



Reconstructing differentiation networks and their regulation from time series single cell expression data

Jun Ding, Bruce Aronow, Naftali Kaminski, et al.

Genome Res. published online January 9, 2018

Access the most recent version at doi:[10.1101/gr.225979.117](https://doi.org/10.1101/gr.225979.117)

P<P	Published online January 9, 2018 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Creative Commons License	This article is distributed exclusively by Cold Spring Harbor Laboratory Press for the first six months after the full-issue publication date (see http://genome.cshlp.org/site/misc/terms.xhtml). After six months, it is available under a Creative Commons License (Attribution-NonCommercial 4.0 International), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .

Advance online articles have been peer reviewed and accepted for publication but have not yet appeared in the paper journal (edited, typeset versions may be posted when available prior to final publication). Advance online articles are citable and establish publication priority; they are indexed by PubMed from initial publication. Citations to Advance online articles must include the digital object identifier (DOIs) and date of initial publication.

To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

Reconstructing differentiation networks and their regulation from time series single cell expression data

Jun Ding¹, Bruce J. Aronow², Naftali Kaminski³, Joseph Kitzmiller⁴, Jeffrey A. Whitsett⁴,
and Ziv Bar-Joseph ^{*1}

¹Computational Biology Department, School of Computer Science, Carnegie Mellon University, Pittsburgh, PA 15213, USA.

²Division of Biomedical Informatics, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA.

³Section of Pulmonary, Critical Care and Sleep Medicine, Yale School of Medicine, New Haven, CT 06520, USA.

⁴Section of Neonatology, Perinatal and Pulmonary Biology, Perinatal Institute, Cincinnati Children's Hospital Medical Center, Cincinnati, OH 45229, USA.

Keywords: Kalman filter, Organogenesis, Regulatory differentiation networks, Single cell RNA-seq, Time-series.

Running Title: Developmental trajectory reconstruct

*To whom correspondence should be addressed. Ziv Bar-Joseph, Email: zivbj@cs.cmu.edu

Abstract

Generating detailed and accurate organogenesis models using single cell RNA-seq data remains a major challenge. Current methods have relied primarily on the assumption that descendant cells are similar to their parents in terms of gene expression levels. These assumptions do not always hold for in-vivo studies which often include infrequently sampled, un-synchronized and diverse cell populations. Thus, additional information may be needed to determine the correct ordering and branching of progenitor cells and the set of transcription factors (TFs) that are active during advancing stages of organogenesis. To enable such modeling we have developed a method that learns a probabilistic model which integrates expression similarity with regulatory information to reconstruct the dynamic developmental cell trajectories. When applied to mouse lung developmental data the method accurately distinguished different cell types and lineages. Existing and new experimental data validated the ability of the method to identify key regulators of cell fate.

Supplemental Website: <http://www.cs.cmu.edu/~jund/scdiff/>.

Introduction

Most methods for reconstructing regulatory networks using high throughput data relied on microarray and RNA-seq studies profiling large populations of cells (Ernst et al., 2007; Schulz et al., 2012; Margolin et al., 2006; Liao et al., 2003). While such approaches have led to many important results, they tend to overlook the heterogeneity of the population being profiled. This may be problematic where a mixture of different cell types, with different regulatory programs, is being profiled, for example in cancer (Dalerba et al., 2011), immune response (Shalek et al., 2013) and development (Treutlein et al., 2014).

Single cell RNA-seq data addresses this problem by profiling the contribution of different cell types to changes in tissue level expression, allowing for much more detailed and accurate models. However, such data has also raised new computational challenges, some of which were recently addressed including issues related to sample quality (Stegle et al., 2015), normalization of single cell data (which is more challenging, especially for lowly expressed genes (Wu et al., 2014; Shapiro et al., 2013)), and the development of clustering methods to identify distinct components within a specific mixture / time point (Buettner et al., 2015; Guo et al., 2016).

Another challenge with single cell RNA-seq data is the analysis of time series. While several methods have been developed for the analysis and modeling of temporal data in population based microarray and RNA-seq experiments (Bonneau et al., 2006; Bar-Joseph et al., 2012; Patil and Nakai, 2014; Young et al., 2014), they all relied on one key assumption: that consecutive time points measure a continuously evolving process. In other words, the assumption is that measurements at time point $t + 1$ are correlated with measurements at the previous time point t (either the $t + 1$ expression levels continuously evolve from the expression of the same genes at time point t (Bar-Joseph et al., 2003) or they are regulated by genes expressed at the previous time point (Bar-Joseph et al., 2012)). While these assumptions may hold for the population as a whole, it clearly does not hold for all individual cells whose functions, proliferation and differentiation vary dynamically within the population. Thus, a key issue when analyzing single cell RNA-seq data is the ability to not only identify different cells within a specific time point (for example, by clustering (Xu and Su, 2015)) but also to link these cells over time by identifying the subsets of cells that belong to the same trajectory. A further challenge is to derive the regulatory networks

that control different cell fates or states that are profiled in the study.

A few recent methods have been developed to address the problem of connecting single cells along a temporal trajectory. Some of these methods are limited and can only reconstruct models with no branching (a single trajectory) (Bendall et al., 2014) or with a single branch point (Setty et al., 2016). While these may be useful for *in vitro* data, they are less appropriate for *in vivo* studies in which multiple types of cells are studied (Treutlein et al., 2014). Other methods either completely ignore the time at which the cell was measured (Trapnell et al., 2014) or rely on the measurement time (Treutlein et al., 2014; Marco et al., 2014), ignoring the fact that cells may be in different developmental states at a single time point. Indeed, both types of methods cannot accurately reconstruct complex developmental trajectories (Rashid et al., 2017) and fail to distinguish between differentiated and undifferentiated cells at a specific time point. While these methods differ in the computational models they use, they generally rely on the same underlying assumption that consecutive cells (or states) in the ordering should be very similar in terms of expression levels of their genes. While this assumption makes sense when sampling rates are very high, they do not always hold for *in vivo* studies (for example, the lung developmental data discussed in this paper which is sampled every two days). In such cases additional information can be used to determine the ordering and branching in the model. One such source of information is the set of transcription factors (TFs) that are active at each developmental stage. If these can be inferred, then states in which the factors are active could be linked to downstream states in which their targets are activated or repressed even if the overall correlation between the two states is not very high. An advantage of such an approach is that in addition to the ordering or branching model we also obtain a network model that describes which regulatory events lead to the different cell fates, the TFs controlling these events and their time of activation.

In the present work, we have utilized single cell RNA expression data during the critical period of lung morphogenesis as the embryo prepares for air-breathing at birth. At this time the lungs consist of many distinct, mesenchymal and endodermally derived epithelial cells that are rapidly dividing and differentiating to form a functional organ. During late gestation, dramatic changes in organ structure and epithelial cell differentiation and function creates a functional gas exchange unit in the alveolar regions of the lung via a process that is highly active, but still not fully understood at the molecular level (Whitsett and Weaver, 2015).

To model the process of lung epithelial growth and differentiation, we present a model that integrates time series single cell RNA-seq data with general protein-DNA interaction data. We applied our model to reconstruct differentiation networks and their regulation based on single cell lung development data. The model accurately distinguishes between cell types and trajectory of branching during cell differentiation. The model predicts several TFs as important regulators at various stages of development. While many of these were known others are novel. We used existing and new data to validate some of these TFs and their activation times.

Results

We developed a new method which learns a probabilistic model for constructing regulatory networks from single cell time series expression data. An overview of the method is presented in Figure 1 (see also Supplemental Figure 1 for an illustration using real data). We initially start by clustering cells within each time point. We use several clustering evaluation metrics to determine the number of clusters at each time point (Methods). While the measured time provides useful information about the state of the different cells, previous studies demonstrated that cells from the same time point can be unsynchronized. To address this issue we allow for cells to be moved to states representing other time points than the ones they were measured in (see below) and we further test each of the clusters to determine whether certain clusters in the same time point are actually composed of cells at different differentiation stages (Methods). Following this analysis, we arrive at an initial model in which we link each cluster to the cluster at the preceding time point that is most similar to it in expression space (Fig 1c). We next iterate between two steps: Reassigning cells to states in the model and determining the set of states and their connectivity (parent - child relationship) until the likelihood does not increase. As part of the model learning, we determine the set of TFs predicted to regulate genes in each of the states. We use this information to improve our model by requiring that factors regulating descendant states be expressed at parental states and by ensuring that genes expressed or repressed in cells assigned to descendant cells that are regulated by the identified TF follow their predicted trajectory (up or down regulation). Both requirements further impact cell assignment and model learning. If, after reassignment, states become empty (no assigned cells) they are removed from the model. Thus, unlike all prior methods for determining trajectories in single cell time series data, our model does not only rely on expression similarity but also takes into account potential regulation which may be important for in-vivo studies in which sampling is not frequent.

Application of the model to lung developmental data

To test our method we first applied it to study cell fate trajectories in mouse lung development. We used a time series dataset with 152 cells from (Treutlein et al., 2014). Of these, 45 cells were profiled at day E14.5, 27 cells at E16.5 and 80 cells at E18.5. Known cell fate markers were used to

determine the cell type for the E18.5 cells (earlier cells are progenitors and may not express these markers). We applied our method to these data and observed that it converged after 5 iterations. While convergence is obviously data dependent, we note that we observed similarly fast convergence when analyzing other single cell datasets (for example, 6 and 8 iterations for different sets of mouse embryonic fibroblasts reprogramming datasets) indicating that the method can likely be widely applied.

As can be seen in Figure 2, the model correctly separated all 4 terminal cell types (AT1, AT2, Club and ciliated) into different terminal states. It identified an additional terminal state (E2_18_0, state numbers are arbitrary) that contains a mixture of AT1 cells and bi-potential progenitors (BP) cells. These latter cells were first identified by Treutlein et al. and predicted to be non-terminal fates that can serve as progenitors to both AT1 and AT2 cells. Treutlein et al. have also observed that intermediate states are present at E18.5 such as Early AT1 (*Pdpr*, *Ager* Positive, *Sftpc* Low) and Early AT2 cells (*Sftpc* positive, *Pdpr* and *Ager* low). This is captured in our model. There are two AT1 states in our model, the first (E2_18_0) contained BP cells and the second (E3_18_4) is more homogeneous. Average expression of *Sftpc* (an AT2 marker) is 7.81 in the AT1 cells from first state and only 4.72 in the AT1 cells from the second (p-value difference of 0.0107 based on rank test) suggesting that the first AT1 group represents an early AT1 group while the second is a more differentiated state. The model also correctly reconstructs parts of the known branching process indicating that ciliated cells are derived from a different set of progenitors at E16.5 (Rawlins et al., 2007).

In addition to the assignment of cells to states, the model also highlights several TFs as playing an important role in cell differentiation. Several of these factors are known to be involved in this process including NKX2-1 (Herriges and Morrisey, 2014), SOX9 (Rockich et al., 2013), FOXA1/FOXA2 (Wan et al., 2004) and GATA6 (Yang et al., 2002).

Increasing the number of E16.5 cells in our model

While the model in Figure 2 agrees with many known aspects of lung cell differentiation, it does not provide enough information about the less studied parts of this process, specifically the roles TF play in driving cells to different fates. Such understanding is a key point since it can provide information about why certain types of cells are absent from diseased lungs and may even provide

directions for treatments in such cases. One of the problems is the fact that relatively few cells were profiled for the intermediate time point (only 27 cells at E16.5). Thus, to improve the model we replaced the E16.5 cells with epithelial cells from a larger study that only focused on E16.5 single cells in lung development (Du et al., 2015). 49 of these cells are epithelial progenitors and were further associated with various potential fates based on marker expression profiles (Guo et al., 2015)(note that these assignments were not used in model learning which is unsupervised but will be discussed below when analyzing the results). Since the data now comes from two different groups, we performed an analysis, using housekeeping genes, to determine if further normalization was needed (Supplemental Methods, Supplemental Figure 2).

A detailed view of epithelial cell differentiation in lung development

Figure 3 presents the revised model using the Du et al. E16.5 data. As can be seen, while the assignment to terminal states in this model is similar to the one in Figure 2, we see differences in the overall structure with a more detailed view of the differentiation process. For example, in this model, we see an earlier separation of ciliated and Club cells, as has previously been observed (Rawlins et al., 2007). In contrast, the separation of AT1 and AT2 cells is through a set of progenitors and their fate is determined later in the process. Once again we see BP cells mainly clustered with AT1 cells (though this time also as progenitors to these cells, for example the link from the third to the fourth level in the model). The cell identities used to evaluate the cell assignment were obtained from the original studies (Treutlein et al., 2014; Guo et al., 2015).

Given the better agreement with prior knowledge along with the more elaborate view, we have studied this model in more details. To validate some of these predicted TFs, we performed Gene Ontology analysis using PANTHER Version 11.0 (Mi et al., 2013). We found that TFs predicted by the model were significantly enriched for GO terms associated with lung epithelial cell differentiation (1.75×10^{-5}) and regulation of epithelial cell proliferation (2.0×10^{-6}) (Supplemental Table 1). While many of the predicted TFs are novel (see also below), several are supported by prior studies. For example, E2F4 is required for the development of ciliated and Club cells (Danielian et al., 2007). We found E2F4 to be ranked as one of the top TFs (1st and 2nd) for the ciliated and Club states. SREBP1 (SREBF1) is required for the development of alveolar epithelial cells (Mason, 2006), consistent with model predictions (E2_16.0 \rightarrow E3_18.0 (AT1,BP mixture)). Similarly, the

model identifies CEBPA/CEBPB as regulating the development of alveolar and airway epithelial cells, an observation that is supported by prior work (Martis et al., 2006; Roos et al., 2012). See Supplemental Table 2 for a complete list of known TFs identified by our model.

Simulation and robustness analysis

We performed simulation studies to test the ability of our method to handle noise and dropouts, which are often prevalent in single cell data (Kharchenko et al., 2014). For each cell, we simulated different dropout rates by setting the expression of randomly chosen 5% -80% genes to zero. Results indicate that our method is robust against such noise. For 5% and 10% simulated dropouts, the predicted differentiation structures is the same as the one presented above. Cell assignments are only slightly worse but generally agree the ones obtained using the original data without simulated dropouts (Supplemental Table 3). When the dropout rate increases to 20%, the predicted differentiation structure changes slightly (Supplemental Figure 3), ciliated cells and Club cells were assigned to the same cluster and AT1 cells were assigned to 3 different clusters, while the overall cell assignment is still in good agreement with the predictions on the original data and also with the known labels. Beyond 40% dropout, we see more changes though scdiff is still able to separate proliferative and non-cycling AT1/AT2 precursors and AT1 and AT2 cells. See Supplemental Figure 3 4 5 6 7 for more details.

To simulate the expression noise, we also added varying levels of random Gaussian noise to the expression of all genes (Supplemental Results). Again, we observe that for low noise levels, the predicted differentiation structures is the same as the one of the original data while for higher levels the structure is only slightly different (Supplemental Table 4, Supplemental Figure 8). We also performed a bootstrap analysis in which we used a subset of the cells to learn the model (randomly sampling 80%, 82.5%, 85%, 87.5% and 90% of all cells). We compared the resulting models (Supplemental Figures 9 10 11 12 13) and observed that both the models and cell assignments are similar to the ones obtained when using the full set of cells (roughly 90% agreement for cell assignment, Supplemental Table 5). We also tested the impact of some of the parameters used by the model and observed that within a reasonable range the changes did not have a large impact on the resulting model (Supplemental Figures 14 and 15).

Comparisons to prior methods

While pseudo time ordering methods differ from scdiff in several aspects (including the use of the profiled time for the initial assignment and the ability to infer continuous vs. discrete ordering), both types of methods attempt to infer models for the progression of cell states in developmental studies. We have thus compared scdiff to pseudo-time ordering methods. Past work has shown that some of these methods, including Monocle (Trapnell et al., 2014), SCUBA (Marco et al., 2014) and Principle component analysis (PCA) fail to accurately model cell assignment and trajectories for the lung development data discussed above (Rashid et al., 2017). Here we further analyze the performance of another method, Diffusion Pseudotime (DPT) (Haghverdi et al., 2016) on the lung (Treutlein et al., 2014) and on mouse embryonic fibroblasts reprogramming data (Treutlein et al., 2016). As can be seen in Figure 4 while DPT finds some structure in both datasets, it fails to correctly separate cell types and identify branching for the lung data and does not accurately order the reprogramming data. In contrast our method is able to both, correctly assign cells to states and identify the progression of time from embryonic cells to developed neurons. See also Supplemental Figure 16 for comparison using bone marrow data (Olsson et al., 2016).

In addition to direct comparisons that focus on the ordering and cell assignments (which is the focus of all prior methods including TASIC (Rashid et al., 2017)) these prior methods do not use protein-DNA interaction data and so cannot directly identify the set of TFs that regulate each branching point. Thus, a major advantage of our method is the ability to rely on such data to infer not just cell assignment but also the regulatory events that drive this process. To assess the impact of this novel aspect of our method, we have applied the method without using TF information (i.e. similar to prior methods which only use expression similarity). Results are presented in Supplemental Figure 17. As can be seen cell assignments and enriched TFs (based on their targets) were different when not using the TF-gene interaction information to construct the model. Six terminal cells (7.5%) are assigned to an incorrect state in this case. We also see differences in the set of significant TFs (calculated as a post processing step when not using them for the learning). Specifically, several TFs that are known to be involved in epithelial lung development are missing from the non-TF model. These include Foxa1 which regulates the lung epithelial differentiation (Besnard et al., 2005), GATA6 which regulates the differentiation of distal

lung epithelium (Yang et al., 2002), RFX3 which affects the airway epithelium development (Didon et al., 2013), and others.

Staining experiments agree with predicted TF activity time

To test model predictions for the activity of TFs we used staining experiments in developing mouse lungs. We selected a number of factors that were either novel predictions or for which the prediction of their regulatory timing was novel. These include over expression of the hypoxia-inducible factors (HIF1A), which has been identified in a variety of developmental, physiologic, and pathogenic processes within the lung (Shimoda and Semenza, 2011; Tibboel et al., 2015), SOX9 which has multiple roles in the lung epithelium, including the regulation the extracellular matrix (Rockich et al., 2013) and known epithelial cell markers.

HIF1A is predicted to be regulating AT1 and AT2 states and its targets are predicted to be down regulated in these states indicating that over-expression (OE) of *Hif1a* at E18.5 may affect AT1/AT2 cell differentiation or function. A modest sacculation defect, increased dilation and regional difference in the SFTPC (AT2 marker) and HOPX (AT1 marker) distribution were observed in the staining results (Figure 5 a,b) at E18.5, consistent with model predictions.

SOX9 (Supplemental Figure 18) was highly expressed in peripheral regions of proliferating AEC progenitor cells at E16.5. SOX9 staining decreased dramatically by E18.5 matching our prediction. In our model, SOX9 is predicted to have a relatively high activity in edges from states at E16.5 and also the predicted expression of *Sox9* is highest in early proliferative progenitor states, consistent with the loss of SOX9 staining at E18.5.

Perturbation experiments support model predictions

While the overall model structure provides some insights about the process of epithelial cell differentiation, an important advantage of TF based assignment is the ability of the model to make specific predictions about possible perturbation experiments and their outcome. Specifically, if a TF is predicted to regulate a specific path in the model (for example, the edge from E1_16_0 to E2_16_0) but not the other fates that descend from the same parent state, then a possible prediction is that the Knock-out (KO) or over-expression (OE) of that TF (depending on the impact the TF has on its down stream genes) would impact the specific path it regulates and the

cell fates associated with it, but much less so for other fates. We have thus collected available KO and ChIP-CHIP data for 3 TFs identified in our model and compared the results to the model for the specific TFs analyzed. For each such expression experiment, we compare the correlation of the WT and KO Differentially Expressed (DE) genes to the average expression profile for those genes in each of states. For each state in our model, the KO correlation can either be similar to the WT correlation, in which case we cannot infer a large impact of the TF, or be different than the WT correlation, in which case we can infer that the TF is impacting the expression of genes in the state/cell sub-type.

The first TF we looked at was CREB1, which was predicted to regulate both the AT1 (E2_16.4 → E3_18.0 (AT1 28, BP 8)) and AT2 (E2_16.4 → E3_18.2 (AT2 11, BP 4)) edges. It was also found to be regulating ciliated (E1_16.1 → E2_18.4) and Club (E1_16.1 → E3_18.1) cells. We used WT and KO data for *Creb1* from an experiment profiling lung epithelial cells at E17.5 from (Bird et al., 2011). We identified 273 Differentially Expressed (DE) genes between the *Creb1* Knock-out (KO) and Wild-type (WT) samples and then calculated the correlation between the expression of these DE genes in the WT/KO *Creb1* study and the expression of the predicted states in our model. As predicted by the model, and as can be seen in the first three columns of Table1, *Creb1* KO had the strongest impact on AT1 and AT2 gene expression. Specifically, while WT *Creb1* data exhibited a highly significant correlation with AT1 and AT2 cells (correlation coefficient = 0.401), a KO of *Creb1* led to much lower correlation (correlation coefficient = 0.124). This suggests that CREB1 may indeed be required for the differentiation of AT1 and AT2 cells. These results were supported by (Antony et al., 2016; Besnard et al., 2011), demonstrating a severe lack of AEC (Alveolar epithelial cells) in *Creb1* deleted mice. Similarly, WT *Creb1* data shows a strong correlation with Club cells (correlation coefficient=0.352) while the correlation with the KO *Creb1* experiment is much lower. We have observed much weaker correlation between ciliated cells and *Creb1* WT data and no correlation with the *Creb1* KO data.

HMGA2 was identified as a TF required for proper cell differentiation in our model. While it is known to be involved in lung cancer (Di Cello et al., 2008), its role in lung development is much less clear. HMGA2 was predicted to regulate the Proliferative Bi-potential Precursor edge (E0_14.0 → E1_16.1 (Proliferative Bi-potential Precursor 9)) and its descendent AT1 edge (E2_16.4 → E3_18.0 (AT1 28, BP 8)) and AT2 (E2_16.4 → E3_18.2 (AT2 11, BP 4)). We looked at

Hmga2 KO experiments performed at E18.5 from (Singh et al., 2015). We identified a set of 298 differentially expressed (DE) genes. WT *Hmga2* expression levels are highly correlated with AT1 and AT2 states (correlation coefficient 0.509 for AT1 and 0.440 for AT2) Table 2. This correlation disappears for the KO *Hmga2* experiment, which supports the model predictions. The model also predicts HMGA2 as a regulator of Club cell differentiation (E1.16.1, the direct parent of the Club cells state). We did not observe an impact of *Hmga2* gene deletion on the correlation with ciliated cells.

NKX2-1 is a critical factor regulating lung epithelial differentiation (Minoo et al., 1995). Our model predicts NKX2-1 (TTF1) to be active at the early stage of lung epithelial cell differentiation. It was predicted as the regulating factor for edge E0.14.0 (Common Ancestor) \rightarrow E1.16.2 (Proliferative AT2 Early Precursor) and edge E1.16.1 (Proliferative Bi-potential Precursor) \rightarrow E2.16.4 (Non-cycling AT1 Precursor and Non-cycling AT2 Precursor). In order to validate this prediction, we downloaded ChIP-ChIP experiment for NKX2-1 performed in lung epithelial cells from (Tagne et al., 2012). We compared DE genes in each state (defined as genes whose expression in the descendant state is different from their expression in the parent state) the observed targets of NKX2-1 from the ChIP-ChIP experiment using Hyper-geometric test (Supplemental Table 6). The experimental results match the model well. Edges predicted to encode active NKX2-1 TF are much more enriched for targets of NKX2-1 (for example, p-value of 0.007 for the edge from E0.14.0 to the E1.16.2 node based on hyper-geometric distribution). In contrast, several of the other edges, which were not predicted to be regulated by NKX2-1, do not overlap with targets indicating that the model can discriminate between active factors for specific fates.

Staining and OE experiments further support model predictions

We performed staining experiments in developing mouse lungs to see if factors identified by the model based on expression and regulation are indeed active at the protein level at time predicted. For this we looked at SOX9 which has multiple roles in the lung epithelium, including the regulation the extracellular matrix (Rockich et al., 2013) and at the hypoxia-inducible factor (HIF1A), which has been identified in a variety of developmental, physiologic, and pathogenic processes within the lung (Shimoda and Semenza, 2011; Tibboel et al., 2015)

SOX9 (Supplemental Figure 18) was highly expressed in peripheral regions of proliferating AEC

progenitor cells at E16.5. SOX9 staining decreased dramatically by E18.5 similar to its assignment in the model and its expression in these points. In our model, SOX9 is predicted to have a relatively high activity in edges from states at E16.5, consistent with the loss of SOX9 staining at E18.5.

HIF1A is predicted to be regulating AT1 and AT2 states (Figure 5c) and its targets are predicted to be down regulated in these states indicating that over-expression (OE) of *Hif1a* at E18.5 may affect AT1/AT2 cell differentiation or function. HIF1A and its targets are down regulated at later stages of the model (Figure 5 d). A modest sacculation defect, increased dilation and regional differences in SFTPC (AT2 marker) and HOPX (AT1 marker) distribution were observed in the staining results (Figure 5 a,b). As predicted, increased activity (OE) of *Hif1a* disrupted sacculation and impaired epithelial cell differentiation, consistent with model predictions.

Given the staining results obtained for HIF1A, we performed additional experiments to test the impact of its expression on down stream genes. As mentioned above, HIF1A was predicted as a regulator for both AT1 edge (E2_16.4 \rightarrow E3_18.0 (AT1 28, BP 8)) and AT2 edge (E2_16.4 \rightarrow E3_18.2 (AT2 11, BP 4)). However, unlike the other TFs mentioned above, we observed a decline in the expression levels of *Hif1a* at target states indicating that the activity of this TF activator needs to be reduced during lung development (Figure 5). Following (Bridges et al., 2012) we used a cDNA construct that constitutively activates *Hif1a* in normoxic conditions. We compared two versions of over expressed (OE) *Hif1a*, Single Transgenic Samples -STG and Double Transgenic Samples-DTS. We identified 223 DE genes between STG and DTG and used this to examine the correlation between states in our model and *Hif1a* OE. The results are presented in Table 3.

Our results support the role of HIF1A as a regulator of (repressed) lung development. Although the over-expression results for *Hif1a* (Table 3) did not show significant correlation difference between STG and DTG mice samples, the staining experiments (Figure 5) as well as the direction of correlation coefficient change in Table 3 support the model prediction.

Analyzing additional datasets

To further test if our method can be generally applied to analyze progression pathways from single cell RNA-seq data we have also used it to analyze time series single cell datasets from mouse embryonic fibroblasts reprogramming (Treutlein et al., 2016) and from mouse bone marrow (Olsson et al., 2016). The mouse embryonic fibroblasts reprogramming dataset focused on two settings, the

first studied cells treated by ASCL1 and the second looked at cells treated with a combination of ASCL1, POU3F2 (previously known as BRN2) and MYT1L. Our model identified ASCL1 as a key regulator for both conditions even though the information about the specific gene perturbation experiment was not used in the learning process. Several other known factors were identified. The model accurately assigned cells to states (Supplemental Figure 19 and Figure 4) and provided a map for the differentiation of mouse embryonic fibroblasts (MEFs) to multiple cell fates (see Supplemental website for interactive model). Similarly, for the mouse bone marrow data (Olsson et al., 2016) our predicted model correctly determines that HSCP cells differentiate to Mono and Gran cells through a series of intermediate states. See the Supplemental Results and Supplemental Figures 20 and 21 for more details about performance comparison on additional datasets.

Discussion

We developed and tested a computational method for reconstructing dynamic regulatory networks from single cell time series data. Unlike prior methods for pseudo-temporal ordering of such data, our method uses static information about targets of TFs to improve both the learning of a branching model and to identify TFs that regulate various stages in the process. Applying our method to single cell lung development data from multiple laboratories allowed us to reconstruct developmental pathways for a number of different types of lung epithelial cells. As we show, the reconstructed models both capture known biology (in terms of cell groupings and temporal assignment of events) and raise new hypotheses about the roles that certain TFs play in the development of specific cell types. We validated these predictions using both immunofluorescence staining and expression experiments identifying new roles for a number TFs in regulating lung development.

One of the predicted TFs, HIF1A is known to decrease in expression with advancing gestation in the fetal mouse lung (Bridges et al., 2012). To assess the effects of HIF1A on epithelial cell differentiation, an oxygen-stable form of HIF1A was conditionally expressed in respiratory epithelial cells. As predicted, OE of *Hif1a* inhibited maturation of AT1 cell precursors, indicated by decreased intensity and numbers of HOPX stained AT1 cells, and increased proportion of proSP-C stained AT2 cells. Additional support to the model was obtained using immunofluorescence confocal microscopy analysis of fetal mouse lung from the canalicular (E16.5) to saccular stage (E18.5) of lung morphogenesis. At E16.5, lung mesenchyme was prominent and epithelial cells lining peripheral regions of acinar buds stained for both NKX2-1 and SOX9, as predicted by the model (Supplemental Figure 18). At E18.5 the peripheral acinar buds had dilated, mesenchyme thinned, the levels SOX9 were markedly decreased, and HOPX increased, consistent with differentiation of AT1 and AT2 cell progenitors. Phosphohistone H3, a marker of cell proliferation, expressed in the SOX9 positive epithelial progenitors and associated mesenchyme at E16.5, was markedly decreased at E18.5, consistent with the decreased proliferation that occurs with advancing gestation in the mouse lung.

The limited knowledge of TF-DNA interaction is one bottleneck of our method. For example, we did not have the targets for HOPX in our database and thus were unable to predict it as an active regulator. In order to overcome this problem, our method provides the ability to predict top DE

genes for each edge. By combining predicted regulators and top DE genes information, our model is able to identify potential regulators even without accurate target information. For example, the aforementioned missing regulator *Hopx* was predicted as top DE genes (top up-regulated DE genes in AT1 paths and top down-regulated DE gene in AT2 paths), which is consistent with the fact that *Hopx* is an AT1 marker.

To further test if our method can be generally applied to analyze progression pathways from single cell RNA-seq data we have also used it to analyze single cell RNA-seq from mouse embryonic fibroblasts reprogramming data (Treutlein et al., 2016) and time series mouse bone marrow (Olsson et al., 2016). The reconstructed models in both cases agreed with known biology while highlighting several novel TFs as potential regulators. These results highlight the global applicability of the method which we hope can be used to study a wide range of developmental and differentiation processes.

Materials and methods

Single cell RNA-seq datasets

We downloaded time series lung single cell data from (Treutlein et al., 2014) and mouse embryonic fibroblasts reprogramming single cell data from (Treutlein et al., 2016). We also used lung single cell E16.5 data from (Du et al., 2015; Guo et al., 2015). We pre-processed these datasets as was done in the original paper (Treutlein et al., 2014). Specifically, (1) If the FPKM of gene expression is smaller than 1, the gene will be regarded as not expressed. (2) Genes with zero variance across cells are removed. We also tried a more stringent criterion, please refer to the Supplemental Result and Supplemental Figure 22 for details. (3) Transform to Log FKPM.

Initial clustering of single cells

We start by clustering the single cells at each individual time point to get an initial cell assignment. For this, we use a correlation-based method that was shown to be more suited than Euclidean distance when dealing with noisy (and sometimes partial) data (Zimek et al., 2012). We use Spearman correlation to compute a similarity matrix across cells. Next, spectral clustering (Ng et al., 2001) is used to cluster single cells based on the similarity matrix. For larger datasets with thousand of cells, the time complexity of the spectral clustering ($O(n^3)$ where n is the number of cells) may be prohibitive. For such datasets we have also implemented an alternative initial clustering strategy: PCA+K-Means which is much faster and does not significantly impact results. See Supplemental Methods, Results and Supplemental Figures 23 24 for the complete details. To determine the number of clusters for each time point (or states in our initial model) we used several quality assessment scores. We combined these scores using an ensemble strategy similar to random forest to determine the optimal number of Clusters K for each time point. See Supplemental Methods for the discussion of methods used and how they were combined.

Reassigning clusters and initial model construction

The initial clustering was based on the time point associated with each cell in the time series experiment. However, several recent studies indicate that cells may be unsynchronized with respect to their state even if they are collected at the same time point (Trapnell et al., 2014; Goranov et al.,

2009). Thus, some of the clusters at a specific time point may represent states that are either earlier or later than other clusters in the same time. In this work, we developed and used the ‘**Similarity To Ancestor-STA**’ (STA) strategy to infer an initial cluster assignment to various levels in the model. STA computes the Spearman correlation between the expression of all cells (where the expression of cell is defined as the expression vector of all genes in the cell) within every cluster and the expression of the cluster(s) at the first time point. STA of a cluster represents a vector of Spearman correlation values between the expression of each cell within the cluster and the expression of cluster(s) at the first time point. Clusters (except for the ones belonging to the first time points) are sorted based on the average STA of the cluster. Next, we compute the significance of the difference in correlation between consecutive clusters in the ordering using ranksum test $pv = ranksums(STA_X, STA_Y)$ for a pair of clusters X and Y . If we find a point in the ordering where the difference is significant ($p - value < .05$), we assign the clusters that follow that break to a new level. This process is continued for all levels until reaching the last cluster (see also Supplemental Methods).

Once we determined the set of levels in the model and the clusters associated with each level, we next connect clusters in each level to the most similar cluster (in terms of correlation) at the level right above it. By connecting all clusters to their parents, we get a graph (clusters as Nodes, parent-child relationship as Edges) which structurally represents the differentiation model.

Predicting TFs regulating differentiation pathways

An important aspect of our method is the ability to both reconstruct and analyze the differentiation pathways based on the set of TFs that regulate various state transitions. TFs whose targets are active in later stages of the process are likely active at earlier stages (in order to activate or repress their targets) and so expression levels of TF at a specific time point can be used to determine cell assignment and state connections at the next time point. We discuss below how we use TFs to impact these aspects. Here we discuss how we identify a set of TFs that are used to seed the model and the transition and emission probabilities.

We used the TF-gene interaction data from (Schulz et al., 2012; Ernst et al., 2007). “TF-gene interaction data” refers to the information about potential targets for TFs. The data is in a form of matrix with each entry denoting a TF-gene pair. Values are either binary (yes/no evidence for the interaction or probabilities (between 0 and 1) depending on the source used to infer the

interaction. See (Schulz et al., 2012) for complete information on how the data is collected and processed. Following the initial model construction, we first identify a set of differentially expressed (DE) genes (from parent cluster to current cluster) for each cluster (state) in our model . Using this set, we identify TFs that are enriched for DE targets in each state using the hyper-geometric distribution. TFs with $p - value < 0.1$ are kept as the candidate regulators. Next, we check which of the candidate TFs is expressed in the parent node of the state (expressed in at least 20% cells of the cluster). TFs that are both significantly enriched and expressed are used in the expression progression Kalman filter model as discussed below. A ranking was also provided beside the p-value for each regulating TF to demonstrate the relative regulating power at each specific edge.

To select a subset of TFs for each of the edges in the model, we use a Lasso regression method (Tibshirani, 1996) which uses the TF-gene interaction data to predict the expression values for target genes in the down stream state. We first classify genes in that state as up-regulated \uparrow , down-regulated \downarrow or not-changing \approx comparing with the parent state (Supplemental Methods). Next, a logistic regression classifier that uses the interactions between selected TFs and the genes as input is trained with the target of maximizing the ability to predict the level of the target gene expression based on the interaction data alone. The idea behind this is that TFs that are active would be selected by the Lasso method since they provide useful information about their targets whereas those that are inactive or less significant would have very small coefficients and be removed from the model. See Supplemental Methods for complete details.

A Kalman filter model for differentiation progression

To model expression changes and regulation during single cell differentiation, we use a Kalman Filter model. Similar to Hidden Markov Model (HMMs), when using a Kalman filter we need to estimate transition and emission models, though unlike the unconstrained version in HMMs these take a specific, linear, form. Our Kalman Filter model assumes that gene expression at cluster s is related to the expression of its parent cluster P_s based on the following transition model:

$$X_s = A_s X_{P_s} + B_s + w_s \quad (1)$$

$$w_s \sim N(0, Q) \quad (2)$$

Where X_s denotes the gene expression vector at cluster s , X_{P_s} denotes gene expression of the parent cluster of s , w_s is the process noise, which is assumed to be drawn from a zero mean Gaussian noise. A is the linear transition matrix and B is the offset matrix. We set A to be the identity matrix to denote the fact that genes in descendant states are expected to be similar to genes in their parent states, as was done in prior methods for modeling temporal progression in single cell studies (Trapnell et al., 2014; Marco et al., 2014; Shin et al., 2015; Juliá et al., 2015; Bendall et al., 2014). However, unlike these prior methods, our model allows for a divergence in gene expression between parent and child states for genes that are regulated by TFs that are predicted to be active in the parent state. This is the goal of the B matrix. To encode this we use the logistic regression model discussed above. Once the model is learned we have a set of active TFs for each state. We then use these TFs and the parameters learned for them to assign a label to each gene in the descendant state. Following these assignments each gene in state s is either up-regulated \uparrow , down-regulated \downarrow , or not-regulated \approx . Note that these labels are a function of the parent and so if we reassign a state to another parent in the model (see below) they may change, allowing the model to refine assignments in cases where cell memberships change.

Next, we label gene expression changes in the descendant states as follows. If gene g is determined to be ‘up-regulated’, then its expected expression value in cluster s will be the expression of its parent cluster P_s multiplied by a up-regulation scaling factor U ($1/U$ for down regulation). If g is predicted to be not-regulated, then its expected expression is the same as the one at P_s .

To handle dropouts we use a variant of the zero-inflated negative binomial model (Kharchenko et al., 2014) which we adjust to handle continuous values (in our case, Gaussian emission distribution). Specifically, similar to zero inflated models we use a mixture model for the expression emission probability. This enables us to account for dropouts (zero’s in the expression matrix) without fully penalizing the cells when computing their likelihood of being emitted from the state.

Specifically, we set the emission probability to:

$$P(g|s) = w_g p_1(g|s) + (1 - w_g) p_2(g|s) \quad (3)$$

$$p_1(g|s) \sim N(X_s^g, \sigma_s) \quad (4)$$

$$p_2(g|s) = \begin{cases} k, & \text{if } g = 0. \\ 0, & \text{otherwise.} \end{cases} \quad (5)$$

Where X_s^g is the mean of expression of gene g in cluster s as discussed above and σ_s denotes the variance of gene g . $(1 - w_g)$ is the fraction of dropped out genes for that cluster obtained by maximum likelihood estimation (MLE) as the ratio of cells with non-zero values for that gene in the cluster. k is an arbitrary probability value which is the same for all dropped genes in all clusters. The Kalman Filter model for STA was defined similarly as the expression model described above. Please refer Supplemental Methods for complete details.

Learning parameters for the Kalman Filter model

We used the initial assignments discussed above to fit gene expression transition and emission models. Given current assignments, we can compute the MLE of the transition and emission noise variance. As mentioned above we also use the initial assignments and structure to determine parameters for the logistic regression model which in turn determine the set of values used in the B transition offset matrix for each parent-child relationship. See Supplemental Methods for complete details.

Model refinement and cell re-assignments

Once we learned the initial transition and emission parameters, we can determine the global likelihood based on the assignment of cells to different states in the model.

$$\log(\text{Likelihood}(c_1, c_2, \dots, c_n, A|M)) = \sum_{i=1}^n \log P(c_i, s_i|M) \quad (6)$$

$$= \sum_{i=1}^n [\log(P(s_i)P(c_i|s_i))] \quad (7)$$

$$= \sum_{i=1}^n [\log(P(s_i)) + \log(P(STA_{c_i}|s_i)) + \log(P(G_i|s_i))] \quad (8)$$

$$= \sum_{i=1}^n \{ \log(P(s_i)) + \log(P(STA_{c_i}|s_i)) + \sum_{g_k \in g^i} \log(P(g_k|s_i)) \} \quad (9)$$

$$\log(P(s_i)) = \log\left(\prod_{q \in Q_i} p(q|q_p)\right) \quad (10)$$

$$= \sum_{q \in Q_i: \text{path to } s_i} \log(q|q_p) \quad (11)$$

Here n is the number of all cells, A represents the current assignments of cells to states, g^i is the set of all genes for cell i (c_i) and s_i is the state to which cell i is assigned. $P(G_i|s_i)$ is the expression probability of c_i , which indicates the agreement of gene expression between c_i and the state s_i , $P(g_k|s_i)$ is the expression probability of gene g_k , which represents the probability that g_k is emitted by state s_i . $P(STA_{c_i}|s_i)$ is the time probability of c_i , which indicates the agreement of STA values between c_i and state s_i . Q_i is the path from the root node to state (node) s_i , including the root node: $P(\text{root}|\text{root}_{parent}) = P(\text{root}|None) = P(\text{root})$. $\log(q|q_p)$ modeled the transition relations and was estimated based on the current assignment: $P(q|q_p) = \frac{|C_q|}{|CP_{q_p}|}$. $|C_q|$ is number of cells at state q . CP_{q_p} denotes the number of cells, which are from all children states of q_p (parent state of q).

We next attempt to improve the likelihood of the model by refining the model structure (i.e. changing parent - descendant assignments) and reassigning cells to states in the model. To reassign cells, we compute the maximal probability for cell c_i , $P(c_i|s)$ for all states s in the model. Specifically

we find:

$$\text{Assign}(c_i) = \arg \max_s P(c_i, A|M) \quad (12)$$

$$= \arg \max_s P(c_i, s) \quad (13)$$

$$= \arg \max_s P(s)P(STA_{c_i}|s)P(G_i|s) \quad (14)$$

$$= \arg \max_s P(s)P(STA_i|s) \prod_{g \in g_i} P(g_k|s) \quad (15)$$

$$= \arg \max_s \log(P(s)) + \log(P(STA_i|s)) + \sum_{g_k \in g^i} \log(P(g_k|s)) \quad (16)$$

Note that assignment can lead to states becoming empty. If this happens these states are removed from the model. After re-assigning cells to states, we further refine the model by updating nodes (states) and edges (parent relationship). We remove states that become empty and re-compute the edges (fromNode, toNode, regulating TFs) by updating the parent for each remaining state. For this, we use the (re) assigned cells to re-compute a set of DE genes for that state, test which potential parent state in the preceding level maximizes the transition function for that state (based on the Logistic Regression model computed for the parent) and select the parent with the highest likelihood. Once a parent is assigned, we recompute the set of TFs for the edge by using the new set of DE genes for each state (if reassignment of cell changed the set).

Identifying epithelial cells in a large cohort of E16.5 lung cells

In this work we integrated data from multiple prior lung development single cell studies. One of the datasets we used was the LunGENS (Du et al., 2015; Guo et al., 2015), which profiled 49 single cells from fetal mouse lung at E16.5 using RNA sequencing of cells separated-using the Fluidigm C1.

Perturbation and imaging experiments

All mouse experiments were performed under AAALAC approved protocols reviewed at Cincinnati Children’s Hospital Medical Center (CCHMC). For immunofluorescence confocal microscopy, lung tissue from embryos (E16.5 and E18.5) was fixed in 4% PFA (PBS). Tissue was sectioned at 5 microns for paraffin and 7 microns for frozen samples. Slides were incubated with

antisera versus NKX2-1 (Catalog number: RB1231; rabbit, Seven Hills Bioreagents), SOX9 (Catalog number: AB5335; rabbit, Millipore), SFTPC (Catalog number: SC-7706; goat, Santa Cruz), HOPX (Catalog number: SC-30216; rabbit, Santa Cruz), and ACTA2 (Catalog number: A5228; mouse, Sigma-Aldrich) or phosphohistone H3 (Catalog number: SC-12927; goat Santa Cruz). Detailed methodologies are provided in the lung Image website accessible at <https://research.cchmc.org/lungimage>. Sections were imaged on a Nikon A1Rsi confocal microscope. For studies of HIF1A, tissue was obtained from fetuses (E18.5) of transgenic mice engineered to express a HIF1A mutant protein under control of the human SFTPC-rtTA promoter construct by expressing (tetO)₇/CMV/HIF1A/ODD/N803, a normoxia stable form of HIF1A. Administration of doxycycline activates expression of the transgene in fetal respiratory epithelial cells. Dams were treated with doxycycline from E16.5 until E18.5, the time of sacrifice. Doxycycline treated single transgenic and double transgenic fetuses were identified by genotyping. Imaris and Nikon Elements software was used to export images, and Adobe Photoshop used to adjust levels of fluorescence for data display.

Mouse studies

An activated form of HIF1A, HIF1A(TPM), was expressed under conditional control of the SFTPC-rtTA, (otet)₇-HIF1A TPM. Double and single transgenic littermates were compared from dams treated with doxycycline chow from E12.5 until sacrifice, as previously reported (Bridges et al., 2012). Confocal immunofluorescence microscopy was performed for ACTA2, HOPX, and SFTPC.

Fetal mouse studies: C57BL6 mice were time mated to obtain litters at E14.5, E16.5 and E18.5 for immunofluorescence staining for SFTPC (proSP-C), HOPX, ACTA2 (α SMA), SOX9, phosphohistone H3 (pHisH3) and NKX2-1 (TTF1) as described in the LungMAP data repository at www.lungmap.net.

Software availability

scdiff is primarily written in Python, available as an open source tool at GitHub (<https://github.com/phoenixding/scdiff>). This GitHub repository includes detailed instructions on how to use the method. The scdiff source code is also available as the Supplemental code. All the data and

results on the Supplemental website (<http://www.cs.cmu.edu/~jund/scdiff/>) are provided as the Supplemental Materials (Supplemental website).

Acknowledgment

This work is supported in part by National Institutes of Health [grant number U01HL122626-01 to Ziv Bar-Joseph and U01HL122642 to Jeffrey A. Whitsett] and by the National Science Foundation [grant number DBI-1356505 to Ziv Bar-Joseph].

We thank Yina Du for RNA data support and Dr. James Bridges for HIF1A(TDM) tissue. We also thank Easwaran Ramamurthy for testing our software.

References

- Antony N, McDougall A, Mantamadiotis T, Cole T, and Bird A. 2016. Creb1 regulates late stage mammalian lung development via respiratory epithelial and mesenchymal-independent mechanisms. *Scientific reports* **6**.
- Bar-Joseph Z, Gerber G, Simon I, Gifford DK, and Jaakkola TS. 2003. Comparing the continuous representation of time-series expression profiles to identify differentially expressed genes. *Proceedings of the National Academy of Sciences* **100**: 10146–10151.
- Bar-Joseph Z, Gitter A, and Simon I. 2012. Studying and modelling dynamic biological processes using time-series gene expression data. *Nature Reviews Genetics* **13**: 552–564.
- Bendall SC, Davis KL, Amir EaD, Tadmor MD, Simonds EF, Chen TJ, Shenfeld DK, Nolan GP, and Peer D. 2014. Single-cell trajectory detection uncovers progression and regulatory coordination in human b cell development. *Cell* **157**: 714–725.
- Besnard V, Wert S, Kaestner K, and Whitsett J. 2005. Stage-specific regulation of respiratory epithelial cell differentiation by foxa1. *American Journal of Physiology-Lung Cellular and Molecular Physiology* **289**: L750–L759.
- Besnard V, Wert SE, Ikegami M, Xu Y, Heffner C, Murray SA, Donahue LR, and Whitsett JA. 2011. Maternal synchronization of gestational length and lung maturation. *PLoS one* **6**: e26682.
- Bird AD, Flecknoe SJ, Tan KH, Olsson PF, Antony N, Mantamadiotis T, Hooper SB, and Cole TJ. 2011. camp response element binding protein is required for differentiation of respiratory epithelium during murine development. *PLoS One* **6**: e17843.
- Bonneau R, Reiss DJ, Shannon P, Facciotti M, Hood L, Baliga NS, and Thorsson V. 2006. The inferelator: an algorithm for learning parsimonious regulatory networks from systems-biology data sets de novo. *Genome biology* **7**: 1.
- Bridges JP, Lin S, Ikegami M, and Shannon JM. 2012. Conditional hypoxia inducible factor-1 α induction in embryonic pulmonary epithelium impairs maturation and augments lymphangiogenesis. *Developmental biology* **362**: 24–41.
- Buettner F, Natarajan KN, Casale FP, Proserpio V, Scialdone A, Theis FJ, Teichmann SA, Marioni JC, and Stegle O. 2015. Computational analysis of cell-to-cell heterogeneity in single-cell rna-sequencing data reveals hidden subpopulations of cells. *Nature biotechnology* **33**: 155–160.
- Dalerba P, Kalisky T, Sahoo D, Rajendran PS, Rothenberg ME, Leyrat AA, Sim S, Okamoto J, Johnston DM, Qian D, et al.. 2011. Single-cell dissection of transcriptional heterogeneity in human colon tumors. *Nature biotechnology* **29**: 1120–1127.
- Danielian PS, Kim CFB, Caron AM, Vasile E, Bronson RT, and Lees JA. 2007. E2f4 is required for normal development of the airway epithelium. *Developmental biology* **305**: 564–576.
- Di Cello F, Hillion J, Hristov A, Wood LJ, Mukherjee M, Schuldenfrei A, Kowalski J, Bhattacharya R, Ashfaq R, and Resar LM. 2008. Hmga2 participates in transformation in human lung cancer. *Molecular Cancer Research* **6**: 743–750.

- Didon L, Zwick RK, Chao IW, Walters MS, Wang R, Hackett NR, and Crystal RG. 2013. Rfx3 modulation of foxj1 regulation of cilia genes in the human airway epithelium. *Respiratory research* **14**: 1.
- Du Y, Guo M, Whitsett JA, and Xu Y. 2015. lunggens: a web-based tool for mapping single-cell gene expression in the developing lung. *Thorax* **70**: 1092–1094.
- Ernst J, Vainas O, Harbison CT, Simon I, and Bar-Joseph Z. 2007. Reconstructing dynamic regulatory maps. *Molecular systems biology* **3**: 74.
- Goranov AI, Cook M, Ricicova M, Ben-Ari G, Gonzalez C, Hansen C, Tyers M, and Amon A. 2009. The rate of cell growth is governed by cell cycle stage. *Genes & development* **23**: 1408–1422.
- Guo M, Bao EL, Wagner M, Whitsett JA, and Xu Y. 2016. Slice: determining cell differentiation and lineage based on single cell entropy. *Nucleic Acids Research* p. gkw1278.
- Guo M, Wang H, Potter SS, Whitsett JA, and Xu Y. 2015. Sincera: a pipeline for single-cell rna-seq profiling analysis. *PLoS Comput Biol* **11**: e1004575.
- Haghverdi L, Buettner M, Wolf FA, Buettner F, and Theis FJ. 2016. Diffusion pseudotime robustly reconstructs lineage branching. *Nature Methods* **13**: 845–848.
- Herriges M and Morrisey EE. 2014. Lung development: orchestrating the generation and regeneration of a complex organ. *Development* **141**: 502–513.
- Juliá M, Telenti A, and Rausell A. 2015. Sincell: an r/bioconductor package for statistical assessment of cell-state hierarchies from single-cell rna-seq. *Bioinformatics* **31**: 3380–3382.
- Kharchenko PV, Silberstein L, and Scadden DT. 2014. Bayesian approach to single-cell differential expression analysis. *Nature methods* **11**: 740–742.
- Liao JC, Boscolo R, Yang YL, Tran LM, Sabatti C, and Roychowdhury VP. 2003. Network component analysis: reconstruction of regulatory signals in biological systems. *Proceedings of the National Academy of Sciences* **100**: 15522–15527.
- Marco E, Karp RL, Guo G, Robson P, Hart AH, Trippa L, and Yuan GC. 2014. Bifurcation analysis of single-cell gene expression data reveals epigenetic landscape. *Proceedings of the National Academy of Sciences* **111**: E5643–E5650.
- Margolin AA, Nemenman I, Basso K, Wiggins C, Stolovitzky G, Favera RD, and Califano A. 2006. Aracne: an algorithm for the reconstruction of gene regulatory networks in a mammalian cellular context. *BMC bioinformatics* **7**: S7.
- Martis PC, Whitsett JA, Xu Y, Perl AKT, Wan H, and Ikegami M. 2006. C/ebp α is required for lung maturation at birth. *Development* **133**: 1155–1164.
- Mason RJ. 2006. Biology of alveolar type ii cells. *Respirology* **11**: S12–S15.
- Mi H, Muruganujan A, Casagrande JT, and Thomas PD. 2013. Large-scale gene function analysis with the panther classification system. *Nature protocols* **8**: 1551–1566.
- Minoo P, Hamdan H, Bu D, Warburton D, Stepanik P, et al.. 1995. Ttf-1 regulates lung epithelial morphogenesis. *Developmental biology* **172**: 694–698.

- Ng AY, Jordan MI, and Weiss Y. 2001. On spectral clustering: Analysis and an algorithm. *NIPS* **14**: 849–856.
- Olsson A, Venkatasubramanian M, Chaudhri VK, Aronow BJ, Salomonis N, Singh H, and Grimes HL. 2016. Single-cell analysis of mixed-lineage states leading to a binary cell fate choice. *Nature* .
- Patil A and Nakai K. 2014. Timexnet: Identifying active gene sub-networks using time-course gene expression profiles. *BMC systems biology* **8**: 1.
- Rashid S, Kotton D, and Bar-Joseph Z. 2017. Tasic: Determining branching models from time series single cell data. *Bioinformatics* doi:10.1093/bioinformatics/btx173.
- Rawlins EL, Ostrowski LE, Randell SH, and Hogan BL. 2007. Lung development and repair: contribution of the ciliated lineage. *Proceedings of the National Academy of Sciences* **104**: 410–417.
- Rockich BE, Hrycaj SM, Shih HP, Nagy MS, Ferguson MA, Kopp JL, Sander M, Wellik DM, and Spence JR. 2013. Sox9 plays multiple roles in the lung epithelium during branching morphogenesis. *Proceedings of the National Academy of Sciences* **110**: E4456–E4464.
- Roos AB, Berg T, Barton JL, Didon L, and Nord M. 2012. Airway epithelial cell differentiation during lung organogenesis requires *c/ebp α* and *c/ebp β* . *Developmental Dynamics* **241**: 911–923.
- Schulz MH, Devanny WE, Gitter A, Zhong S, Ernst J, and Bar-Joseph Z. 2012. Drem 2.0: Improved reconstruction of dynamic regulatory networks from time-series expression data. *BMC systems biology* **6**: 104.
- Setty M, Tadmor MD, Reich-Zeliger S, Angel O, Salame TM, 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 AK, Satija R, Adiconis X, Gertner RS, Gaublomme JT, 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**: 236–240.
- Shapiro E, Biezuner T, and Linnarsson S. 2013. Single-cell sequencing-based technologies will revolutionize whole-organism science. *Nature Reviews Genetics* **14**: 618–630.
- Shimoda LA and Semenza GL. 2011. Hif and the lung: role of hypoxia-inducible factors in pulmonary development and disease. *American journal of respiratory and critical care medicine* **183**: 152–156.
- Shin J, Berg DA, Zhu Y, Shin JY, Song J, Bonaguidi MA, Enikolopov G, Nauen DW, Christian KM, Ming GL, et al.. 2015. Single-cell rna-seq with waterfall reveals molecular cascades underlying adult neurogenesis. *Cell Stem Cell* **17**: 360–372.
- Singh I, Ozturk N, Cordero J, Mehta A, Hasan D, Cosentino C, Sebastian C, Krüger M, Looso M, Carraro G, et al.. 2015. High mobility group protein-mediated transcription requires dna damage marker γ -h2ax. *Cell research* **25**: 837–850.
- Stegle O, Teichmann SA, and Marioni JC. 2015. Computational and analytical challenges in single-cell transcriptomics. *Nature Reviews Genetics* **16**: 133–145.

- Tagne JB, Gupta S, Gower AC, Shen SS, Varma S, Lakshminarayanan M, Cao Y, Spira A, Volkert TL, and Ramirez MI. 2012. Genome-wide analyses of nkx2-1 binding to transcriptional target genes uncover novel regulatory patterns conserved in lung development and tumors. *PloS one* **7**: e29907.
- Tibboel J, Groenman FA, Selvaratnam J, Wang J, Tseu I, Huang Z, Caniggia I, Luo D, van Tuyl M, Ackerley C, et al.. 2015. Hypoxia-inducible factor-1 stimulates postnatal lung development but does not prevent o₂-induced alveolar injury. *American journal of respiratory cell and molecular biology* **52**: 448–458.
- Tibshirani R. 1996. Regression shrinkage and selection via the lasso. *Journal of the Royal Statistical Society. Series B (Methodological)* pp. 267–288.
- Trapnell C, Cacchiarelli D, Grimsby J, Pokharel P, Li S, Morse M, Lennon NJ, Livak KJ, Mikkelsen TS, and Rinn JL. 2014. The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells. *Nature biotechnology* **32**: 381–386.
- Treutlein B, Brownfield DG, Wu AR, Neff NF, Mantalas GL, Espinoza FH, Desai TJ, Krasnow MA, and Quake SR. 2014. Reconstructing lineage hierarchies of the distal lung epithelium using single-cell rna-seq. *Nature* **509**: 371–375.
- Treutlein B, Lee QY, Camp JG, Mall M, Koh W, Shariati SAM, Sim S, Neff NF, Skotheim JM, Wernig M, et al.. 2016. Dissecting direct reprogramming from fibroblast to neuron using single-cell RNA-seq. *Nature* pp. 1–15.
- Wan H, Kaestner KH, Ang SL, Ikegami M, Finkelman FD, Stahlman MT, Fulkerson PC, Rothenberg ME, and Whitsett JA. 2004. Foxa2 regulates alveolarization and goblet cell hyperplasia. *Development* **131**: 953–964.
- Whitsett JA and Weaver TE. 2015. Alveolar development and disease. *American journal of respiratory cell and molecular biology* **53**: 1–7.
- Wu AR, Neff NF, Kalisky T, Dalerba P, Treutlein B, Rothenberg ME, Mburu FM, Mantalas GL, Sim S, Clarke MF, et al.. 2014. Quantitative assessment of single-cell rna-sequencing methods. *Nature methods* **11**: 41–46.
- Xu C and Su Z. 2015. Identification of cell types from single-cell transcriptomes using a novel clustering method. *Bioinformatics* p. btv088.
- Yang H, Lu MM, Zhang L, Whitsett JA, and Morrisey EE. 2002. Gata6 regulates differentiation of distal lung epithelium. *Development* **129**: 2233–2246.
- Young WC, Raftery AE, and Yeung KY. 2014. Fast bayesian inference for gene regulatory networks using scanbma. *BMC systems biology* **8**: 1.
- Zimek A, Schubert E, and Kriegel HP. 2012. A survey on unsupervised outlier detection in high-dimensional numerical data. *Statistical Analysis and Data Mining* **5**: 363–387.

Figure legends

Figure 1: **Learning differentiation models from single cell RNA-seq data.** (a) Initial clusters are determined using Spectral Clustering. T0, T1, T2 represent the measurement time. (b) Initial ‘differentiation time’ is estimated for clusters based on difference with clusters for the first time point. DT0, DT1, DT2, DT3 denote the estimated differentiation stage. (c) Differentiation paths are constructed by connecting clusters at lower levels to their most similar parent at the level above them. (d) Regulating TFs are determined for each edge. TFs are colored based on their expression change along the edge. Increased expression: Red, Decreased expression: Blue, Stable expression: Green. Shades represent the extent of the expression change. (e) Initial model. (f) Iterating between cells and state reassignments and parameter learning until convergence. (g) Final model.

Figure 2: **Differentiation model on Treutlein et al. lung single cell dataset.** diff: Differentiation stage 0% (most undifferentiated)-100% (most differentiated) from the first state. Number of cells assigned to each state was given inside the pie chart. Cell types for each state (node) are based on the assignments of Treutlein et al. and so are only available for cells in the last time point (E18.5), cells from earlier time points are labeled using time point (E14, E16). TFs are associated with paths they are predicted to regulate and color-coded based on their expression change (fold change) along the regulated path. TF p-values are based on the set of targets associated with the states they are predicted to regulate.

Figure 3: **Differentiation model using data from both Treutlein et al. and Du et al.** Cell types taken from both, Treutlein et al. and Du et al. PT2: Proliferative AT2 Early Precursor; PT1: Proliferative AT1 Early Precursor; PB: Proliferative Bi-Potential Precursor; NT1: Noncycling AT1 Precursor; NT2: Noncycling AT2 Precursor. Differentiation scores, TFs color and p-value have the same meaning as in Figure 2.

Figure 4: **Performance comparison with Diffusion Pseudotime (DPT).** (a) Top row: DPT analysis of the Treutlein et al. mouse lung single cell data (the data used in Figure 2). Left: Cells colored by their type as determined by Treutlein et al. Right cell assignment by DPT using default parameters. While DPT finds some structure in the data it is unable to separate the AT1 and Club cells and does not show any major branching prior to E18.5 (b) Middle row: DPT analysis of mouse embryonic fibroblasts reprogramming data (Treutlein et al., 2016) setting 2. While the DPT model finds a branch leading from the Mouse Embryo Fibroblast (MEF) cells to the neurons it does not order correctly the intermediate day 2 and day 5 cells (note that d5 cells are mostly on the other branch and only d2 are close to neurons). (c) In contrast, a model based on our method for the same data (bottom) correctly places most d2 cells in the second level with d5 cells closer to the neurons. See also text for discussion.

Figure 5: Increased HIF1A activity disrupts sacculation and influences AT1/AT2 cell distribution. (a,b) Experimental results for HIF1A staining. HIF1A (Three Point mutant) was expressed under conditional control of SFTPC-rtTA, (otet)₇-HIF1A-TPM. Doxycycline was provided to the dam from E12.5 to E18.5. Single transgene (STG) controls, lacking HIF1A-TPM expression, (n=3) were compared with double transgenic (DTG) mice expressing HIF1A-TPM under doxycycline control in airway epithelial cells, n=4. Staining of lung tissue for ACTA2 (smooth muscle actin), HOPX (an AT1 cell marker) and SFTPC (proSP-C, and AT2 cell marker) are shown. (c) Model prediction for HIF1A. HIF1A is identified as a top regulator of AT1 cells (ranked as the 12th TF) with a lower, though still significant impact on AT2 cells (ranking as the 33rd TF). (d) mRNA expression of *Hif1a* in the different states reconstructed by the model. As predicted by the model, OE of *Hif1a* influences AT1/AT2 cell distribution with a larger impact on AT1 cells when compared to AT2 cells.

Tables

Table 1: Spearman correlation of DE genes between *Creb1* Knock-out data and predicted clusters

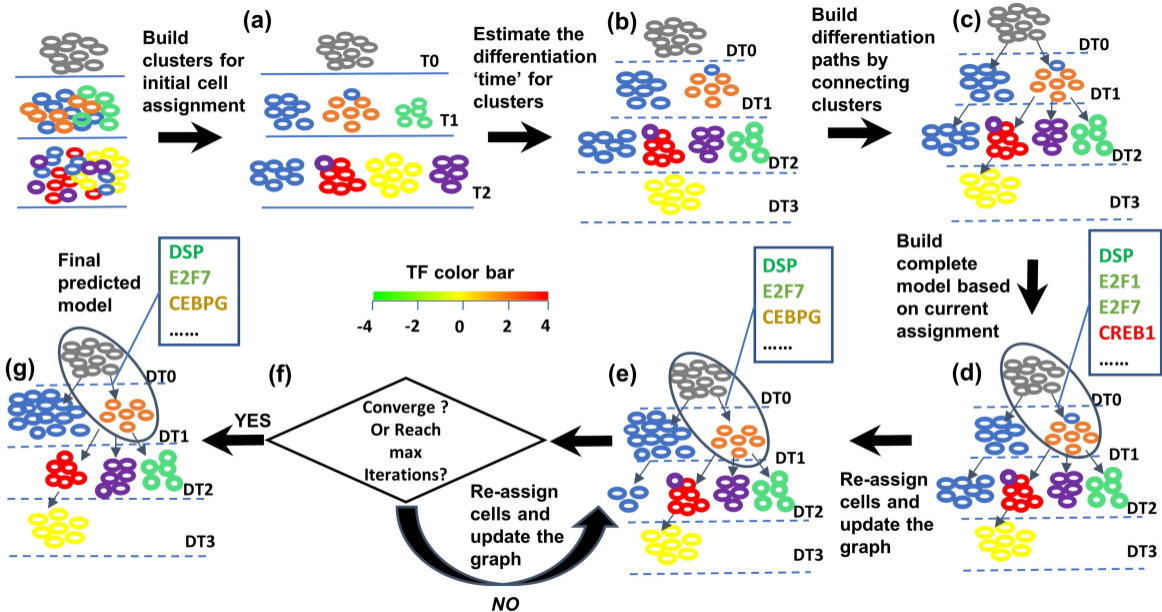
	E3.18.0 AT1(28),BP(8)	E3.18.2 AT2(11),BP(4)	E4.18.3 AT1(13),BP(1),Club(1),AT2(1)	E2.18.4 ciliated(3)	E3.18.1 Club(10)
<i>Creb1</i> WT	(0.449, 1.1×10^{-12})	(0.371, 2.6×10^{-9})	(0.383, 2.32×10^{-9})	(0.209, 0.0015)	(0.352, 4.83×10^{-8})
<i>Creb1</i> KO	(0.196, 0.00306)	(0.038, 0.569)	(0.137, 0.0390)	(0.0516, 0.439)	(0.0637, 0.339)

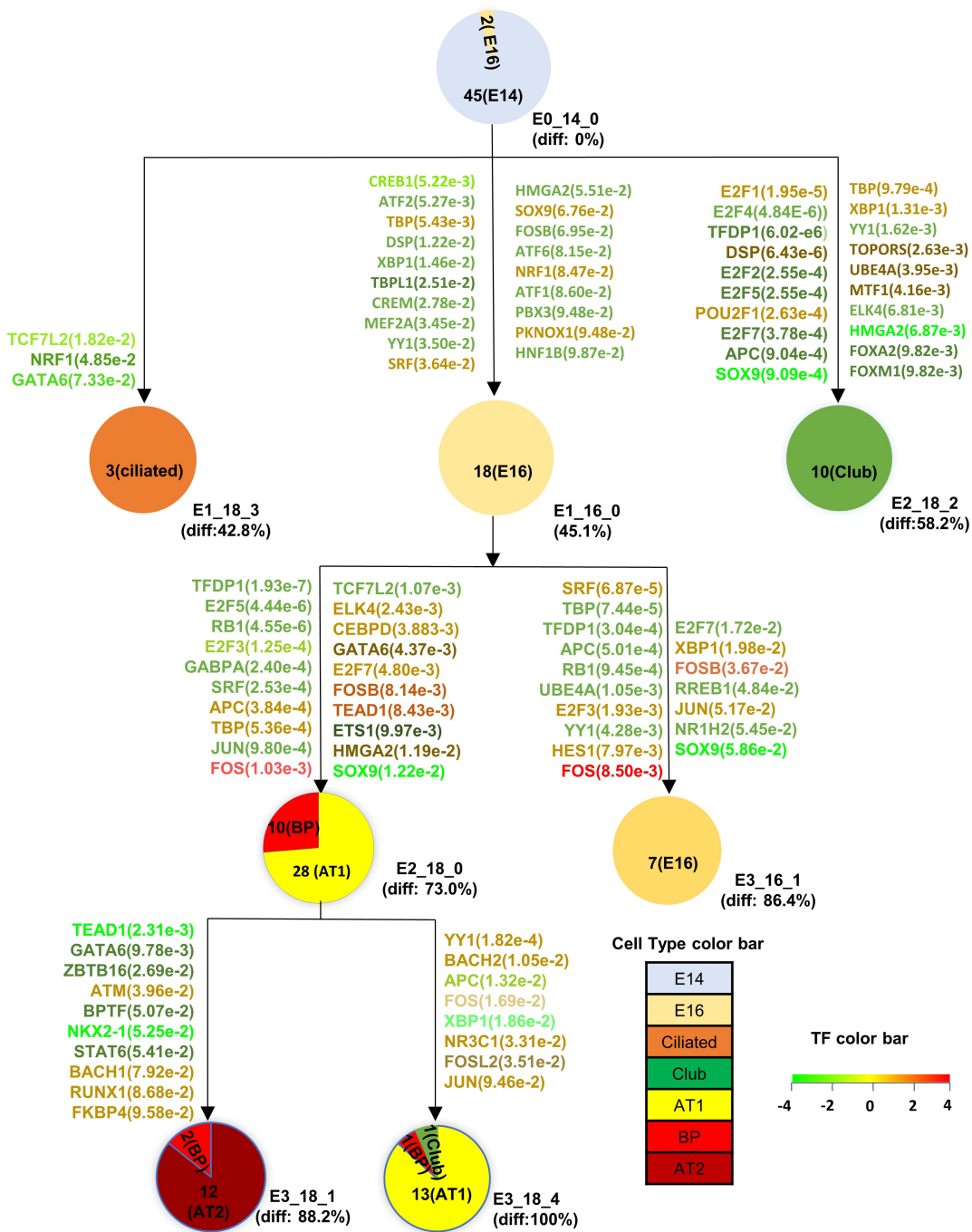
Table 2: Spearman correlation of DE genes between *Hmga2* Knock-out data and the predicted states

	E3.18.0 AT1(28),BP(8)	E3.18.2 AT2(11),BP(4)	E4.18.3 AT1(13),BP(1),Club(1),AT2(1)	E2.18.4 ciliated(3)	E3.18.1 Club(10)
<i>Hmga2</i> WT	(0.509, 4.78×10^{-21})	(0.440, 1.6×10^{-15})	(0.397, 1.1×10^{-12})	(0.176, 0.00231)	(0.441, 1.351×10^{-15})
<i>Hmga2</i> KO	(0.0889, 0.126)	(-0.0433, 0.457)	(0.0582, 0.317)	(0.167, 0.00378)	(0.0154, 0.792)

Table 3: Spearman correlation of DE genes between *Hif1a* OE experiment and predicted clusters

	E3.18.0 AT1(28),BP(8)	E3.18.2 AT2(11),BP(4)	E4.18.3 AT1(13),BP(1),Club(1),AT2(1)	E2.18.4 ciliated(3)	E3.18.1 Club(10)
<i>Hif1a</i> STG	(0.134, 0.0495)	(0.197, 0.00383)	(0.167, 0.0144)	(0.137, 0.0455)	(0.194, 0.00442)
<i>Hif1a</i> DTG	(0.04447, 0.516)	(0.0133, 0.847)	(-0.046, 0.503)	(-0.05112, 0.456)	(-0.0486, 0.479)





45(E14)
E0_14_0
(diff: 0%)

