramming condition (Ascl1 versus BAM) encoded in point shape and colour. b, Overview of quadratic programming. Fractional identities are calculated assuming a linear combination of different cell fates. c, Biplots showing the fractional fibroblast identity as a function of fractional neuron (left) and fractional myocyte (right) identity for each cell with points shaped and colour coded based on reprogramming time point and condition. **d**, Correlation of transcriptomes from days 0, 2, 5, and 20/22 cells (Ascl1-only and BAM-induced) with bulk RNA-seq from MEFs, cortical neurons and myocytes. Bottom bars show Tau-eGFP fluorescence intensity. e, Bar plot quantifying the number of cells with a maximum correlation to bulk RNA-seq data from each of the observed fates (d). f, Immunofluorescent detection of Tau-eGFP (green), DAPI (blue), Myh3 (red) and Tubb3 (cyan) for day 22 cells that were infected



with Ascl1 co-infected with Brn2 or Myt1l. See Fig. 3e for respective data for cells infected with Ascl1-only or all three BAM factors. Images are representative for four biological replicates. Right, mean fractions of eGFP⁺ cells that express either Tubb3 or Myh3. Only Tubb3⁺ cells with a neuronal morphology were counted. Co-expression of Ascl1 with Brn2 and/or Myt1l increases fraction of Tau-eGFP⁺ cells that are also Tubb3⁺, while decreasing the number of cells that are Myh3⁺. Six or seven images were analysed for each of four biological replicates. Error bars, s.e.m. g-i, qRT-PCR of selected myogenic (g), neuronal (h), and fibroblast (i) markers using day 22 cells that are infected with Ascl1 only or co-infected with Brn2 or Myt1l or both and FAC-sorted by Tau-eGFP (n = 3, biological replicates; error bars, s.e.m.). Myogenic genes were significantly downregulated in Tau-eGFP⁺ cells that were co-infected with Brn2 and/or Myt1l compared to those infected with Ascl1 alone, while some neuronal genes are significantly upregulated (Map2, Gria) (**P* < 0.05, ***P* < 0.01, ****P* < 0.001, two-tailed *t*-test).



Extended Data Figure 7 | Comparison of Monocle and quadratic programming with respect to ordering of neuronal cells through the reprogramming path. a, Biplot showing the total number of transcripts per cell for all cells on the MEF-to-iN cell lineage as a function of the fraction neuron identity of each cell (see Fig. 4). The total number of transcripts decreases during the reprogramming process. b, Cells (depicted as circles) are arranged in the 2D independent component space based on the expression of genes used for quadratic programming in Fig. 4a. Lines connecting cells represent the edges of a minimal spanning tree with the bold black line indicating the longest path. Time points are colour coded. c, Monocle plots with single cells coloured based on gene expression that distinguishes the stages of iN cell reprogramming. d, Biplot shows the correlation between ordering of cells based on pseudo-time (Monocle) and fractional identity (quadratic programming). Time points are colour coded. Pearson correlation coefficient = 0.91.





Extended Data Figure 8 | Neuronal maturation proceeds through expression of distinct transcriptional regulators. a, Correlogram showing transcriptional regulators (TRs) highly correlated within MEFs as well as the initiation phase and the maturation phase of reprogramming. b, Heat map shows expression of TRs that control the two stages of MEF to iN cell reprogramming (Fig. 4d) in cells ordered based on fractional neuron identity. Each row represents a single cell, each column a gene. Experimental time point (green/blue sidebar) and fractional neuron identity (yellow/red sidebar) are shown at the top. c-e, Pseudo-temporal expression dynamics of exemplary TRs marking the initiation stage (c) and the maturation stage (d) of iN cell reprogramming as well as MEF identity (e). Transcript levels of the TRs are shown across all single

cells on the MEF-to-iN cell lineage ordered based on fractional neuron identity. Growth curves based on a model-free spline method were fitted to the data. f, qRT–PCR of selected TRs from initiation and maturation subnetworks from Fig. 4d. Uninfected MEF controls and day 2–12 Ascl1-infected cells were assayed for all selected TRs, and day 22 Ascl1-alone and BAM-infected cells were additionally assayed for maturation TRs. Cells for day 5 to day 22 samples were FAC-sorted into Tau–eGFP⁺ and Tau–eGFP⁻ populations (n = 4 for all populations, biological replicates; error bars, s.e.m.). g, Western blot for selected TRs from the initiation subnetwork presented in panel b. β -Actin was used as a loading control (Supplementary Data 8).