

Stereotyped Subclones Revealed by High-Density Single-Cell Lineage Tracing Support Robust Development

Xiaoyu Zhang, Zizhang Li, Jingyu Chen, Wenjing Yang, Xingxing He, Peng Wu, Feng Chen, Ziwei Zhou, Chenze Ren, Yuyan Shan, Xiewen Wen, Vassily A. Lyubetsky, Leonid Yu. Rusin, Xiaoshu Chen, and Jian-Rong Yang*

Robust development is essential for multicellular organisms. While various mechanisms contributing to developmental robustness are identified at the subcellular level, those at the intercellular and tissue level remain underexplored. This question is approached using a well-established *in vitro* directed differentiation model recapitulating the *in vivo* development of lung progenitor cells from human embryonic stem cells. An integrated analysis of high-density cell lineage trees (CLTs) and single-cell transcriptomes of differentiating colonies enabled the resolution of known cell types and developmental hierarchies. This dataset showed little support for the contribution of transcriptional memory to developmental robustness. Nevertheless, stable terminal cell type compositions are observed among many subclones, which enhances developmental robustness because the colony can retain a relatively stable composition even if some subclones are abolished by cell death. Furthermore, it is found that many subclones are formed by sub-CLTs resembling each other in terms of both terminal cell type compositions and topological structures. The presence of stereotyped sub-CLTs constitutes a novel basis for developmental robustness. Moreover, these results suggest a unique perspective on individual cells' function in the context of stereotyped sub-CLTs, which can bridge the knowledge of the atlas of cell types and how they are organized into functional tissues.

1. Introduction

Developmental robustness, also known as canalization,^[1] refers to the phenomenon that biological development outcomes remain largely unchanged despite environmental or genetic perturbations.^[2,3] In addition to being an essential feature of complex organisms, developmental robustness also has profound implications for evolution^[4,5] and disease.^[6] Decades of studies have identified a variety of mechanisms that contribute to developmental robustness, including chaperone proteins,^[7] microRNAs,^[8–10] morphology-stabilizing genes,^[11,12] feedback loops,^[13] molecular redundancies^[14] and defect-buffering cellular plasticity.^[15] While significant advances have been made at the molecular/intracellular level, other mechanisms that ensure robust development at the intercellular/tissue levels remain poorly understood. A couple examples include the nonlinear relationship between key regulators' gene expression and embryonic structures,^[16] and the robustness to cell death observed for determinative developmental cell lineages.^[17]

X. Zhang, W. Yang, X. He, P. Wu, F. Chen, Z. Zhou, C. Ren, Y. Shan, X. Chen, J.-R. Yang
Advanced Medical Technology Center
The First Affiliated Hospital
Zhongshan School of Medicine
Sun Yat-sen University
Guangzhou 510080, China
E-mail: yangjianrong@mail.sysu.edu.cn

X. Zhang, Z. Li, J. Chen, W. Yang, P. Wu, F. Chen, Z. Zhou, C. Ren, Y. Shan, J.-R. Yang
Department of Genetics and Biomedical Informatics
Zhongshan School of Medicine
Sun Yat-sen University
Guangzhou 510080, China
X. He, X. Chen
Department of Immunology and Microbiology
Zhongshan School of Medicine
Sun Yat-sen University
Guangzhou 510080, China
X. Wen
University Research Facility in 3D Printing, & State Key Laboratory of Ultra-precision Machining Technology, Dept. of ISE
the Hong Kong Polytechnic University
Hong Kong 999077, China

 The ORCID identification number(s) for the author(s) of this article can be found under <https://doi.org/10.1002/advs.202406208>

© 2025 The Author(s). Advanced Science published by Wiley-VCH GmbH. This is an open access article under the terms of the [Creative Commons Attribution](#) License, which permits use, distribution and reproduction in any medium, provided the original work is properly cited.

DOI: 10.1002/advs.202406208

The developmental process encompasses both the history of cell divisions and state transitions.^[18,19] It is thus possible to examine development, as well as its robustness, from two perspectives. In the first, cellular states, such as single-cell transcriptomes, were recorded during various developmental stages and used to construct a continuum of states known as an epigenetic landscape^[20,21] or state manifolds.^[18] In the second, all cell divisions since the zygote or some progenitor cells can be recorded and used to construct a cell lineage tree (CLT).^[22] This CLT-based perspective, however, has been much less studied due to the difficulty in obtaining CLTs in complex organisms. Nonetheless, recent technological advancements in CLT reconstruction, particularly those utilizing genomic barcoding,^[19] have led to new opportunities for joint analyses of these two perspectives. For example, scGESTALT simultaneously determined cell states by single-cell transcriptomics and the corresponding CLT via lineage barcodes.^[23] Similar methods^[18,19] provide a combined view of single-cell states and CLTs, enabling CLT-based analyses of robustness for different developmental models.

One of the main manifestations of developmental robustness is the generation of adequate numbers of cells of various types in an appropriate cellular composition, especially when they work together as a functional unit. For example, the *Drosophila* peripheral nervous system contains thousands of identical mechanosensory bristles,^[24] each consisting of exactly one hair cell, one socket cell, one sheath cell and one neuron.^[25] Another well-known example is the functional unit of the endocrine pancreas, the islet, which has been shown in mice to consist predominantly ($\approx 90\%$) of β cells at the core and α and δ cells in the periphery.^[26] In identifying potential mechanisms that contribute to such a manifestation of developmental robustness, two CLT-based studies are particularly relevant. In the first, it was found that development of mammalian organs is preceded by significant mixing of multipotent progenitor cells.^[27] Therefore, most organs have a polyclonal origin that ensures sufficient number of cells even some progenitors failed.^[27] In the second, CLT of cortical development revealed stereotyped development giving rise to monophyletic clades of mixed cell types.^[28] On the basis of these observations, we hypothesized that the combination of polyclonal origin and stereotyped development facilitates the robust development of adequate numbers of cells with an appropriate cellular composition. It is imperative to note that as our hypothesis revolves around the above-mentioned functional units, CLTs with sufficient density (fraction of cells sampled) are essential, otherwise stereotyped development cannot be detected when

only a very small fraction of cells was sampled from each functional unit. In addition, a high resolution CLT would also reveal how stereotyped development occurs, such as mitotic-coupling versus population-coupling development^[18] and whether epigenetic memory^[29] plays a role.

To this end, we obtained the single-cell transcriptomes and high density (capturing $> 10\%$ cells in the colony) CLTs of three in vitro cell cultures that mimic the in vivo development of human embryonic stem cells (hESCs) into lung progenitors.^[30] According to a joint analysis with another in vitro culture that retained stemness, single-cell transcriptomes were clearly separated into clusters of undifferentiated and various differentiated cell types, and the CLTs showed significant signals of divergence among subclones consistent with known sequential involvement of Bmp/TGF- β , Wnt and other endoderm differentiation related pathways. Multiple monophyletic groups of cells with stable cellular compositions were revealed by this CLT, directly supporting the existence of polyclonal stereotyped development. Based on the assumption that cells work collectively as functional units composed of similar compositions of various cell types, the stereotyped polyclonal developmental programs observed produce subpopulations with properly mixed cell types, thereby ensuring the formation of more functional units in the event of random cell deaths compared to non-stereotyped development, and therefore enhances robustness. Furthermore, we found that some sub-CLTs with similar topological structures and terminal cell type compositions are significantly overrepresented, suggesting that at least some stereotyped development is driven by a mitotic-coupling process. Together, we demonstrate the existence of stereotyped lineage trees, a feature of CLTs that likely contributes to stable cellular composition and therefore developmental robustness.

2. Results

2.1. Reconstructing High-Density Cell Lineage Trees for Directed Differentiation of Primordial Lung Progenitors

We aimed to determine the CLT of embryonic stem cells undergoing in vitro directed differentiation toward lung progenitors according to a well-established protocol recapitulating in vivo development.^[30] This in vitro model of directed differentiation was chosen for several reasons. First, cells cultured in a small petri dish have a relatively homogenous environment, so that transcriptome divergence caused by environmental factors, or phylogeny-independent convergence due to niche-specific signals is unlikely. Second, the development trajectory of embryonic stem cells to the lung is well-known, such that the in vitro cell culture can be monitored to ensure that they closely mimic physiological situation. Indeed, our implementation of the protocol can reach the alveolar epithelial cells (AEC2s) fate after 20 days of directed differentiation (Figure S1A and Video S1, Supporting Information). Third, in vitro culture allows us to induce Cas9 expression and therefore initiate the editing of the lineage barcode concurrently with the directed differentiation (Figure S1B,C, Supporting Information). Last but not least, it allows better control over the number of cells within the colony assayed for single-cell transcriptomes and CLTs. In particular, our cell culture begins with ≈ 10 hESCs and ends with ≈ 5000 cells on day

V. A. Lyubetsky, L. Y. Rusin
Kharkevich Institute for Information Transmission Problems Russian
Academy Sciences
Moscow 127051, Russia

V. A. Lyubetsky
Department of Mathematical Logic and Theory of Algorithms
Faculty of Mechanics and Mathematics
Lomonosov Moscow State University
Moscow 119991, Russia

X. Chen, J.-R. Yang
Key Laboratory of Tropical Disease Control
Ministry of Education
Sun Yat-sen University
Guangzhou 510080, China

10 (Figure S1B, Supporting Information), of which a relatively high percentage can be captured in downstream experimental pipelines of 10x Chromium. The ten-day directed differentiation covers three critical phases of lung development, including definitive endoderm (DE), anterior foregut endoderm (AFE) and NKX2-1⁺ primordial lung progenitor (PLP)^[30] (Figure 1A; Figure S1A,B, Supporting Information).

To assess the CLT of the cultured cells, we employed a modified scGESTALT method,^[23,31] which combines inducible cumulative editing of a lineage barcode array by CRISPR-Cas9 with large-scale transcriptional profiling using droplet-based single-cell RNA sequencing. Briefly, we initiated the editing of the lineage barcode concurrently with the directed differentiation using a Cas9 inducible by Doxycycline (Dox, Figure S1C, Supporting Information). We used an EGFP-fused cell lineage barcode that consists of 13 editing sites, each of which is targeted by one of four mCherry-fused sgRNAs each containing 2 to 3 mismatches in order to avoid large deletions resulting from excessive editing (Figure 1A; Figure S1D–F, Supporting Information). These sgRNAs were designed to not target any part of the normal human genome other than the integrated lineage barcode (Table S1, Supporting Information, see Experimental Section). The hESCs carrying the lineage tracing system were subjected to the ten-day directed differentiation, then the colonies were processed for cDNA libraries using the standard 10x Chromium protocol. Each cDNA library was split into two halves, with the first half subjected to conventional RNA-seq for single-cell transcriptomes, and the other half subjected to amplification of the lineage barcode followed by PacBio Sequel-based HiFi sequencing of the lineage barcode (Figure 1A).

We obtained single-cell transcriptomes of 3576/4400/1456/5659 cells respectively from three differentiating colonies CBRAD5-A1/G11/G2 and one parallel non-differentiating hESC colony, all of which appeared to have good quality (Figure S2A,B and Table S1, Supporting Information). The UMAP clustering of the single-cell transcriptomes revealed a large fraction of cells from differentiating/CBRAD5 colonies separated with those from hESC colonies, clearly indicating their differentiated cell states (Figure 1B). NKX2-1⁺ cells accounted for an average of 11.03% of the cells in each differentiating sample, which is highly consistent with previous findings.^[32] We identified 12 major functional clusters within the sampled cells (Figure 1C,D and See Experimental Section). According to the average expression of pluripotent gene (*NANOG*, *POU5F1*), endoderm progenitor gene (*GATA6*) and lung progenitor gene (*NKX2-1*, *SHH*, *CD47*), these clusters were defined as *NANOG*^{hi}*POU5F1*^{hi} (C1), *NANOG*^{low}*POU5F1*^{hi} (C2), *NANOG*^{low}*POU5F1*^{low} (C3), *NANOG*^{hi/low}*POU5F1*^{hi} (C4), *CD47*^{hi} (C5), *CD47*^{low} (C6), *GATA6*^{hi}*SHH*^{hi}*CD47*^{low} (C7), *GATA6*^{low}*NKX2-1*^{neg}*SHH*^{neg}*CD47*^{neg} (C8), *GATA6*^{hi}*NKX2-1*^{hi}*CD47*^{hi} (C9), *GATA6*^{hi} (C10). Below, they are also more broadly categorized into the less differentiated spontaneous state (R1 and R2) or pluripotent state (C1/C2/C3/C4), and the more differentiated progenitor state (C5/C6/C7/C8/C9/C10). These clusters displayed transcriptomic states largely compatible with known cell types occurred during the directed differentiation^[33] (Figure 1D,E; Figure S2C,D, Supporting Information), and were differentially distributed between hESC and CBRAD5 samples (Figure 1F), thereby suggesting successfully induced differenti-

ation and accurate measurement of single-cell transcriptomes. After confirming the sequencing quality of PacBio (Table S3 and Figure S3A, Supporting Information), the CLT of each sample was constructed based on the lineage barcode using maximum likelihood method (Figure 1A,G; See also Experimental Section, Figures S3 and S4 and Tables S4–S6, Supporting Information). The hierarchical population structures of the colonies were complex and intricate. In support of the accuracy of the CLT, cells more closely related to one another displayed more similar lineage barcode alleles (Figure 1H), and are more likely to share yet-to-decay transcripts of ancestral lineage barcode (Figure 1I). In conclusion, our experiment reliably captured the coarse-grained phylogenetic relationship of the cells within each colony.

2.2. The Cell Lineage Trees Recapitulate Key Features of the Transcriptome Divergence

To better elucidate the divergence between the single-cell transcriptomes in the context of the observed clusters, we identified differentially expressed genes (DEGs) in previously published microarray-based transcriptome^[32] data of samples from six timepoints of directed differentiation toward PLP (Figure 2A). Note here that despite being sampled on day12, the neural NKX2-1⁺ transcriptome has been shown to be most similar to that of day0 hESCs.^[32] The Gene Ontology (GO) terms enriched with these microarray-based stage-specific DEGs (Table S7, Supporting Information) were then individually examined for overall activities in our single-cell transcriptomes by the member genes' average expression levels in each cluster (Figure 2B and See Experimental Section). For pluripotent stage cells (C1/C2/C3/C4), significantly enhanced activities were found among GO terms enriched with DEGs of day 0/3 samples (including neural NKX2-1⁺) (Figure 2B). The same observations were made for progenitor stage cells C6/C10 in GO terms related to day3 samples, as well as C7/C9 cells in GO terms related to day6/day15 lung samples (Figure 2B). These results indicate that the single-cell transcriptomes recapitulated major differentiation stages of the in vitro PLP differentiation.

Our data also permit us to resolve divergence among sub-CLTs. It is commonly understood that the developmental process involves an increase in transcriptional divergence among cells and a reduction of developmental potentials in individual cells. Analyzing single-cell transcriptomes among sub-CLTs should reveal these patterns with fine resolution, especially when using high-density CLTs as we obtained. As an initial assessment for whether there is transcriptional divergence among sub-CLTs in the differentiating samples, we calculated for each sub-CLT, the CV (coefficient of variation) of the pseudotime estimates^[34] (see Experimental Section) of all its tips. When compared with their null expectations generated by randomly shuffling all tips, majority of these CVs were significantly smaller (Figure 2C), suggesting cells in the same sub-CLT are more similar than expected by the full range of transcriptional variation, an observation directly supports the transcriptional divergence among sub-CLTs.

For a more detailed analyses, we quantified the developmental potential of an internal node by the multivariate variance

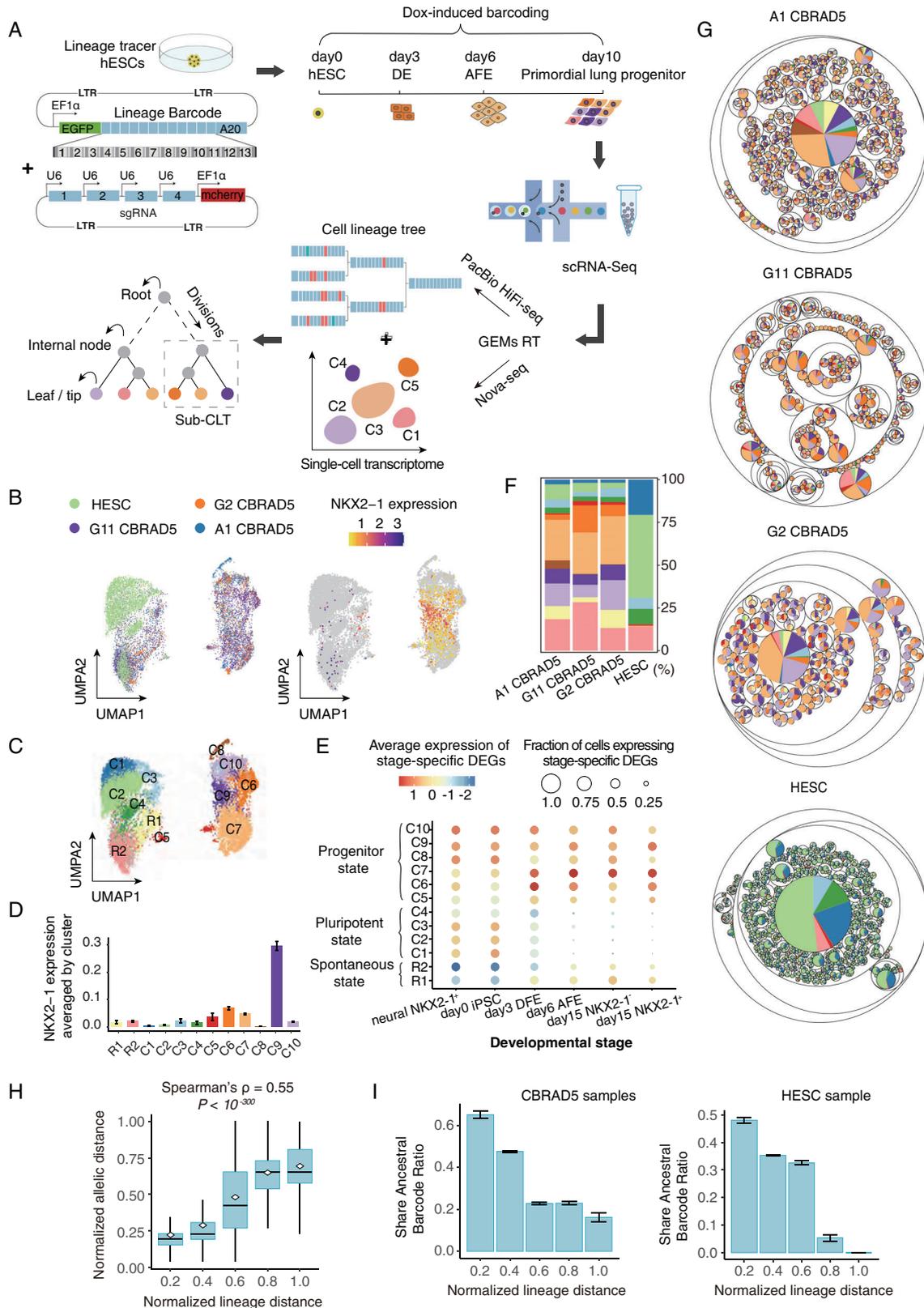


Figure 1. Cell lineage tracing for directed differentiation of primordial lung progenitors A) Schematic diagram illustrating the overall experimental process. The 10-day directed differentiation from several Lineage Tracer hESCs to primordial lung progenitors (PLP) was conducted along out with simultaneous lineage tracing utilizing inducible CRISPR-Cas9 editing of an expressed lineage barcode (13 editable sites). The resulting colony was assayed for single-cell transcriptomes by Nova-seq and lineage barcode by PacBio HiFi-seq, which were used to reconstruct CLTs with single-cell transcriptomes

among its descendant single-cell transcriptomes, which then allowed us to perform PERMANOVA-based statistical tests (PERmutational Multivariate Analysis Of VAriance, see Experimental Section) for the transcriptomic divergence. Briefly, by subtracting from the developmental potential of a focal node by the sum of the potentials of all its daughter nodes, we estimated the degree of divergence that occurred during the growth of the focal node (Figure 2D). Using the degree of divergence seen in the hESC sample as the null distribution, an average of $\approx 65\%$ internal nodes of the CBRAD5 samples displayed significant divergence (Figure 2E), whereas only $\approx 5\%$ internal nodes displayed divergence in the HESC sample. When such degree of divergence is depicted against normalized depths (see Figure S5A, Supporting Information and Experimental Section) of the corresponding nodes, the CBRAD5 samples consistently showed rapid divergence that is not seen in HESC samples (Figure 2E). Please note that divergence here is not equivalent to differentiation, since two sister cells differentiating into the same fate would not reveal any divergence for their mother cell. In other words, divergence implies asymmetric division creating daughter cells of different developmental potentials, whereas differentiation can occur during symmetric division giving rise to a pair daughter cells that both activate a particular function or differentiate in the same direction.

By applying the above analysis to gene subsets associated with specific GO terms, it is possible to elucidate the progression of divergence in the corresponding cellular functions. As shown in several key GO terms including Wnt signaling (Figure 2B), the cumulative growth in the fraction of internal nodes with significant divergence at various normalized depths is also highly reproducible among CBRAD5 samples, and it differs from the hESC sample (Figure 2F; Figure S5B, Supporting Information). Additionally, we examined whether our CLT data could resolve the temporal order of divergence completion for different cellular functions. To this end, we traced all root-to-tip paths on the CLTs and calculated the average depth of the last (furthest from the root) internal node exhibiting significant divergence on a GO term. As a result, the normalized depths of divergence completion appear consistent with known temporal orders of key developmental events (Figure 2G). Collectively, these results indicate that our dataset of single-cell transcriptomes and CLTs allowed the elucidation of cellular development with reasonable resolution.

2.3. Transcriptional Memory has Limited Contribution to Developmental Canalization

Following confirmation of the CLT data's resolution, we began searching for contributors to developmental robustness using CLTs. A first hypothesis is that transcriptional memory may have constrained gene expression variation during development, which would canalize transcriptomic state during development and contribute to robustness. In this context, transcriptional memory is the phenomenon of cells closely related on the CLT displaying similar expression levels due to the inheritance of the same cellular contents (proteins/transcripts) and/or epigenetic states from recent common ancestors.^[29,31,35,36] Nevertheless, gene expression can also be restricted by transcriptional regulation that has nothing to do with cellular inheritance, such as negative feedback^[37] and denoising promoters.^[38] If the transcriptional memory dominates the experimented differentiation, one would expect all cells of the same type would have been clustered into an exclusive sub-CLT, which is clearly not the case (Figure 1G). For a quantitative analysis, we reasoned that the CV of single-cell expression levels within real sub-CLTs should reflect the combined effect of transcriptional memory and inheritance-independent regulation (Figure 3A top), whereas that of CLTs randomized by shuffling cells of the same type at different lineage positions should reflect only inheritance-independent regulation but not transcriptional memory (Figure 3A bottom). It is therefore possible to isolate the contribution of transcriptional memory to the expression constraint by contrasting the CV of real CLTs with that of randomized CLTs (Figure 3A and Experimental Section), which is hereinafter referred to as the "memory index". We note that this definition of memory index is similar to that used in previous transcriptional memory-related studies.^[29,39]

For each cell type, we calculated an overall memory index for each gene in each sub-CLT (Figure 3B; Figure S6A,B, Supporting Information). The top (10%) memory indices (Figure 3C) were found to be enriched in pluripotent cell types (C1/C2/C3/C4) as compared to progenitor cell types (C6/C7/C9/C10) (t -test $P = 0.0039$, Figure 3D), suggesting that transcriptional memory is more important to maintaining pluripotency than differentiation. Because transcriptional memory is mediated by cellular contents inherited from mother to daughter cells, such as transcription factors, we hypothesized that these genes with top memory indices should exhibit significant overlap with those regu-

assigned to tips. B) The variation among single-cell transcriptomes captured in the four samples (one non-differentiating "HESC" sample and three differentiating samples) as shown by UMAP. A data point represents a cell, which is colored based on its source sample on the left panel and the expression level of NKX2-1 (the marker for PLP) on the right panel. C) Major clusters of the single-cell transcriptomes are differentially colored and labeled by their corresponding cell types. D) The average expression levels of NKX2-1 in each cell type. The error bars indicate the standard error among single cells. E) In the 12 major cell types (y axis), differentially expressed genes (DEGs) found in bulk samples of specific developmental stages preceding PLP (x axis) were examined for their average expression levels (dot color) and fraction of cells that expressed the gene (dot size). See also Figure S2C (Supporting Information). F) For each of the four samples (x axis), the percentage of cells belonging to each type was shown. The cell types are colored identically to those in panel C. G) Reconstructed CLTs are visualized as circle packing charts for the four samples. Circles represent sub-CLTs, whose sizes indicate the number of terminal cells in the sub-CLTs, while the color (same as panel C) indicates the fraction of terminal cells belonging to each cell type. See Figure S3B (Supporting Information) for their tree representation. H) A pair of cells' normalized lineage distance (the number of internal nodes on the path from one cell to the other, divided by the maximal lineage distance found in the sample) is highly correlated with the normalized allelic distance of their lineage barcodes (the total number of target sites that differed from the reference, divided by the maximum value of 26). All cell pairs were separated into five groups based on their normalized lineage distance (x axis), and the distribution of normalized allelic distances (y axis) within each group is shown in the form of a standard boxplot, with the mean value indicated by the white point. On top, Spearman's ρ and P value for raw data are indicated. I) The probability of finding a common ancestral allele (as yet-to-decay transcripts) between a pair of single-cell tips decreased as their normalized lineage distance (x axis) increased. The error bars indicate the standard error estimated by bootstrapping the cell pairs for 1000 times.

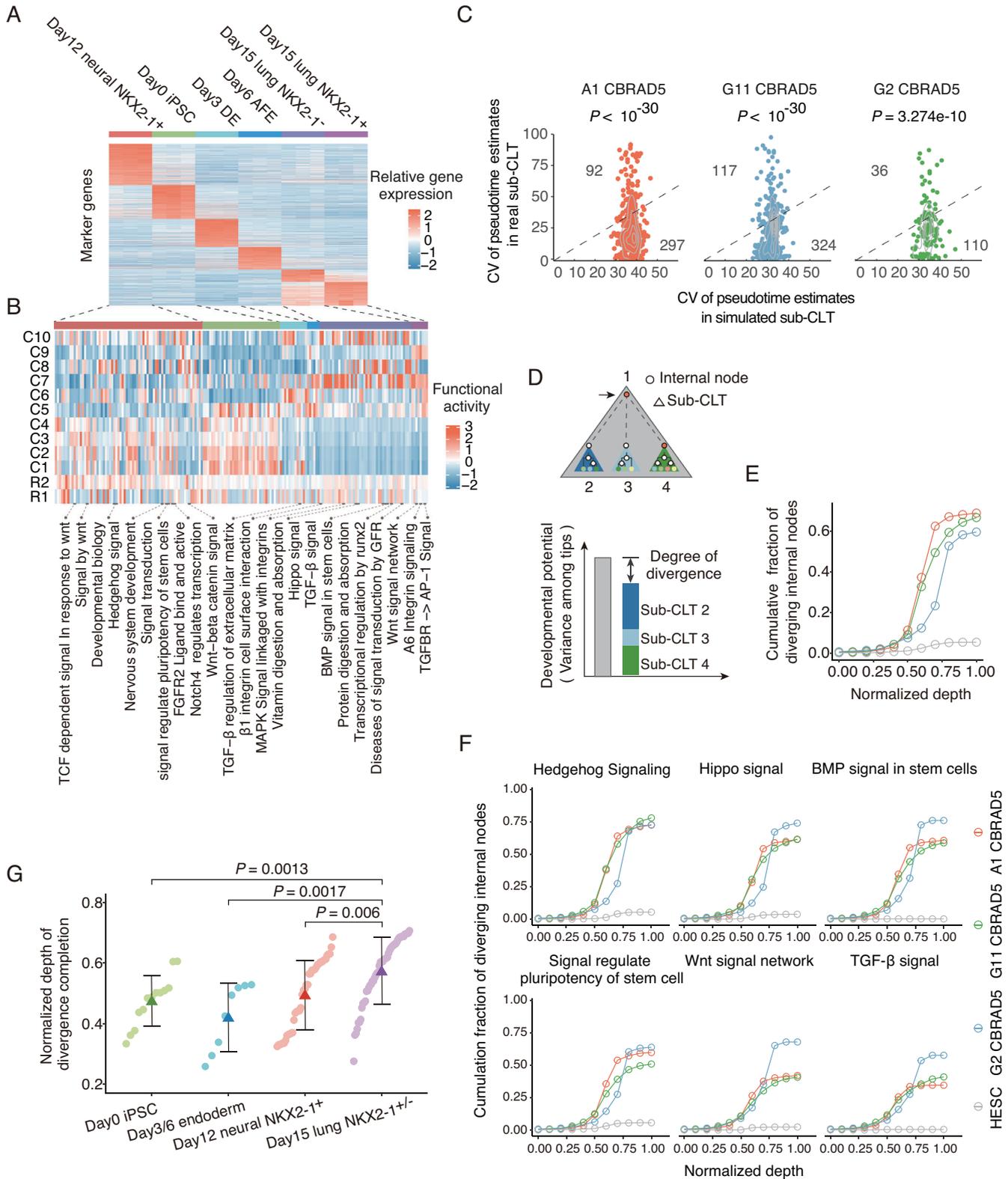


Figure 2. The transcriptome divergence among cell type clusters and among subclones A) Heatmap for expression levels of DEGs extracted from microarray-based transcriptomes of specific developmental stages (color bars on top) of the directed differentiation.^[32] B) Functional activities of GO terms (x axis. Full list in Table S7, Supporting Information) enriched with stage-specific DEGs were shown for every cluster (y axis) identified in our samples. Here functional activity as indicated by the color scale was estimated by the average Z-score-transformed expression of all genes annotated with the GO term. Some important GO terms are boxed and labeled by dashed lines, and are further analyzed in panel (F) and Figure S5B (Supporting

lated by some related transcription factors. Thus, we tested these genes for enrichment in genes responsive to genetic perturbation of individual transcription factors^[40] (see Experimental Section), and made three observations. First, some transcription factors with known involvement in the experimented differentiation, such as NANOG/MYC/TFE3 in the pluripotent C1^[41,42] and MECOM/KLF5/GATA6 in progenitor C6/C7/C9,^[32,43–46] indeed exhibit significant enrichment of the genes with top memory index. This suggests that transcriptional memory contributed to the regulatory role of these transcription factors, for instance, the ability of NANOG/MYC/TFE3 to maintain pluripotency and the ability of MECOM/KLF5/GATA6 to promote differentiation into lung progenitors. Second, the union of the top ten transcription factors from every cell types is significantly enriched in pluripotency or differentiation regulation-related GO terms (Figure S6C, Supporting Information), again supporting the contribution of transcriptional memory to related pathways. Closer examination of specific transcription factors assists in resolving specific regulatory functions mediated by transcriptional memory, such as Pyruvate metabolism^[47–49] being regulated by NANOG (Figure S6D, Supporting Information). Third, the enrichment was generally stronger for pluripotent cell types than it was for progenitor cell types (Figure 3E), a pattern again suggesting that transcriptional memory only played a minor role in differentiation, which is at least not as significant as in maintaining pluripotency.

2.4. Stable Cell Type Compositions Across Sub-Cloned

Observations above indicate that terminal cells within a sub-CLT have restricted fates that are not dominated by transcriptional memory from the common ancestor (root of the sub-CLT). This observation automatically prompted an assessment of the cell fate restriction imposed by inheritance-independent regulation, as well as its contribution to the robustness of developmental processes. We reasoned that inheritance-independent regulation should result in multiple similarly restricted sub-CLTs dispersed across the entire CLT. Therefore, we calculated the terminal cell type composition for each sub-CLT found in the CBRAD5 samples and compared it with the overall composition of the corresponding full CLT (see Experimental Section). Intriguingly, the cell type compositions of sub-CLTs are usually more similar to those of the full CLTs than expected in randomized CLTs (Figure 4A–C). A closer examination of some sub-CLTs

reveals a highly stable terminal cell type composition. For example, there are 35 sub-CLTs that generated subclones with highly stable (<10% deviation) proportions of 0.13, 0.39, 0.13, and 0.18 respectively for C6, C7, C9, and C10 (the top four most abundant progenitor cell types), which corresponds to the average proportion of these cell types in the three differentiating samples (Figure 4D). This observation suggests that a stereotyped developmental program may exist that produces subclones with highly similar compositions of cell types derived from multiple ancestral cells.

The observed polyclonal stereotypic development can be understood from two perspectives. On the one hand, the consistent execution of such a developmental program across subclones may be by itself a manifestation of robust genetic and/or molecular regulation. On the other hand, stable cell type compositions across subclones might enhance developmental robustness. We examined this latter perspective by simulating a CLT for the development of a single cell into an “organoid” consisting of 1024 cells (i.e., 10 cell cycles) comprised of four types (namely α , β , γ , and δ) of cells in a 1:1:2:4 ratio. These cells formed 128 functional units each consisting of one α cell, one β cell, two γ cells, and four δ cells. Normally developed organoid consisting of 128 functional units (assuming sufficient cellular migration) are considered 100% functional. Meanwhile, CLT perturbed by random cell deaths (see below), which results in the loss of some ancestral cells and all their descendants, has a functional capacity defined as the fractional survival rate of functional units with proper cellular composition. This design was inspired by the observation that functional units in living tissues, such as mouse pancreatic islets, display a highly stable cell type composition as the outcome of normal development.^[26] To generate the normal (death-free) CLT with the predetermined number of cells of each type, two models were used. The first “random” model assigns each cell to a random tip of the CLT regardless of its cell type (Figure 4E left). A second “stereotyped” model defines all eight-tip sub-CLTs as strictly consisting of one α cell, one β cell, two γ cells, and four δ cells, but different placements of these cells are allowed on the tips (Figure 4E right). A total of 1000 normal CLTs were generated under each model, and the functional capacity of each CLT was determined by exposing all (internal or terminal) cells to various rates of random death. When compared to the random model, we found that CLTs generated with the stereotyped models always formed more functional units, or in other words, were more robust against cell deaths (Figure 4F). Such enhanced developmental robustness is more evident at higher rate of cell death (Figure 4F). Collectively, these results suggest that the ob-

Information). C) A coefficient of variation (CV) was calculated using pseudotime estimates of single-cell transcriptomes within a sub-CLT. These CVs were plotted for all real sub-CLTs (y axis) and corresponding randomized sub-CLTs generated by shuffling all tips (x axis) in each differentiating sample (name on top). As the dashed line indicates $x = y$, sub-CLTs with CVs lower than random expectation (i.e., restricted variation) will appear below it. Each panel includes the number of CLTs above and below the dashed line, which was also tested against the binomial expectation (50% below the line) and yielded the P values on top. D) Schematic diagram for the PERMANOVA-based estimation of transcriptome divergence for an internal node (see Experimental Section). E) Cumulative fraction (y axis) of internal nodes exhibiting significant transcriptome divergence as the normalized depth (x axis) considered increased. Results from different samples were shown with different colors, as indicated by the color legend. F) Same as panel E except that the analyses were limited to specific GO terms indicated on top of each panel. G) We calculated the normalized depths (y axis) at which the divergence of specific functions is completed. GO terms enriched of marker genes in representative developmental stages (x axis and colors) were examined. Dots represent GO terms and triangles represent the average depth within the same-color group. Number of GO terms are 37 for day0 iPSC, 19 for day3/day6 endoderm, 71 for day12 neural NKX2-1⁺, and 52 for day15 lung NKX2-1^{+/−}. Significant P values from between-group Wilcoxon Rank Sum test are labeled on top.

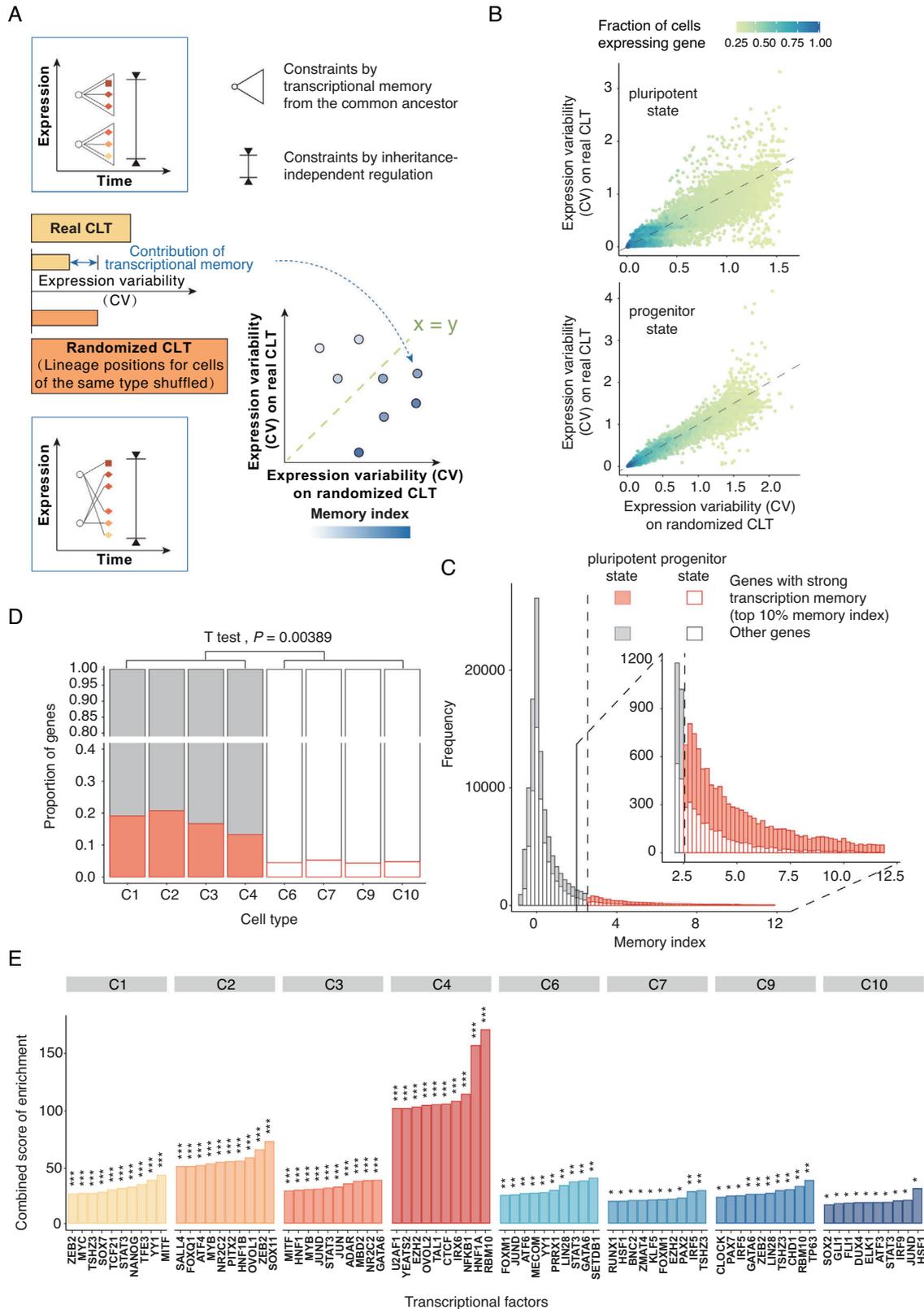


Figure 3. Limited contribution of transcriptional memory in differentiation A) Schematic diagram for the CLT-based estimation of transcriptional memory. B) Expression variability in the real CLT (y axis) compared to that in the randomized CLT (y axis). Each dot represents a gene in a cell type. Dot color shows the fraction of cells within the cell type that express the gene, as indicated by the color scale on top. C) A stacked histogram showing the distribution of the memory indices calculated. A filled bar represents those estimated from pluripotent cell types and an empty bar represents those estimated from

served stable cell type composition among subclones contributed to developmental robustness.

2.5. Stereotyped Cell Lineage Trees Underlie Stable Cell Type Compositions

We next seek further evidence for the existence of stereotyped developmental programs based on the CLT data at hand. Specifically, we hypothesized the existence of multiple sub-CLTs with highly similar topology and terminal cell types. Note that the similarity in sub-CLT topology is an additional requirement beyond the similarity of cellular compositions observed above, and the similarity in both topology and cellular composition is compatible with previously proposed “mitotic coupling” mode of cell state-lineage relationship.^[18] As recurrent sub-sequences of biological sequences, such as transcription factor binding sites, are usually referred to as “sequence motifs”, we call our target recurrent sub-CLTs “tree motifs” or simply “motifs”. In fact, some tree motifs in development have been well characterized. For example, the *Drosophila* peripheral nervous system contains thousands of identical mechanosensory bristles.^[24] Each of the bristles is formed by a sub-CLT rooted at a sensory organ precursor cell. This sub-CLT encompasses two cell cycles, the first of which produces PIIa and PIIb cells. Then PIIa divides to yield one shaft cell and one socket cell, followed by PIIb, which gives rise to one neuron and one sheath cell.^[24] Therefore, this specific tree motif appears thousands of times in *Drosophila*'s developmental CLT. Furthermore, the meiosis process, in which one germ cell divides into four gametes, is another example of a tree motif in developmental CLTs.

Just as sequence motifs are identified by comparisons between (sub-)sequences, tree motifs should also be identified through comparisons between (sub-)CLTs. In order to identify potential tree motifs in the CLT of the differentiating samples, we utilized Multifurcating Developmental cEll Lineage Tree Alignment (mDELTA), an algorithm we previously developed for quantitative comparisons and alignments between CLTs^[50] (Figure 5A, see Experimental Section and Text S1, Supporting Information). Using a dynamic programming scheme analogous to that employed by classical algorithms looking for similarities between biological sequences (e.g., the Smith-Waterman algorithm), the mDELTA algorithm searches for pairs of homeomorphic sub-CLTs^[50] within two given full CLTs. As a result, mDELTA identified a large number of highly similar sub-CLT pairs between and within differentiating samples (Figure 5B). Some of the most frequently occurring sub-CLTs exhibited a consistent structure, comprising multiple layers of internal cells, a stable composition of terminal cell types, and appeared 20 to 40 times in the three differentiating samples, whose summed mDELTA alignment scores are significantly higher than expectation assessed in 1000 randomized trees (Figure 5C). Groups of such highly simi-

lar sub-CLTs represent strong candidates of tree motifs on the developmental CLT, and strongly supports the existence of a stereotyped developmental program that contributes to developmental robustness.

3. Discussion

In the current study, we have reconstructed high density developmental CLTs for in vitro directed differentiation from hESC to primordial lung progenitors. In comparison with CLTs of non-differentiating hESC colonies, differentiation CLTs showed a clear signal of transcriptomic divergence that recapitulates known involvements of key developmental regulatory pathways. Using CLTs, we investigated mechanisms that might have contributed to developmental robustness at the intercellular level. Although transcriptional memory appeared to have limited effects on canalizing cell fates within subclones, we found that multiple subclones exhibit stable compositions of terminal cell types, which enables sufficient numbers of cells in proper composition to be generated, and thus, a more robust development. By using a CLT alignment algorithm, we further showed that the observed stable cell type composition is underlied by stereotyped sub-CLTs with similar topology and terminal cell fate. Our results demonstrated the existence of stereotyped sub-CLTs, which support robust development.

During development, cell death in various forms is pervasive, making robustness against cell death an essential aspect of the developmental process. Cell death can be triggered by external stimuli/stress such as in necrosis, or can occur autonomously in response to internal signals that are more or less stochastic, such as shown in the cell-to-cell random variations in Bcl-2 regulated apoptosis.^[51] We chose to model developmental robustness against random cell death instead of cell deaths induced by specific external stimuli or stress, in part because the former is more general than any environment specific stimuli. Furthermore, it has previously been demonstrated, also by CLT-based analyses, that robustness against random cell death should enhance robustness against specific types of cell deaths, including those caused by genetic perturbations.^[17] It would, however, be very interesting to investigate the developmental robustness against cell deaths or other adverse events triggered by specific types of environmental stress.

As a preliminary assessment on how the stereotyped CLT occurs, we treated the cell type composition of all descendent tips as a quantitative trait of the ancestral cells (internal nodes of the CLT) and regressed the difference of this trait between two ancestral nodes (that is not descendent of each other) onto their relatedness on the cell lineage (see Experimental Section). This method, known in the genetics literature as a Haseman–Elston Regression,^[52,53] is an unbiased estimator of heritability. In all of our samples, cell type compositions displayed heritability to some

progenitor cell types. Genes exhibiting strong transcriptional memory, i.e., those with a memory index ranking among the top 10% (dashed line), were red, while others were gray. The inset shows a zoomed-in view of the large memory index region. D) Among different cell types, the fraction (height of bar) of genes exhibiting high memory indices was compared. The bars are colored similarly to those in panel C. E) Gene sets responsive to perturbation of individual transcription factors (x axis) were tested for the enrichment of genes exhibiting strong signal of transcriptional memory (see Experimental Section). The top ten transcription factors with the highest combined enrichment score (y axis) were shown for each cell type. The statistical significance of enrichment according to Fisher's exact test is indicated as *: $P < 0.05$; **: $P < 0.01$; ***: $P < 0.001$. See also Figure S6 (Supporting Information).

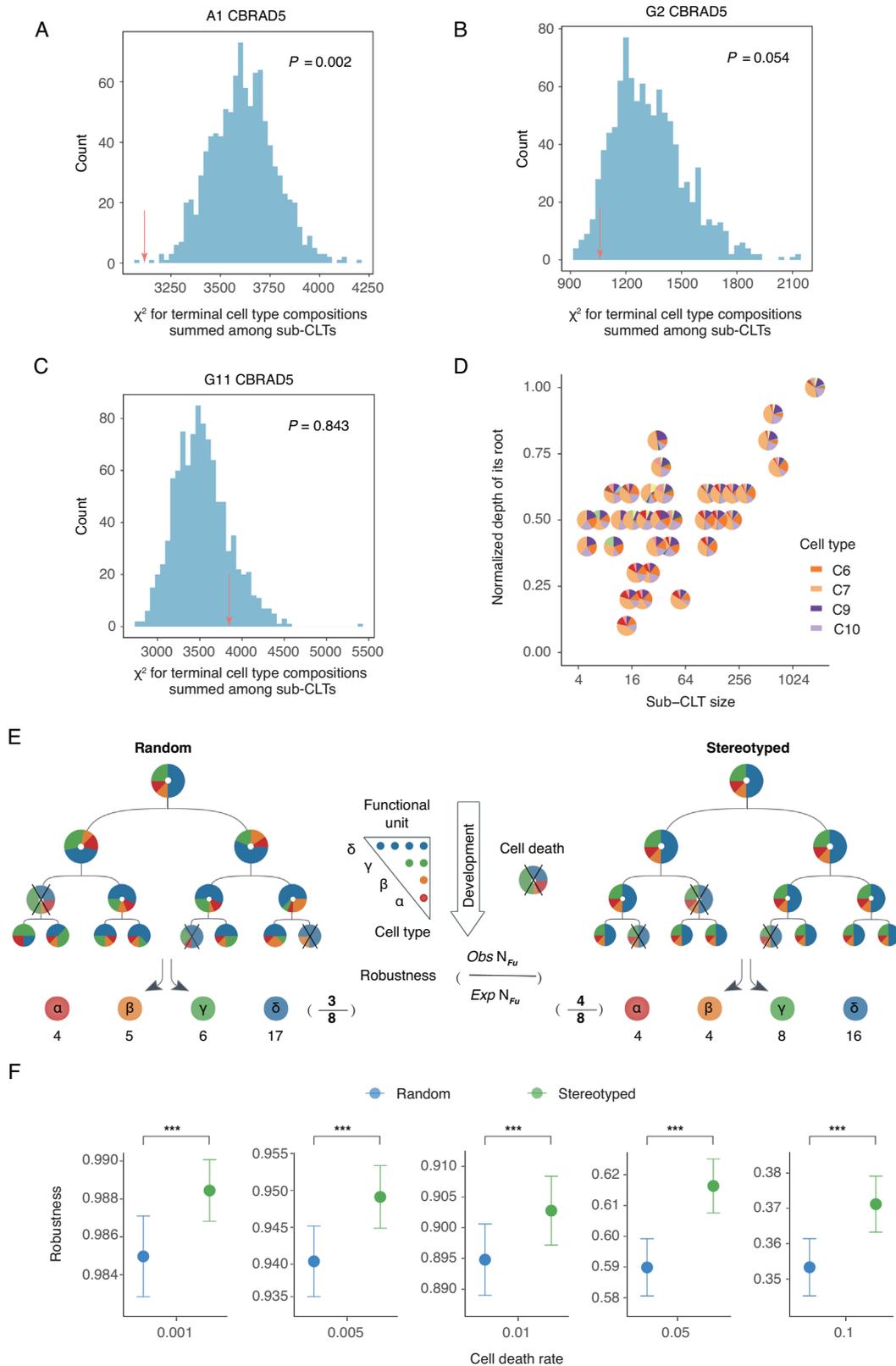


Figure 4. Stable cell type composition across sub-clones supports robust development A–C) In each panel for each of the CBRAD5 samples (names on top of the panel), the diversity of compositions of terminal cell types within sub-CLTs were estimated by a summed chi-square value (χ^2) (see Experimental Section) as indicated by the red arrows. The same summed χ^2 values were calculated for 1000 randomized CLTs, whose distribution was shown as a blue histogram. The probability of a summed χ^2 value being smaller than the observation (red arrow) is indicated by the P values in the

degree, with the heritability in the differentiating samples being significantly greater than that in the non-differentiating sample (Figure S7, Supporting Information). Furthermore, similarly estimated heritability of single-cell transcriptome for each sample were lower than that of cell type composition (Figure S7, Supporting Information). This result is unlikely to be explained by the higher measurement accuracy of cell type composition compared to single-cell transcriptomes for two reasons. First, the cell type itself is inferred based on single-cell transcriptomes. Second, the heritability of cell type composition in the non-differentiating sample is almost equal to that of the single-cell transcriptome, suggesting similar measurement accuracy for these two traits. Thus, we concluded that descendent cell type composition is a heritable trait of ancestral cells. This trait is likely inherited from their earlier common ancestors by a mechanism independent of transcriptional memory, and is therefore expected to be pervasive in a CLT.

Beyond the specific mechanisms underlying developmental robustness, our findings suggest a novel perspective regarding cell types within the context of stereotyped sub-CLTs. In particular, just as letters can be better understood within the context of words, and nucleotides/amino acids can be better understood within the context of sequence motifs, stereotyped sub-CLTs can potentially bridge our knowledge of the atlas of cell types and their organization into functional tissues. Indeed, Elowitz and colleagues^[54] recently identified statistically overrepresented patterns of cell fates on lineage trees as indicative of progenitor states or extrinsic interactions. The analysis was done using their newly proposed Lineage Motif Analysis, which differs from the method presented here that examined cell type composition and topological structure on incomplete CLTs, as their method uses the fully resolved CLTs and only analyzes cell type composition. Nevertheless, similar to our proposition here, they considered lineage motifs as a way of breaking complex developmental processes down into simpler components.^[54]

3.1. Limitations of the Study

There are a couple limitations of our study that are worth discussing here. First, our study was based on an in vitro directed differentiation model, which maintains cells in a low density 2D monolayer. This choice is a compromise between the feasibility for reconstruction of high density CLTs and a model that closely reflects the in vivo development. We believe our experiment reasonably recapitulates the in vivo situation because clear morphology of alveolar can be achieved on the 20th day of the directed differentiation (Figure S1A and Video S1, Supporting Information), and it has previously been shown that 3D-conditions did not increase gene expression in comparison to 2D for the differentiation up to the NKX2-1⁺ PLP stage, particular for the NKX2-1.^[55]

For development up to a later stage with more cells, organoid or in vivo models should ideally be combined with single-cell transcriptomes of a larger throughput (in terms of number of cells) in order to assess the question at a broader scale. Nevertheless, our main conclusion of polyclonal stereotyped development is most likely NOT an artefact of in vitro development, because none of the media components can create such pattern, and the number of ancestor hESCs seeding the colony is not correlated with the frequency of recurrence of lineage motifs. Second, various antibiotic treatments were used during the construction of the cell lineage and the directed differentiation, which may have an impact on the pluripotency and differentiation of the cells. This possibility, however, should be largely negligible in our system because, on the one hand, previous study^[56] showed that antibiotics such as doxycycline increase cell survival rate without apparent negative side effects, and on the other, our experiment produced a proportion of NKX2-1⁺ cells that were comparable to those obtained under a standard in vitro differentiation condition.^[30,57] Third, we used Dox-inducible Cas9 to switch the lineage tracing system on and off, which could lead to unintended premature editing before differentiation begins, or temporarily interrupted editing during differentiation. However, the editing status of the lineage barcode in most terminal cells (Figure S4, Supporting Information) seem to provide adequate information capacity for the resolution of their lineage relationship and thereby the stereotyped sub-CLTs, as most (3525/4076 = 86.48%) unique barcodes still have at least two intact editable sites. Additionally, we reasoned that the imprecise induction of Cas9 would produce noisy lineage tracing data that should obscure but not strengthen the biological signal. As our current lineage tracing results already support our argument about stereotyped development, the underlying biological signal should thus be even stronger. Fourth, we have not inferred detailed molecular processes and/or trajectories of gene expression changes in the stereotyped sub-CLT, as can be done for the nematode *Caenorhabditis elegans*,^[50] whose temporal changes in gene expression have been recorded by microscopic image.^[58,59] In the near future, this may be possible when the algorithms for inferring ancestral states based on cell lineage trees become sufficiently accurate.^[19,60]

4. Experimental Section

Design of the Lineage Tracer hESC Cell Line: To design the lineage barcode and corresponding sgRNA, randomized 20-bp candidate sgRNA sequences with >3 substitutions relative to any human genome fragments were first generated. Among these candidates, the spacer sequence 5'-TATTCGCGACGGTTCGTACG-3' was selected as sgRNA1. A total of 13 protospacer sequences were designed based on sgRNA1 according to the following criteria: i) each protospacer contained 2–3 mismatches with sgRNA1, ii) there was no recurrence of any sequence of 9 bp or longer, and iii) consecutive repeats of the same nucleotide for more than 2 bp

panel. D) For 35 sub-CLTs in CBRAD5 samples, the normalized depths of their roots (γ axis) and the sizes of the sub-CLTs (x axis) were plotted. These sub-CLTs display highly similar terminal cell type compositions (less than 10% deviation from 0.13, 0.39, 0.13, and 0.18 respectively for C6, C7, C9, and C10) E) A schematic diagram showing a simple model of the functional robustness of the random (left) versus stereotyped (right) development against random cell deaths (indicated by "X"). The robustness is quantified by the number of functional units (with cell type compositions indicated in the triangle) that can be formed by terminal cells surviving cell deaths, as exemplified at the bottom. F) Robustness (γ axis) of the random (blue) versus stereotyped (green) development under different rate of cell death (x axis), as estimated by the model in E. The statistical significance of student's t -test is indicated as ***: $P < 0.001$.

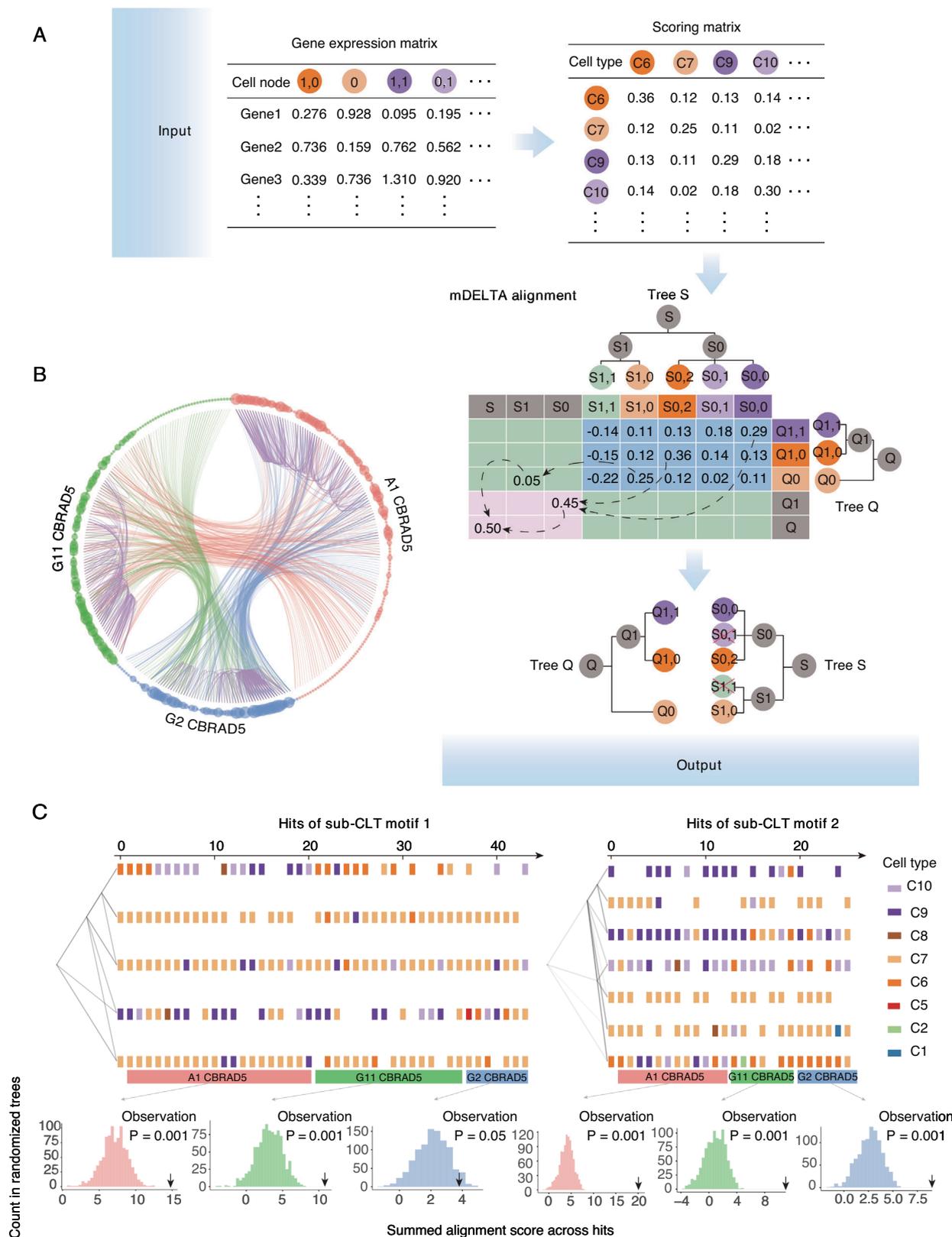


Figure 5. Stereotyped subtrees underlying the stable cell type composition A) The input (top) for DELTA includes two CLTs (query and subject) and the expression profiles of all terminal cells on these CLTs. DELTA uses a dynamic programming procedure (middle) to compare the two CLTs and identify homeomorphic sub-CLTs. The procedure has three phases, including (i) a cell pair scoring stage, (ii) a forward stage that maximizes the alignment scores by finding the best correspondence between terminal cells, and (iii) a backtracking stage for extracting the alignment behind the maximized scores.

were completely absent. The 13 protospacers (along with PAM, or protospacer adjacent motif) were organized according to decreasing CFD (cutting frequency determination) scores into the full lineage barcode.^[61,62] The next three sgRNAs, sgRNA2, sgRNA3, and sgRNA4, were designed to perfectly match the 9th, 12th, and 13th protospacers, but with lower CFD scores (<0.55) for other protospacers, because these three protospacers were rarely edited in preliminary experiments using only sgRNA1. To facilitate capture by poly-dT reverse transcription primers on 10x gel beads, the full lineage barcode with a 20-nt poly-dA(A20) 3' tail was inserted into the 3' UTR of an EGFP driven by an EF1 α promoter.

Lineage tracer hESC cell lines were constructed by genomic integration of the lineage barcode, Doxycycline-inducible Tet-on Cas9 and the sgRNAs. Briefly, the lineage barcode vector (pLV-EF1A>EGFP:T2A:Bsd:V1(Barcode), VectorBuilder, no:VB1709 11-1008qmt) was constructed by the Gateway system and then transfected into H9 hESCs with MOI = 0.15. The EGFP-fused lineage barcode was confirmed to exist as a single copy in the genome and to be highly expressed after blasticidin selection (15 $\mu\text{g ml}^{-1}$, InvivoGen, no. ant-bl-1) and flow cytometry sorting. Then the Tet-on inducible Cas9 vector (PB-Tet-ON-T8>Cas9:T2A:puro-PGK:rtTA, donated by Professor Jichang Wang, Zhongshan School of Medicine, Sun Yat-sen University) was co-transfected with hyPBBase (VectorBuilder, no: VB190515-1005npr) in a ratio of 1 μg :100ng for 1×10^7 /ml cells by Neon transfection system (Life, MPK5000). In order to ensure adequate Cas9 expression for efficient editing, double reinforced selection of Puromycin (1.0 $\mu\text{g ml}^{-1}$, InvivoGen, no. ant-pr-1) and Doxycycline (Dox, 1.0 $\mu\text{g ml}^{-1}$, sigma, D9891-1G) was applied for 7 days. Lastly, the sgRNA vector (pLV-U6>sgRNA1>U6>sgRNA2>U6>sgRNA3>U6>sgRNA4-EF1 α >Mcherry:T2A:Neo, VB1912 11-3149jwe) was constructed by Golden Gate ligation and transfected at MOI = 30. H9 hESC cells with high expression of sgRNAs (fused with mCherry) were enriched by G418 selection (1000 $\mu\text{g ml}^{-1}$, InvivoGen, ant-gn-1) for 11 days and flow cytometry sorting. Expression levels of Cas9, lineage barcode and sgRNA1 transcripts were detected by RT-qPCR with primers listed in Table S8 (Supporting Information).

The editing efficiency of the lineage tracer hESC cell line was evaluated by inducing Cas9 expression in mTesR media with 1.0 $\mu\text{g ml}^{-1}$ Dox for five days. gDNA was extracted from all cells using DNeasy Blood & Tissue Kits (Qiagen, no.69504). Using primers gDNA-V1-F and gDNA-V1-R (Table S8, Supporting Information), the lineage barcode was amplified from gDNA using Phanta Max Super-Fidelity DNA Polymerase (Vazyme, No. P505), which was then cloned into pCE-Zero vector (Vazyme, No. C115). The efficiency of editing was then evaluated by colony PCR and Sanger sequencing for 50 recombinant clones.

Additionally, editing efficiency was examined in the context of the directed differentiation experiment, in which only a small number of initial cells were used to form each colony. In 96-well dishes, matrigel (Corning, No. 354277) was plated and each well was seeded with < 10 log-phased lineage tracer hESC cells manually by micromanipulation. For 11 days, the cells were cultured in 100 μl of mTesR media, to which 10 μl of cloneR

(Stemcell, No.05888) were added on day0 and day2, and 1.0 $\mu\text{g ml}^{-1}$ Dox+mTesR media was added and refreshed every 48 h since day2. Normally surviving colonies after the 11-day culture were harvested by GCDR (Stemcell, No.07174). Next, 50 ng of genomic DNA was extracted from each colony using the QIAamp DNA Micro Kit (Qiagen, No.56304) and PCR amplified for the lineage barcode. The Cas9-induced mutations accumulated during colony formation were then identified by Sanger sequencing, TA cloning-based sequencing or Illumina HiSeq PE250 sequencing. Specifically, the raw HiSeq data were trimmed by fqtrim (<https://ccb.jhu.edu/software/fqtrim/>) with default parameters. The paired reads were merged by FLASH^[63] using 30 bp of overlapping sequence and 2% mismatches. Sequences alignable to the human reference genome by Bowtie2 with default parameters,^[64] or to primer sequences of gDNA-V1-F and gDNA-V1-R with two mismatches, were removed as they likely represented non-specifically amplified sequences. MUSCLE^[65] aligned the sequenced lineage barcode to the wild-type lineage barcode using default parameters. The editing events of each sequence were identified according to a previous method.^[61]

Validating Directed Differentiation from hESC to Lung Progenitor and Alveolosphere: Using the BU3 NGST (NKX2-1-GFP; SFTPCtdTomato) iPS cell line (donated by Professor Darrell N. Kotton, Department of Medicine, Boston University), the protocol of directed differentiation was tested toward lung progenitor and alveolosphere published by Kotton and colleagues.^[30] Briefly, in six-well dishes pre-coated with Matrigel (Stemcell, No.356230), 2×10^6 cells maintained in mTESR1 media were differentiated into definitive endoderm using the STEMdiff Definitive Endoderm Kit (StemCell, No.05110), adding supplements A and B on day 0, and supplements B only on day 1 to day 3. Flow cytometry was used to evaluate the efficiency of differentiation to definitive endoderm at day 3 using the endoderm markers CXCR4 and c-KIT (Anti-human CXCR4 PE conjugate, Thermo Fisher, MHCXCR404, 1:20; Anti-human c-kit APC conjugate, Thermo Fisher, CD11705, 1:20; PE Mouse IgG2a isotype, Thermo Fisher, MG2A04, 1:20; APC Mouse IgG1 isotype, Thermo Fisher, MG105, 1:20) based on the method of Sahabian and Olmer.^[66] After the endoderm-induction stage, cells were dissociated for 1–2 min at room temperature with GCDR and passaged at a ratio between 1:3 to 1:6 into 6 well plates pre-coated with growth factor reduced matrigel (Stemcell, No.356230) in “DS/SB” anteriorization media, which consists of complete serum-free differentiation medium (cSFDM) base, including IMDM (Thermo Fisher, No.12440053) and Ham's F12 (Corning, No. 10-080-CV) with B27 Supplement with retinoic acid (Gibco, No.17504044), N2 Supplement (Gibco, No.17502048), 0.1% bovine serum albumin Fraction V (Sigma, A1933-5G), monothioglycerol (Sigma, No. M6145), Glutamax (ThermoFisher, No. 35050-061), ascorbic acid (Sigma,A4544), and primocin with supplements of 10 μM SB431542 (“SB”); Tocris, No.1614) and 2 μM Dorsomorphin (“DS”); Sigma, No. P5499). In the first 24 h following passage, 10 μM Y-27632 was added to the media. After anteriorization in DS/SB media for three days (72 h, from day 3 to day 6, refreshed every 48 h), cells were cultured in “CBRa” lung progenitor-induction media for nine days (from day 6 to day 15, refreshed every 48 h). This CBRa me-

The output (lower right) is one or more aligned sub-CLTs ordered by decreasing alignment scores. See Experimental Section and Text S1 (Supporting Information) for more details. B) A circular plot of the top 100 sub-CLT pairs found by mDELTA in each of the six pairwise comparisons among the CLTs from the three differentiating samples. In the outer circle, each sub-CLT is represented by a dot, with the color indicating its source sample. Each pair of homeomorphic sub-CLTs identified by mDELTA is shown by curved links between two corresponding dots, where inter-sample pairs/links are colored the same as the sample used as the query CLT, and intra-sample pairs/links are colored purple. A dot's size indicates how many links it has. Only sub-CLTs with at least one link are included. C) Two highly recurrent tree motifs found in the three samples are shown by “densitree” plots. For each motif, all sub-CLTs homeomorphic to a specific reference sub-CLT (the focal motif) are extracted from mDELTA results in panel B. In each plot, the mDELTA-aligned topological structure of each sub-CLT (including the reference sub-CLT) is drawn with transparency on the left so that common topologies can be seen as darker lines. Each column of tiles on the right shows the mDELTA-aligned terminal cell types (colored as the label on top) on one of the homeomorphic sub-CLTs. The x axis on top indicates the number of sub-CLTs hit as homeomorphic to the reference sub-CLT, which is displayed at $x = 0$. The hits were separated by their source sample as indicated by color block below the densitree plot. To obtain permutation test-based statistical significance (P value) for a motif's hits within a sample, the summed mDELTA alignment scores across these hits of the focal motif were compared with its expectation as assessed by 1000 random trees of the sample (terminal cells shuffled). The six results for the hits of the two motifs in the three samples were shown respectively for each color block below the densitree plot. The probability of a summed alignment score from a randomized tree being larger than the observation (black arrow) is indicated by the P values in the panel.

dia consists of cSFDM containing 3 μm CHIR99021 (Tocris, No.4423), 10 ng mL^{-1} rhBMP4 (R&D, 314-BP-010), and 100 nm retinoic acid (RA, Sigma, No. R2625). At day 15 of differentiation, single-cell suspensions were prepared by incubating the cells at 37 $^{\circ}\text{C}$ in 0.05% trypsin-EDTA (Gibco, 25200056) for 7–15 min. The cells were then washed in media containing 10% fetal bovine serum (FBS, ThermoFisher), centrifuged at 300 g for 5 min, and resuspended in sort buffer containing Hank's Balanced Salt Solution (ThermoFisher), 2% FBS, and 10 μm Y-27632. The efficiency of differentiation into NKX2-1⁺ lung progenitors was evaluated either by flow cytometry for NKX2-1-GFP reporter expression, or expression of surrogate cell surface markers CD47^{hi}/CD26^{lo}. Cells were subsequently stained with CD47-PerCP/Cy5.5 and CD26-PE antibodies (Anti-human CD47 PerCP/Cy5.5 conjugate, Biolegend, Cat#323110, 1:200; Anti-human CD26 PE conjugate, Biolegend, Cat#302705, 1:200; PE mouse IgG1 isotype, Biolegend, Cat#400113, 1:200, PerCP/Cy5-5 mouse IgG1 isotype, Biolegend, Cat#400149, 1:200) for 30 min at 4 $^{\circ}\text{C}$, washed with PBS, and resuspended in sort buffer based on the method of Hawkins and Kotton.^[66] Cells were filtered through a 40 μm strainer (Falcon) prior to sorting. The CD47^{hi}/CD26^{lo} cell population was sorted on a high-speed cell sorter (MoFlo Astrios EQs) and resuspended in undiluted growth factor-reduced 3D matrigel (Corning 356230) at a dilution of 20–50 cells μL^{-1} , with droplets ranging in size from 20 μL (in 96 well plate) to 1 ml (in 10 cm dish). Cells in 3D matrigel suspension were incubated at 37 $^{\circ}\text{C}$ for 20–30 min, followed by the addition of warm media. The differentiation into distal/alveolar cells after day 15 was performed in “CK+DCI” medium, consisting of cSFDM base, with 3 μm CHIR (Tocris, No.4423), 10 ng mL^{-1} rhKGF (R&D, No.251-KG-010) (CK), and 50 nm dexamethasone (Sigma, No. D4902), 0.1 mM 8-Bromoadenosine 3',5'-cyclic monophosphate sodium salt (Sigma, No.B7880) and 0.1 mM 3-Isobutyl-1-methylxanthine (IBMX; Sigma, No.I5879) (DCI). Immediately after replating cells on day 15, 10 μm Y-27632 was added to the medium for 24 h. Upon replating on day 15, alveolospheres developed in 3D Matrigel culture outgrowth within 3–7 days, and were maintained in CK+DCI media for weeks. These spheres were analyzed by Z stack live images of alveolospheres taken and processed on the Leica DMI8 fluorescence microscope.

Directed Differentiation Followed by Simultaneous Assessment of Single-Cell Transcriptomes and Cell Lineage Tree: Based on the results from the full directed differentiation experiment above, it is aimed to evaluate single-cell transcriptomes and CLTs simultaneously for directed differentiation from hESCs to PLP, a stage at which the colony had <10000 cells, allowing to sample a large proportion of cells. To prepare suitable ancestor hESCs, the cell colonies outgrowth after 5–7 days, plated in 96-well dishes with microscopic selection for GFP⁺ mCherry⁺, were digested with GCDR to form ≈ 50 μm aggregates, and cultured in mTesR media until day 5. Combining selection and induction by Dox (1.0 $\mu\text{g mL}^{-1}$) and puro (1.0 $\mu\text{g mL}^{-1}$) from day 5 to day 7, the normally survived GFP⁺ mCherry⁺ colonies were capable of Cas9 expression and marked by primary editing events (to distinguish ancestor cells), as confirmed by DNA extraction and barcode PCR and sanger sequencing. The cell colonies with primary editing events were digested by GCDR for cell counting (≈ 4000 cells) and resuspended at a density of 10 cells μL^{-1} . One microliters cell suspension was added into each well of 96-well dishes plated with 1:10 diluted Matrigel (Corning, No.354277) for culture in mTesR media with ClonR (10:1) (Stemcell, No.05888) added in the first 48 h to promote the survival of very few stem cell. Directed differentiation was then initiated by applying both Dox (1.0 $\mu\text{g mL}^{-1}$, for editing the lineage barcode) and the STEMdiff Definitive Endoderm Kit to the normally survived colonies. Later stages of directed differentiation followed the differentiation protocols described above, with the exception that it was stopped on the tenth day after its initiation (Figure S1B, Supporting Information). Finally, colonies with intermediate size (≈ 5000 cells as approximated by colony size and cell counts) and $\geq 50\%$ GFP⁺ mCherry⁺ cells were digested with 0.05% trypsin-EDTA for 1 min at 37 $^{\circ}\text{C}$, washed in PBS containing 10% fetal bovine serum (FBS, ThermoFisher), centrifuged at 500 g for 5 min, and resuspended in single cell resuspension buffer containing PBS and 0.04% BSA. Using the standard 10x Chromium protocol, cDNA libraries were prepared from these single cell suspensions. Each cDNA library was split into two halves, with the first half subjected to conventional RNA-seq for single-cell transcrip-

omes, and the other half subjected to amplification of the lineage barcode followed by PacBio Sequel-based HiFi sequencing of the lineage barcode (Figure 1A).

Analysis of scRNA-seq: Following the 10x Genomics official guidelines, the Cell Ranger^[67] pipeline was used to map raw reads to the human reference genome (GRCh38) by STAR^[68] and obtained the read counts for each gene. Using Seurat v3.2.1,^[69] cells with <10% mitochondrial reads and >200 expressing unique features detected were retained. Then highly variable genes were detected by Single-cell Orientation Tracing (SOT),^[70] which were then subjected to Principle Component Analysis, followed by batch effect correction by Harmony.^[71] Then cells were clustered based on the cell-cell distance calculated by FindNeighbors and FindClusters using the Harmony-normalized matrix of gene expression. Then, runUMAP (default parameters except for “resolution” set to 0.6) was used for visualization and FindAllMarkers to obtain differentially expressed genes (DEGs) among clusters. To identify cell types, microarray data were downloaded from Gene Expression Omnibus (GEO),^[32,72] and DEGs (Wilcoxon Rank Sum test, $P < 0.01$) were detected for different stages of differentiation toward PLP. The clusters were scored based on the average expression and numbers of expressed stage-specific DEGs. Finally, 12 cell clusters were named based on the inferred order of appearance in the differentiation progress. To confirm these cell clusters were not biased by the chosen clustering method, another algorithm called PHATE was applied using its default parameters.^[73] Nine of the ten PHATE clusters were uniquely matched (i.e., the majority of the cells) to one of the currently defined cell clusters in Figure 1C, suggesting that these results are robust to the choice of cell clustering method.

Construction of Cell Lineage Trees: Based on the PacBio HiFi sequencing results, the CLTs were built and assessed for quality using PacBio HiFi reads following the previous pipeline.^[31] Briefly, using HiFi-seq raw sequences, consensus sequences were called separately from positive and negative strand subreads from each zero-mode waveguide (ZMW). Only consensus sequences with at least three subreads and identifiable barcode primers (Table S8, Supporting Information, allowing up to two mismatches) were reserved. From the consensus sequences, 10x cell barcodes and UMIs were extracted and matched to those from scRNA-seq, with one mismatch allowed. Lineage barcode sequences were then extracted from the consensus sequences, grouped by identical cell barcode and UMI, then merged by MUSCLE alignment followed by selecting the nucleotide with the highest frequency at each site. After MUSCLE alignment of the merged sequence to the reference lineage barcode, the editing events were called.^[61] Then, for each lineage barcode allele from the same cell, the frequency was calculated as the total number of UMIs of the allele and its ancestral allele. Here, the ancestral allele of a specific allele was defined as any allele in which the observed editing events were a subset of the editing events in the focal allele. Finally, the lineage barcode allele of a cell was defined as the allele with the highest frequency, prioritizing the alleles with more editing events if the frequencies were equal. For each sample, all cells with a lineage barcode and a single-cell transcriptome were used to construct a multifurcating lineage tree based on the lineage barcode using the maximum likelihood (ML) method implemented by the IQ-TREE LG model.^[74]

Transcriptome Divergence Among Cell Type Clusters: To elucidate the transcriptomic divergence among the observed clusters in the context of the directed differentiation toward PLP, stage-specific DEGs were extracted with the top 10% most extreme fold-change relative to other stages (Figure 2A, using microarray data^[32] mentioned above), and identified the Gene Ontology terms enriched (BH-adjusted $P < 0.05$, Fisher's exact test) with these stage-specific DEGs. After eliminating GO terms that had very few expressed genes, 179 GO terms (Table S7, Supporting Information) were focused on. For each cell, the activities of the specific cellular functions represented by these GO terms were estimated by the AddModuleScore function of Seurat, which basically calculated the average Z-score transformed expression levels of all genes annotated by the GO term. All cells within a cluster were then combined to determine the average activity of the GO term for the cluster (Figure 2B).

Transcriptome Divergence Among Sub-CLTs: As for the divergence among sub-CLTs, estimation of pseudotime was conducted via

Monocle^[34] with all cells on differentiating CLTs pooled together. After Principal Component Analysis of all cells from all samples combined, the transcriptomic divergence (D_T) between any two cells was quantified by one minus Pearson's Correlation Coefficient of the top 100 principal components. The developmental potential of an ancestor cell (an internal node on the CLT) was then calculated by the summed squared D_T of all pairs of its descendant cells. The reduction of developmental potential (Δ_{DP}) during the growth of an internal node to its daughter nodes was calculated by the focal internal node's Δ_{DP} subtracted by the summed Δ_{DP} of all its daughter nodes (Figure 2D). The statistical significance of an observed Δ_{DP} was estimated by contrasting the observation with its null distribution generated by random assignment of single-cell transcriptomes from hESC samples to the focal CLT (Figure 2D). It is emphasized here that the null distribution should be estimated by the single-cell transcriptomes from the non-differentiating hESC sample, since using those from the differentiating CBRAD5 samples would introduce actual divergence into the null and thus lead to an underestimated statistical significance. It was also worth noting that this method was very similar to the commonly used nonparametric method of permutational multivariate analysis of variance (PERMANOVA^[75]), except that Pearson's correlation-based divergence replaces the distance-based divergence used in canonical PERMANOVA, as the correlation-based metric had consistently been shown to result in superior performance for single-cell transcriptomes.^[76,77] This PERMANOVA-based method was also applied to subsets of genes within the transcriptome. For example, only genes annotated with a specific GO term (Table S7, Supporting Information) were used. A significant divergence for a specific GO term did not necessarily indicate a significant divergence in the whole transcriptome, since genes annotated with the GO term might had a small effect on the transcriptome as a whole. As a result, internal nodes with transcriptomic divergence did not necessarily represent a larger fraction than nodes with divergence on a specific GO term.

In order to perform a retrospective analysis of divergence progression, a normalized temporal scale was needed that was comparable across samples. In theory, this scale could be derived from the mutation rate of the lineage barcode and/or the topological depth of a node (i.e., the number of nodes between the root and the focal node). Considering the variability in Cas9 editing efficiency over barcodes, as well as long inter-site deletions, the mutation rate-based scale was discarded. For the topological depth scale, due to both biological and experimental stochasticity, the reconstructed CLTs and their nodes had very different depths, despite the fact that they were supposed to correspond to the ten-day directed differentiation. Assuming that the internal nodes were evenly sampled on all root-to-tip paths throughout the CLT, the actual depth of a node should be reflected equally by its depth from the root and (indirectly) by the depth from the focal node to its descendent tips. Based on this logic, the normalized depth of a node was defined as $d = (d_r/d_t + (1 - d_s/d_t)) / 2$, where d_r is the focal node's depth from root, d_t is the max depth found in the CLT, and the d_s is the max depth from the focal node to its descendent tips (Figure S5A, Supporting Information). Here, via division by d_t , all depths were scaled from 0 to 1, with 0 being the root and 1 being the tips with maximal raw depth within the CLT.

Transcriptional Memory Index: Previously proposed methods^[29,39] were followed to calculate transcriptional memory index. In each cell type and for each gene expressed in >10% of cells of this type, the CV of the expression levels was calculated among all terminal cells of this type within a sub-CLT (containing at least two cells of this type). The minimal CV among all sub-CLTs, i.e., $\min(CV)$, was then used to represent the expression variability of the focal gene in this cell type. It was also calculated for each of 1000 randomized CLTs created by reassigning all cells of the same type to a new lineage position that was originally occupied by the same cell type. These 1000 $\min(CV_{Random})$ from randomized CLTs were averaged, i.e., $\text{mean}(\min(CV_{Random}))$, to yield a null expectation for the observed $\min(CV)$. Finally, the memory index was defined as $M = (\min(CV) - \text{mean}(\min(CV_{Random}))) / \min(CV)$. Note that the final division by $\min(CV)$ was different from the previously defined memory index,^[29,39] but allows comparisons between genes with very different baseline CVs or expression levels.

To test the hypothesized role of transcription factors in mediating transcriptional memory, lists of gene sets responsive to perturbations of individual transcription factors ("TF_Perturbations_Followed_By_Expression" in Enrichr^[40]) were obtained. The genes with highest memory indices (top 10% across all cell types) were assessed for enrichment in each of these TF-responsive gene sets using Enrichr.^[40] The "combined score" (Figure 3E) was calculated by Enrichr, which takes into account both the statistical significance and the magnitude of enrichment (combined score of enrichment $c = \log(p) * o$, where p is the P value from Fisher's exact test and o is the odds ratio of the enrichment^[40]). Similar analyses were conducted on other gene sets (Figure S6C,D, Supporting Information) to detect enrichment for GO terms.

Composition of Terminal Cell Types Compared Among Sub-CLTs and the Full CLTs: To compare the terminal cell type composition of one sub-CLT with its expectation, a 2-by- n contingency table was constructed for the n cell types appearing in the entire CLT. The first row of the contingency table lists the observed count of terminal cells for each cell type within the focal sub-CLT. The second row of the table lists the expected count of each cell type as determined by the fractional cell type composition of the entire CLT multiplied by the size of the focal sub-CLT. Then $\chi^2 = \sum_{i=1}^n (O_i - E_i)^2 / E_i$ was calculated for the focal sub-CLT, where O_i and E_i are the observed and expected count for cell type i . Then χ^2 values from all sub-CLTs with roots of normalized depth < 0.7 (because internal nodes closer to terminal cells produce sub-CLTs that are too small for meaningful statistics) were summed up to represent the diversity of cell type compositions among sub-CLTs (x axis of Figure 4A–C). In other words, a small summed χ^2 indicates uniform/stereotyped composition of cell types among sub-CLTs. To assess the null distribution of the summed χ^2 , 1000 control CLTs were created by randomly reassigning all cells on the tree to a different terminal node, while keeping the topology of the tree unchanged.

Robustness of Random versus Stereotyped Development: Without loss of generality, a functional unit was defined as consisting of four cell types, namely α , β , γ , and δ , in a 1:1:2:4 ratio. 1000 binary CLTs were simulated, each consisting of 1024 terminal cells (128 α cells, 128 β cells, 256 γ cells, 512 δ cells) generated through ten cell cycles, under two developmental models. The first "random" model randomly assigns the four types of cells onto the tips of the tree. A second "stereotyped" model strictly assigns α , β , γ , and δ cells in a 1:1:2:4 ratio onto each sub-CLT consisting of eight tips (three cell cycles). A predefined fraction (0.001, 0.005, 0.01, 0.05, or 0.1, as on x axis of Figure 4F) of the 2047 (1024 terminal and 1023 internal) cells were chosen and removed along with all their descendent cells to mimic random cell deaths. Assuming sufficient cell migration to allow formation of the functional unit as long as there are enough terminal cells of the proper type, the robustness is thus quantified by the number of functional units that could be formed by all terminal cells surviving cell deaths. A simple example shown in Figure 4E.

Comparison and Alignment of Sub-CLTs by mDELTA: Let vectors/nodes be denoted as V and edges connecting nodes as E . Given a query tree $Q = (V, E)$ and a subject tree $S = (V', E')$, an isomorphic alignment is a bijection $A: V \rightarrow V'$, such that for every pair of nodes with $v, u \in V$, it have $(v, u) \in E \Leftrightarrow (A(v), A(u)) \in E'$. Based on two types of biologically informed tree editing operations, namely pruning and merging (see Text S1, Supporting Information), a homeomorphic subtree alignment A between Q and S is defined as an isomorphic alignment between Q' and S' , where Q' is the result of zero or more pruning and merging in Q , and S' is the result of zero or more pruning and merging in S . Here all the pruning in Q and S are collectively denoted as $\pi(A)$, and all merging in Q and S are collectively denoted as $\mu(A)$. If the alignment score between two nodes $v \in V$ and $v' \in V'$ is further denoted as $a(v, v')$, the cost for pruning a subtree \hat{T} as $p(\hat{T})$, the cost for merging an internal node \hat{v} with its mother node as $m(\hat{v})$. The score of a homeomorphic subtree alignment A between Q and S can then be expressed as

$$w(Q, S, A) = \sum_{(v, v') \in A} a(v, v') - \sum_{\hat{T} \in \pi(A)} p(\hat{T}) - \sum_{\hat{v} \in \mu(A)} m(\hat{v}) \quad (1)$$

The algorithm of mDELTA find the optimal A (with optimal/highest possible w) given Q , S , a , p , and m by a dynamic programming procedure. a was defined based on similarity of single-cell transcriptomes, p based on the number of pruned terminal cells, and m based on the number of merged internal nodes. Detailed computational procedures of mDELTA can be found in Text S1 (Supporting Information).

Heritability of Quantitative Traits in the CLT: In order to gauge the heritability of quantitative traits on the CLT, the correlation between the relatedness and the phenotypic divergence of a pair of nodes was calculated. When the relatedness was defined by genomic relatedness like DNA sequence identity, this analysis was the same as the classic statistical genetics method called Haseman-Elston Regression.^[52] Thus, the correlation coefficient from this analysis was considered a proxy for phenotypic heritability among nodes on the CLT. However, it should be emphasized that since the DNA sequences of all cells in this dataset were presumably nearly identical, the relatedness between nodes was therefore defined by their distance on the CLT instead (see below), and the resulting correlation coefficients could not be interpreted as traditional heritability as they are in Haseman-Elston Regression. Specifically, the relatedness between any two nodes on the CLT inversely was defined by the number of cell divisions separating them, which was then estimated by contrasting the number of their descendent cells with the number of descendent cells of their latest common ancestor. Following previous Haseman-Elston Regression applications,^[53] the relatedness between nodes was then scaled so that the mean relatedness between any pair of nodes was 0 and the maximal relatedness was 1. As such scaling is equivalent to calculating relatedness relative to a different population,^[53] comparing the heritability of one trait relative to that of another trait would not be affected as long as both traits were analyzed in the same focal population (the focal CLT). On the phenotype side, two quantitative traits were examined, the single-cell transcriptomes of terminal nodes and the descendent cell type compositions of internal nodes. Here, the single-cell transcriptomes of terminal nodes were first processed by Principle Component Analyses, then all principle components of a cell were used to represent its transcriptome. As for the descendant cell type compositions of an internal node, each internal node was represented by a vector comprising M elements, where M is the total number of cell types identified in the dataset, and each element represents the percentage of descendent cells of that type. The phenotypic divergence between two nodes was calculated as the Euclidian distance between the multidimensional quantitative traits. Lastly, the Spearman's Correlation Coefficient were reported between the relatedness and the phenotypic divergence between all relevant node pairs in Figure S7 (Supporting Information) as a proxy for the heritability of quantitative traits.

Ethical Approval: All experiments were conducted in vitro using previously established human embryonic stem cell line H9. There were, therefore, no ethical concerns related to the use of human or animal subjects in this study. The research conducted adheres to the highest standards of scientific ethics and was in compliance with all relevant regulations and guidelines regarding the use of human embryonic stem cells.

Statistical Analysis: Due to the computational complexity and diversity of the presented analyses, methods of statistical analysis, including preprocessing and presenting data, as well as statistical methods used, had been described in-context in Experimental Section, figure legends and Text S1 (Supporting Information). These analyses were conducted using custom codes developed in R and Python, available at https://github.com/ChenJy0212/mdelta_full (the mDELTA algorithm) and <https://github.com/ZhangxyOk/Stereotyped-CLT> (all other custom codes).

Supporting Information

Supporting Information is available from the Wiley Online Library or from the author.

Acknowledgements

This work was supported by the National Key R&D Program of China (grant numbers 2021YFA1302500 and 2021YFF1200904 to J.-R. Y.), the National Natural Science Foundation of China (grant numbers 3212022 and 32361133555 to J.-R. Y.), Fundamental Research Funds for the Central Universities, Sun Yat-sen University (grant number 23ptpy70 to X.Z.), the Research Grants Council of the Hong Kong Special Administrative Region, China (grant number 21206223 to X.W.), and the Russian Science Foundation (grant number 24-44-00099 to V.A.L.)

Conflict of Interest

The authors declare no conflict of interest.

Author Contributions

X.Z., Z.L., and J.C. contributed equally to this work. J.-R.Y. conceived the idea, and designed and supervised the study. X.Z., Z.L., W.Y., X.H., P.W., F.C., Z.Z., and X.C. conducted experiments and acquired data. X.W., V.A.L., L.Y.R., X.C., and J.-R.Y. contributed new devices/reagents/analytic tools. X.Z., Z.L., J.C., W.Y., P.W., F.C., X.H., C.R., Y.S., X.C., and J.-R.Y. analyzed the data. X.Z., Z.L., and J.-R.Y. wrote the paper with inputs from all the authors.

Data Availability Statement

The new data generated in this study were deposited to NCBI BioProjects under accession number PRJNA1099925.

Keywords

cell lineage tree, developmental robustness, human embryonic stem cell, lung primordial progenitor

Received: June 5, 2024

Revised: March 12, 2025

Published online:

- [1] C. H. Waddington, *Nature* **1942**, 150, 563.
- [2] H. Kitano, *Nat. Rev. Genet.* **2004**, 5, 826.
- [3] H. F. Nijhout, *BioEssays* **2002**, 24, 553.
- [4] M. L. Siegal, J. Y. Leu, *Annu. Rev. Ecol. Evol. Syst.* **2014**, 45, 495.
- [5] M. A. Felix, A. Wagner, *Heredity (Edinb)* **2008**, 100, 132.
- [6] G. Gibson, K. A. Lacey, *Annu. Rev. Genet.* **2020**, 54, 189.
- [7] S. L. Rutherford, S. Lindquist, *Nature* **1998**, 396, 336.
- [8] E. Hornstein, N. Shomron, *Nat. Genet.* **2006**, 38, S20.
- [9] C. I. Wu, Y. Shen, T. Tang, *Genome Res.* **2009**, 19, 734.
- [10] V. Siciliano, I. Garzilli, C. Fracassi, S. Criscuolo, S. Ventre, D. di Bernardo, *Nat. Commun.* **2013**, 4, 2364.
- [11] S. F. Levy, M. L. Siegal, *PLoS Biol.* **2008**, 6, 264.
- [12] N. Mo, X. Zhang, W. Shi, G. Yu, X. Chen, J.-R. Yang, *Mol. Biol. Evol.* **2021**, 38, 1874.
- [13] H. El-Samad, *Cell Syst.* **2021**, 12, 477.
- [14] R. Kafri, M. Levy, Y. Pilpel, *Proc. Natl. Acad. Sci. USA* **2006**, 103, 11653.
- [15] L. Xiao, D. Fan, H. Qi, Y. Cong, Z. Du, *Cell Syst.* **2022**, 13, 619.
- [16] R. M. Green, J. L. Fish, N. M. Young, F. J. Smith, B. Roberts, K. Dolan, I. Choi, C. L. Leach, P. Gordon, J. M. Cheverud, C. C. Roseman, T. J. Williams, R. S. Marcucio, B. Hallgrímsson, *Nat. Commun.* **2017**, 8, 1970.

- [17] J. R. Yang, S. Ruan, J. Zhang, *PLoS Genet.* **2014**, *10*, 1004501.
- [18] D. E. Wagner, A. M. Klein, *Nat. Rev. Genet.* **2020**, *21*, 410.
- [19] Z. Li, W. Yang, P. Wu, Y. Shan, X. Zhang, F. Chen, J. Yang, J.-R. Yang, *J. Genet. Genomics* **2024**, *51*, 35.
- [20] C. H. Waddington, *The Strategy of the Genes: A Discussion of Some Aspects of Theoretical Biology*, George Allen & Unwin, London **1957**.
- [21] E. Pujadas, A. P. Feinberg, *Cell* **2012**, *148*, 1123.
- [22] J. E. Sulston, E. Schierenberg, J. G. White, J. N. Thomson, *Dev. Biol.* **1983**, *100*, 64.
- [23] B. Raj, J. A. Gagnon, A. F. Schier, *Nat. Protoc.* **2018**, *13*, 2685.
- [24] D. P. Furman, T. A. Bukharina, *Curr. Genomics* **2008**, *9*, 312.
- [25] U. Koch, R. Lehal, F. Radtke, *Development* **2013**, *140*, 689.
- [26] G. Kilimnik, J. Jo, V. Periwal, M. C. Zielinski, M. Hara, *Islets* **2012**, *4*, 167.
- [27] S. J. Salipante, A. Kas, E. McMonagle, M. S. Horwitz, *Evol. Dev.* **2010**, *12*, 84.
- [28] D. J. Anderson, F. M. Pauler, A. McKenna, J. Shendure, S. Hippenmeyer, M. S. Horwitz, *Cell Syst.* **2022**, *13*, 438.
- [29] N. W. Hughes, Y. Qu, J. Zhang, W. Tang, J. Pierce, C. Wang, A. Agrawal, M. Morri, N. Neff, M. M. Winslow, M. Wang, L. Cong, *Mol. Cell* **2022**, *82*, 3103.
- [30] A. Jacob, M. Morley, F. Hawkins, K. B. McCauley, J. C. Jean, H. Heins, C.-L. Na, T. E. Weaver, M. Vedaie, K. Hurley, A. Hinds, S. J. Russo, S. Kook, W. Zacharias, M. Ochs, K. Traber, L. J. Quinton, A. Crane, B. R. Davis, F. V. White, J. Wambach, J. A. Whitsett, F. S. Cole, E. E. Morrissey, S. H. Guttentag, M. F. Beers, D. N. Kotton, *Cell Stem Cell* **2017**, *21*, 472.
- [31] F. Chen, Z. Li, X. Zhang, P. Wu, W. Yang, J. Yang, X. Chen, J.-R. Yang, *Mol. Biol. Evol.* **2023**, *40*, msad113.
- [32] F. Hawkins, P. Kramer, A. Jacob, I. Driver, D. C. Thomas, K. B. McCauley, N. Skvir, A. M. Crane, A. A. Kurmann, A. N. Hollenberg, S. Nguyen, B. G. Wong, A. S. Khalil, S. X. L. Huang, S. Guttentag, J. R. Rock, J. M. Shannon, B. R. Davis, D. N. Kotton, *J. Clin. Invest.* **2017**, *127*, 2277.
- [33] Z. He, A. Maynard, A. Jain, T. Gerber, R. Petri, H.-C. Lin, M. Tra Santel, K. Ly, J.-S. Dupré, L. Sidow, F. Sanchis Calleja, S. M. J. Jansen, S. Riesenberger, J. G. Camp, B. Treutlein, *Nat. Methods* **2022**, *19*, 90.
- [34] C. Trapnell, D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J. Livak, T. S. Mikkelsen, J. L. Rinn, *Nat. Biotechnol.* **2014**, *32*, 381.
- [35] R. Bonasio, S. Tu, D. Reinberg, *Science* **2010**, *330*, 612.
- [36] S. M. Shaffer, B. L. Emert, R. A. Reyes Hueros, C. Cote, G. Harmange, D. L. Schaff, A. E. Sizemore, R. Gupte, E. Torre, A. Singh, D. S. Bassett, A. Raj, *Cell* **2020**, *182*, 947.
- [37] Y. Dublanche, K. Michalodimitrakis, N. Kummerer, M. Foglierini, L. Serrano, *Mol. Syst. Biol.* **2006**, *2*, 41.
- [38] E. Sharon, D. van Dijk, Y. Kalma, L. Keren, O. Manor, Z. Yakhini, E. Segal, *Genome Res.* **2014**, *24*, 1698.
- [39] A. S. Eisele, M. Tarbier, A. A. Dormann, V. Pelechano, D. M. Suter, *bioRxiv* **2022**, 508646.
- [40] M. V. Kuleshov, M. R. Jones, A. D. Rouillard, N. F. Fernandez, Q. Duan, Z. Wang, S. Koplev, S. L. Jenkins, K. M. Jagodnik, A. Lachmann, M. G. McDermott, C. D. Monteiro, G. W. Gundersen, A. Ma'ayan, *Nucleic Acids Res.* **2016**, *44*, W90.
- [41] G. Pan, J. A. Thomson, *Cell Res.* **2007**, *17*, 42.
- [42] A. Tan, R. Prasad, E. H. Jho, *Cell Death Dis.* **2021**, *12*, 343.
- [43] J. Lv, S. Meng, Q. Gu, R. Zheng, X. Gao, J.-D. Kim, M. Chen, B. Xia, Y. Zuo, S. Zhu, D. Zhao, Y. Li, G. Wang, X. Wang, Q. Meng, Q. Cao, J. P. Cooke, L. Fang, K. Chen, L. Zhang, *Nat. Commun.* **2023**, *14*, 2390.
- [44] D. C. Liberti, W. A. Liberti III, M. M. Kremp, I. J. Penkala, F. L. Cardenas-Diaz, M. P. Morley, A. Babu, S. Zhou, R. J. Fernandez III, E. E. Morrissey, *Dev. Cell* **2022**, *57*, 1742.
- [45] Y. Zhang, N. Rath, S. Hannehalli, Z. Wang, T. Cappola, S. Kimura, E. Atochina-Vasserman, M. M. Lu, M. F. Beers, E. E. Morrissey, *Development* **2007**, *134*, 189.
- [46] J. A. Heslop, B. Pournasr, J. T. Liu, S. A. Duncan, *Cell Rep.* **2021**, *35*, 109145.
- [47] W. Liu, G. Chen, *Cell. Mol. Life Sci.* **2021**, *78*, 8097.
- [48] C. Song, F. Xu, Z. Ren, Y. Zhang, Ya Meng, Y. Yang, S. Lingadahalli, E. Cheung, G. Li, W. Liu, J. Wan, Y. Zhao, G. Chen, *Stem Cell Rep.* **2019**, *13*, 338.
- [49] P. Fojtik, M. Senfluk, K. Holomkova, A. Salykin, J. Gregorova, P. Smak, O. Pes, J. Raska, M. Stetkova, P. Skladal, M. Sedlackova, A. Hampl, D. Bohaciakova, S. Uldrijan, V. Rotrekl, *bioRxiv* **2023**, 524871.
- [50] M. Yuan, X. Yang, J. Lin, X. Cao, F. Chen, X. Zhang, Z. Li, G. Zheng, X. Wang, X. Chen, J.-R. Yang, *iScience* **2020**, *23*, 101273.
- [51] J. Skommer, T. Brittain, S. Raychaudhuri, *Apoptosis* **2010**, *15*, 1223.
- [52] J. K. Haseman, R. C. Elston, *Behav. Genet.* **1972**, *2*, 3.
- [53] J. Yang, B. Benyamin, B. P. McEvoy, S. Gordon, A. K. Henders, D. R. Nyholt, P. A. Madden, A. C. Heath, N. G. Martin, G. W. Montgomery, M. E. Goddard, P. M. Visscher, *Nat. Genet.* **2010**, *42*, 565.
- [54] M. Tran, A. Askary, M. B. Elowitz, *Dev. Cell* **2024**, *59*, 812.
- [55] M. Müller, Y. Kohl, A. Germann, S. Wagner, H. Zimmermann, H. von Briesen, *In Vitro Models* **2023**, *2*, 249.
- [56] C. Peaslee, C. Esteva-Font, T. Su, A. Munoz-Howell, C. C. Duwaerts, Z. Liu, S. Rao, K. Liu, M. Medina, J. B. Sneddon, J. J. Maher, A. N. Mattis, *Hepatology* **2021**, *74*, 2102.
- [57] A. Jacob, M. Vedaie, D. A. Roberts, D. C. Thomas, C. Villacorta-Martin, K.-D. Alysandratos, F. Hawkins, D. N. Kotton, *Nat. Protoc.* **2019**, *14*, 3303.
- [58] Y. Li, S. Chen, W. Liu, D. Zhao, Y. Gao, S. Hu, H. Liu, Y. Li, L. Qu, X. Liu, *Nat. Commun.* **2024**, *15*, 358.
- [59] J. I. Murray, T. J. Boyle, E. Preston, D. Vafeados, B. Mericle, P. Weisdepp, Z. Zhao, Z. Bao, M. Boeck, R. H. Waterston, *Genome Res.* **2012**, *22*, 1282.
- [60] K. Ouardini, R. Lopez, M. G. Jones, S. Prillo, R. Zhang, M. I. Jordan, N. Yosef, *bioRxiv* **2021**, 446021.
- [61] A. McKenna, G. M. Findlay, J. A. Gagnon, M. S. Horwitz, A. F. Schier, J. Shendure, *Science* **2016**, *353*, aaf7907.
- [62] J. G. Doench, N. Fusi, M. Sullender, M. Hegde, E. W. Vaimberg, K. F. Donovan, I. Smith, Z. Tothova, C. Wilen, R. Orchard, H. W. Virgin, J. Listgarten, D. E. Root, *Nat. Biotechnol.* **2016**, *34*, 184.
- [63] T. Magoc, S. L. Salzberg, *Bioinformatics* **2011**, *27*, 2957.
- [64] B. Langmead, S. L. Salzberg, *Nat. Methods* **2012**, *9*, 357.
- [65] R. C. Edgar, *Nucleic Acids Res.* **2004**, *32*, 1792.
- [66] A. Sahabian, J. Dahlmann, U. Martin, R. Olmer, *Nat. Protoc.* **2021**, *16*, 1581.
- [67] G. X. Y. Zheng, J. M. Terry, P. Belgrader, P. Ryvkin, Z. W. Bent, R. Wilson, S. B. Ziraldo, T. D. Wheeler, G. P. McDermott, J. Zhu, M. T. Gregory, J. Shuga, L. Montesclaros, J. G. Underwood, D. A. Masquelier, S. Y. Nishimura, M. Schnall-Levin, P. W. Wyatt, C. M. Hindson, R. Bharadwaj, A. Wong, K. D. Ness, L. W. Beppu, H. J. Deeg, C. McFarland, K. R. Loeb, W. J. Valente, N. G. Ericson, E. A. Stevens, J. P. Radich, et al., *Nat. Commun.* **2017**, *8*, 14049.
- [68] A. Dobin, C. A. Davis, F. Schlesinger, J. Drenkow, C. Zaleski, S. Jha, P. Batut, M. Chaisson, T. R. Gingeras, *Bioinformatics* **2013**, *29*, 15.
- [69] Y. Hao, S. Hao, E. Andersen-Nissen, W. M. Mauck, S. Zheng, A. Butler, M. J. Lee, A. J. Wilk, C. Darby, M. Zager, P. Hoffman, M. Stoeckius, E. Papalexi, E. P. Mimitou, J. Jain, A. Srivastava, T. Stuart, L. M. Fleming, B. Yeung, A. J. Rogers, J. M. McElrath, C. A. Blish, R. Gottardo, P. Smibert, R. Satija, *Cell* **2021**, *184*, 3573.

- [70] L. Guo, L. Lin, X. Wang, M. Gao, S. Cao, Y. Mai, F. Wu, J. Kuang, H. e Liu, J. Yang, S. Chu, H. Song, D. Li, Y. Liu, K. Wu, J. Liu, J. Wang, G. Pan, A. P. Hutchins, J. Liu, D. Pei, J. Chen, *Mol. Cell* **2019**, *73*, 815.
- [71] I. Korsunsky, N. Millard, J. Fan, K. Slowikowski, F. Zhang, K. Wei, Y. Baglaenko, M. Brenner, P.-R. Loh, S. Raychaudhuri, *Nat. Methods* **2019**, *16*, 1289.
- [72] T. Barrett, S. E. Wilhite, P. Ledoux, C. Evangelista, I. F. Kim, M. Tomashevsky, K. A. Marshall, K. H. Phillippy, P. M. Sherman, M. Holko, A. Yefanov, H. Lee, N. Zhang, C. L. Robertson, N. Serova, S. Davis, A. Soboleva, *Nucleic Acids Res.* **2013**, *41*, D991.
- [73] K. R. Moon, D. van Dijk, Z. Wang, S. Gigante, D. B. Burkhardt, W. S. Chen, K. Yim, A. V. D. Elzen, M. J. Hirn, R. R. Coifman, N. B. Ivanova, G. Wolf, S. Krishnaswamy, *Nat. Biotechnol.* **2019**, *37*, 1482.
- [74] C. Umkehrer, F. Holstein, L. Formenti, J. Jude, K. Froussios, T. Neumann, S. M. Cronin, L. Haas, J. J. Lipp, T. R. Burkard, M. Fellner, T. Wiesner, J. Zuber, A. C. Obenauf, *Nat. Biotechnol.* **2021**, *39*, 174.
- [75] M. J. Anderson, *Austral Ecol.* **2001**, *26*, 32.
- [76] B. Ozgode Yigin, G. Saygili, *Sci. Rep.* **2023**, *13*, 6567.
- [77] T. Kim, I. R. Chen, Y. Lin, A. Y.-Y. Wang, J. Y. H. Yang, P. Yang, *Brief Bioinf.* **2019**, *20*, 2316.