

1 **High density cell lineage tracing reveals polyclonal stereotyped**  
2 **sub-trees, a contributor to developmental robustness**

3 Xiaoyu Zhang<sup>1,#</sup>, Zizhang Li<sup>1,#</sup>, Jingyu Chen<sup>1,#</sup>, Wenjing Yang<sup>1</sup>, Xingxing He<sup>2</sup>, Peng Wu<sup>1</sup>, Feng  
4 Chen<sup>1</sup>, Ziwei Zhou<sup>1</sup>, Xiewen Wen<sup>2</sup>, Vassily A. Lyubetsky<sup>3,4</sup>, Leonid Yu. Rusin<sup>3</sup>, Xiaoshu Chen<sup>1,5,a</sup>,  
5 Jian-Rong Yang<sup>1,5,\*b</sup>,

6

7 <sup>1</sup> *Department of Genetics and Biomedical Informatics, Zhongshan School of Medicine, Sun Yat-sen*  
8 *University, Guangzhou 510080, China*

9 <sup>2</sup> *University Research Facility in 3D Printing, & State Key Laboratory of Ultra-precision Machining*  
10 *Technology, Dept. of ISE, the Hong Kong Polytechnic University, Hong Kong*

11 <sup>3</sup> *Kharkevich Institute for Information Transmission Problems Russian Academy Sciences, Moscow*  
12 *127051 Russia*

13 <sup>4</sup> *Department of Mathematical Logic and Theory of Algorithms, Faculty of Mechanics and*  
14 *Mathematics, Lomonosov Moscow State University, Moscow 119991 Russia*

15 <sup>5</sup> *Key Laboratory of Tropical Disease Control, Ministry of Education, Sun Yat-sen University,*  
16 *Guangzhou 510080, China*

17 <sup>#</sup>These authors contributed equally to this work.

18

19 <sup>a</sup>ORCID: 0000-0002-5779-5065

20 <sup>b</sup>ORCID: 0000-0002-7807-9455

21 \* Corresponding authors.

22 E-mail: yangjianrong@mail.sysu.edu.cn (Yang JR)

23

## 24 **Abstract**

25 Robust development is essential for multicellular organisms to maintain physiological stability in  
26 the face of environmental changes or perturbations. While various mechanisms contributing to  
27 developmental robustness have been identified at the subcellular level, those at the intercellular and  
28 tissue level remain underexplored. We approach this question using a well-established *in vitro*  
29 directed differentiation model known to recapitulate the *in vivo* development of lung progenitor cells  
30 from human embryonic stem cells. An integrated analysis of high-density cell lineage trees (CLTs)  
31 and single-cell transcriptomes of the differentiating colonies enabled the resolution of known cell  
32 types as well as their developmental hierarchies. Our dataset showed little support for the hypothesis  
33 that transcriptional memory contributes to robust development by constraining single-cell  
34 transcriptomes of closely related cells. We nevertheless observed stable terminal cell type  
35 compositions among many sub-clones. This feature enhances developmental robustness because the  
36 colony could retain a relatively stable cell type composition even if some sub-CLTs are abolished  
37 by necrosis. Furthermore, using a novel computational framework for CLT alignment, we found that  
38 many sub-clones are formed by sub-CLTs resembling each other in terms of both terminal cell type  
39 compositions and topological structures. The existence of such sub-CLTs resembling each other not  
40 only deepens our understanding of developmental robustness by demonstrating the existence of a  
41 stereotyped developmental program, but also suggests a novel perspective on the function of  
42 specific cell types within the context of stereotyped sub-CLTs, just as nucleotides are better  
43 understood in the context of sequence motifs.

## 44 **Keywords**

45 Developmental robustness, Cell lineage tree

46

## 47 **Introduction**

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

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

72       One of the main manifestations of developmental robustness is the generation of adequate  
73 numbers of cells of various types in an appropriate cellular composition, especially when they work  
74 together as a functional unit. For example, the *Drosophila* peripheral nervous system contains  
75 thousands of identical mechanosensory bristles<sup>24</sup>, each consisting of exactly one hair cell, one socket  
76 cell, one sheath cell and one neuron<sup>25</sup>. Another well-known example is the functional unit of the  
77 endocrine pancreas, the islet, which has been shown in mice to consist predominantly (~90%) of  $\beta$   
78 cells at the core and  $\alpha$  and  $\delta$  cells in the periphery<sup>26</sup>. To identify potential CLT characteristics that  
79 contributed to such a manifestation of developmental robustness, two CLT-based studies are  
80 particularly relevant. In the first, it was found that development of mammalian organs is preceded  
81 by significant mixing of multipotent progenitor cells<sup>27</sup>. Therefore, most organs have a polyclonal  
82 origin that ensures sufficient number of cells even some progenitors failed<sup>27</sup>. In the second, CLT of  
83 cortical development revealed stereotyped development giving rise to monophyletic clades of mixed

84 cell types<sup>28</sup>. On the basis of these observations, we hypothesized that the combination of polyclonal  
85 origin and stereotyped development facilitates the robust development of adequate numbers of cells  
86 with an appropriate cellular composition. It is imperative to note that as our hypothesis revolves  
87 around the above-mentioned functional units, CLTs with sufficient resolution (fraction of cells  
88 sampled) are essential, otherwise stereotyped development cannot be detected with only <1% cells  
89 sampled from each functional unit. In addition, a high resolution CLT would also reveal how  
90 stereotyped development occurs, such as mitotic-coupling versus population-coupling  
91 development<sup>18</sup> and whether epigenetic memory<sup>29</sup> plays a role.

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

## 109 **Results**

### 110 **Reconstructing high-density cell lineage trees for directed** 111 **differentiation of primordial lung progenitors**

112 We aimed to determine the CLT of embryonic stem cells undergoing *in vitro* directed  
113 differentiation towards lung progenitors according to a well-established protocol recapitulating *in*  
114 *vivo* development<sup>30</sup>. This *in vitro* model of directed differentiation was chosen for several reasons.  
115 First, cells cultured in a small petri dish have a relatively homogenous environment, so that  
116 transcriptome divergence caused by environmental factors, or phylogeny-independent convergence  
117 due to niche-specific signals is unlikely. Second, the development trajectory of embryonic stem cells

118 to the lung is well-known, such that the *in vitro* cell culture can be monitored to ensure that they  
119 closely mimic physiological situation. Indeed, our implementation of the protocol can reach the  
120 alveolar epithelial cells (AEC2s) fate after 20 days of directed differentiation (**Figure S1A** and  
121 **Video S1**). Third, *in vitro* culture allows us to induce Cas9 expression and therefore initiate the  
122 editing of the lineage barcode concurrently with the directed differentiation (**Figure S1B/C**). Last  
123 but not least, it allows better control over the number of cells within the colony assayed for single-  
124 cell transcriptomes and CLTs. In particular, our cell culture begins with ~10 hESCs and ends with ~  
125 5,000 cells on day 10 (**Figure S1B**), of which a relatively high percentage can be captured in  
126 downstream experimental pipelines of 10x Chromium. The ten-day directed differentiation covers  
127 three critical phases of lung development, including definitive endoderm (DE), anterior foregut  
128 endoderm (AFE) and NKX2-1<sup>+</sup> primordial lung progenitor (PLP)<sup>30</sup> (**Figure 1A**, **Figure S1A/B**).

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

143 We obtained single-cell transcriptomes of 3,576/4,400/1,456/5,659 cells respectively from  
144 three differentiating colonies CBRAD5-A1/G2/G11 and one parallel non-differentiating hESC  
145 colony, all of which appeared to have good quality (**Figure S2A/B**, **Table S2**). The UMAP clustering  
146 of the single-cell transcriptomes revealed a large fraction of cells from differentiating/CBRAD5  
147 colonies separated with those from hESC colonies, clearly indicating their differentiated cell states  
148 (**Figure 1B**). We identified 12 major functional clusters within the sampled cells (**Figure 1C**; See  
149 **Methods**). According to the average expression of pluripotent gene (*NANOG*, *POU5F1*), endoderm  
150 progenitor gene (*GATA6*) and lung progenitor gene (*NKX2-1*, *SHH*, *CD47*), these clusters were  
151 defined as NANOG<sup>hi</sup>POU5F1<sup>hi</sup> (C1), NANOG<sup>low</sup>POU5F1<sup>hi</sup> (C2), NANOG<sup>low</sup>POU5F1<sup>low</sup> (C3),  
152 NANOG<sup>hi/low</sup>POU5F1<sup>hi</sup> (C4), CD47<sup>hi</sup> (C5), CD47<sup>low</sup> (C6), GATA6<sup>hi</sup>SHH<sup>hi</sup>CD47<sup>low</sup> (C7),  
153 GATA6<sup>low</sup>NKX2-1<sup>neg</sup>SHH<sup>neg</sup>CD47<sup>neg</sup> (C8), GATA6<sup>hi</sup>NKX2-1<sup>hi</sup>CD47<sup>hi</sup> (C9), GATA6<sup>hi</sup> (C10). Below,  
154 they are also more broadly categorized into the less differentiated spontaneous state (R1 and R2) or  
155 pluripotent state (C1/C2/C3/C4), and the more differentiated progenitor state

156 (C5/C6/C7/C8/C9/C10). These clusters displayed transcriptomic states largely compatible with  
157 known cell types occurred during the directed differentiation<sup>32</sup> (**Figure 1D**, **Figure S2C**), and were  
158 differentially distributed between hESC and CBRAD5 samples (**Figure 1E**), thereby suggesting  
159 successfully induced differentiation and accurate measurement of single-cell transcriptomes. After  
160 confirming the sequencing quality of PacBio (**Table S3**, **Figure S2D**), the CLT of each sample was  
161 constructed based on the lineage barcode using maximum likelihood method (**Figure 1A/F**; See  
162 also **Methods**, **Figure S2E**, **Table S4/S5/S6**). The hierarchical population structures of the colonies  
163 were complex and intricate. In support of the accuracy of the CLT, cells more closely related to one  
164 another displayed more similar lineage barcode alleles (**Figure 1G**), and are more likely to share  
165 yet-to-decay transcripts of ancestral lineage barcode (**Figure 1H**). In conclusion, our experiment  
166 reliably captured the coarse-grained phylogenetic relationship of the cells within each colony.

## 167 **The cell lineage trees recapitulate key features of the transcriptome** 168 **divergence**

169 To better elucidate the divergence between the single-cell transcriptomes in the context of the  
170 observed clusters, we identified differentially expressed genes (DEGs) in previously published  
171 microarray-based transcriptome<sup>33</sup> data of samples from six timepoints of directed differentiation  
172 towards PLP (**Figure 2A**). Note here that despite being sampled on day12, the neural NKX2-1<sup>+</sup>  
173 transcriptome has been shown to be most similar to that of day0 hESCs<sup>33</sup>. The Gene Ontology terms  
174 enriched with these microarray-based stage-specific DEGs (**Table S7**) were then individually  
175 examined for overall activities in our single-cell transcriptomes by the member genes' average  
176 expression levels in each cluster (**Figure 2B**. See **Methods**). For pluripotent stage cells  
177 (C1/C2/C3/C4), significantly enhanced activities were found among GO terms enriched with DEGs  
178 of day 0/3 samples (including neural NKX2-1<sup>+</sup>)(**Figure 2B**). The same observations were made for  
179 progenitor stage cells C6/C10 in GO terms related to day3 samples, as well as C7/C9 cells in GO  
180 terms related to day6/day15 lung samples (**Figure 2B**). These results indicate that the single-cell  
181 transcriptomes recapitulated major differentiation stages of the *in vitro* PLP differentiation.

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

191 full range of transcriptional variation, an observation directly supports the transcriptional divergence  
192 among sub-CLTs.

193 For a more detailed analyses, we quantified the developmental potential of an internal node by  
194 the multivariate variance among its descendant single-cell transcriptomes, which then allowed us to  
195 perform PERMANOVA-based statistical tests (PERmutational Multivariate Analysis Of VAriance,  
196 see **Methods**) for the transcriptomic divergence. Briefly, by subtracting from the developmental  
197 potential of a focal node by the sum of the potentials of all its daughter nodes, we estimated the  
198 degree of divergence that occurred during the growth of the focal node (**Figure 2D**). Using the  
199 degree of divergence seen in the hESC sample as the null distribution, an average of ~65% internal  
200 nodes of the CBRAD5 samples displayed significant divergence (**Figure 2E**), whereas only ~5%  
201 internal nodes displayed divergence in the HESC sample. When such degree of divergence is  
202 depicted against normalized depths (see **Figure S3A** and **Methods**) of the corresponding nodes, the  
203 CBRAD5 samples consistently showed rapid divergence that is not seen in HESC samples (**Figure**  
204 **2E**). Please note that divergence here is not equivalent to differentiation, since two sister cells  
205 differentiating into the same fate would not reveal any divergence for their mother cell. In other  
206 words, divergence implies asymmetric division creating daughter cells of different developmental  
207 potentials, whereas differentiation can occur during symmetric division giving rise to a pair daughter  
208 cells that both activate a particular function or differentiate in the same direction.

209 By applying the above analysis to gene subsets associated with specific GO terms, it is possible  
210 to elucidate the progression of divergence in the corresponding cellular functions. As shown in  
211 several key GO terms including Wnt signaling (**Figure 2B**), the cumulative growth in the fraction  
212 of internal nodes with significant divergence at various normalized depths is also highly  
213 reproducible among CBRAD5 samples, and it differs from the hESC sample (**Figures 2F** and  
214 **Figure S3B**). Additionally, we examined whether our CLT data could resolve the temporal order of  
215 divergence completion for different cellular functions. To this end, we traced all root-to-tip paths on  
216 the CLTs and calculated the average depth of the last (furthest from the root) internal node exhibiting  
217 significant divergence on a GO term. As a result, the normalized depths of divergence completion  
218 appear consistent with known temporal orders of key developmental events (**Figure 2G**).  
219 Collectively, these results indicate that our dataset of single-cell transcriptomes and CLTs allowed  
220 the elucidation of cellular development with reasonable resolution.

## 221 **Transcriptional memory has limited contribution to developmental** 222 **canalization**

223 Following confirmation of the CLT data's resolution, we began searching for contributors to  
224 developmental robustness using CLTs. A first hypothesis is that transcriptional memory may have  
225 constrained gene expression variation during development, which would canalize transcriptomic

226 state during development and contribute to robustness. In this context, transcriptional memory is  
227 the phenomenon of cells closely related on the CLT displaying similar expression levels due to the  
228 inheritance of the same cellular contents (proteins/transcripts) and/or epigenetic states from recent  
229 common ancestors<sup>29,31,35,36</sup>. Nevertheless, gene expression can also be restricted by transcriptional  
230 regulation that has nothing to do with cellular inheritance, such as negative feedback<sup>37</sup> and denoising  
231 promoters<sup>38</sup>. If the transcriptional memory dominates the experimented differentiation, one would  
232 expect all cells of the same type would have been clustered into an exclusive sub-CLT, which is  
233 clearly not the case (**Figure 1F**). For a quantitative analysis, we reasoned that the CV of single-cell  
234 expression levels within real sub-CLTs should reflect the combined effect of transcriptional memory  
235 and inheritance-independent regulation (**Figure 3A top**), whereas that of CLTs randomized by  
236 shuffling cells of the same type at different lineage positions should reflect only inheritance-  
237 independent regulation but not transcriptional memory (**Figure 3A bottom**). It is therefore possible  
238 to isolate the contribution of transcriptional memory to the expression constraint by contrasting the  
239 CV of real CLTs with that of randomized CLTs (**Figure 3A and Methods**), which is hereinafter  
240 referred to as the "memory index". We note that this definition of memory index is similar to that  
241 used in previous transcriptional memory-related studies<sup>29,39</sup>.

242 For each cell type, we calculated an overall memory index for each gene in each sub-CLT  
243 (**Figure 3B and Figure S4**). The top (10%) memory indices (**Figure 3C**) were found to be enriched  
244 in pluripotent cell types (C1/C2/C3/C4) as compared to progenitor cell types (C6/C7/C9/C10) (*t*-  
245 test  $P=0.0039$ , **Figure 3D**), suggesting that transcriptional memory is more important to maintaining  
246 pluripotency than differentiation. Because transcriptional memory is mediated by cellular contents  
247 inherited from mother to daughter cells, such as transcription factors, we hypothesized that these  
248 genes with top memory indices should exhibit significant overlap with those regulated by some  
249 related transcription factors. Thus, we tested these genes for enrichment in genes responsive to  
250 genetic perturbation of individual transcription factors<sup>40</sup> (see **Methods**), and made two observations.  
251 First, some transcription factors with known involvement in the experimented differentiation, such  
252 as Nanog in the pluripotent C1<sup>41</sup> and Gata6 in progenitor C6/C9<sup>42</sup>, indeed exhibit significant  
253 enrichment of the genes with top memory index. Second, the enrichment was generally stronger for  
254 pluripotent cell types than it was for progenitor cell types (**Figure 3E**), a pattern again suggesting  
255 that transcriptional memory only played a minor role in differentiation, which is at least not as  
256 significant as in maintaining pluripotency.

## 257 **Stable cell type compositions across sub-clones supports robust** 258 **development**

259 Observations above indicate that terminal cells within a sub-CLT have restricted fates that are  
260 not dominated by transcriptional memory from the common ancestor (root of the sub-CLT). This

261 observation automatically prompted an assessment of the cell fate restriction imposed by  
262 inheritance-independent regulation, as well as its contribution to the robustness of developmental  
263 processes. We reasoned that inheritance-independent regulation should result in multiple similarly  
264 restricted sub-CLTs dispersed across the entire CLT. Therefore, we calculated the terminal cell type  
265 composition for each sub-CLT found in the CBRAD5 samples and compared it with the overall  
266 composition of the corresponding full CLT (see **Methods**). Intriguingly, the cell type compositions  
267 of sub-CLTs are usually more similar to those of the full CLTs than expected in randomized CLTs  
268 (**Figure 4A-C**). A closer examination of some sub-CLTs reveals a highly stable terminal cell type  
269 composition. For example, there are 35 sub-CLTs that generated subclones with highly stable (<10%  
270 deviation) proportions of 0.13, 0.39, 0.13 and 0.18 respectively for C6, C7, C9 and C10 (the top  
271 four most abundant progenitor cell types), which corresponds to the average proportion of these cell  
272 types in the three differentiating samples (**Figure 4D**). This observation suggests that a stereotyped  
273 developmental program may exist that produces subclones with highly similar compositions of cell  
274 types derived from multiple ancestral cells.

275 The observed polyclonal stereotypic development can be understood from two perspectives.  
276 On the one hand, the consistent execution of such a developmental program across subclones may  
277 be by itself a manifestation of robust genetic and/or molecular regulation. On the other hand, stable  
278 cell type compositions across subclones might enhance developmental robustness. We examined  
279 this latter perspective by simulating a CLT for the development of a single cell into an "organoid"  
280 consisting of 1,024 cells (i.e., 10 cell cycles) comprised of four types (namely  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ ) of cells  
281 in a 1:1:2:4 ratio. These cells formed 128 functional units each consisting of one  $\alpha$  cell, one  $\beta$  cell,  
282 two  $\gamma$  cells, and four  $\delta$  cells. Normally developed organoid consisting of 128 functional units  
283 (assuming sufficient cellular migration) are considered 100% functional. Meanwhile, CLT perturbed  
284 by random necrosis (see below), which results in the loss of some ancestral cells and all their  
285 descendants, has a functional capacity defined as the fractional survival rate of functional units with  
286 proper cellular composition. This design was inspired by the observation that functional units in  
287 living tissues, such as mouse pancreatic islets, display a highly stable cell type composition as the  
288 outcome of normal development<sup>26</sup>. To generate the normal (necrosis-free) CLT with the  
289 predetermined number of cells of each type, two models were used. The first "random" model  
290 assigns each cell to a random tip of the CLT regardless of its cell type (**Figure 4E** left). A second  
291 "stereotyped" model defines all eight-tip sub-CLTs as strictly consisting of one  $\alpha$  cell, one  $\beta$  cell,  
292 two  $\gamma$  cells, and four  $\delta$  cells, but different placements of these cells are allowed on the tips (**Figure**  
293 **4E** right). A total of 1,000 normal CLTs were generated under each model, and the functional  
294 capacity of each CLT was determined by exposing all (internal or terminal) cells to various rates of  
295 random necrosis. When compared to the random model, we found that CLTs generated with the  
296 stereotyped models always formed more functional units, or in other words, were more robust  
297 against necrosis (**Figure 4F**). Such enhanced developmental robustness is more evident at higher  
298 rate of necrosis (**Figure 4F**). Collectively, these results suggest that the observed stable cell type

299 composition among subclones contributed to developmental robustness.

## 300 **Stereotyped cell lineage trees underlie stable cell type compositions**

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

317 Just as sequence motifs are identified by comparisons between (sub-)sequences, tree motifs  
318 should also be identified through comparisons between (sub-)CLTs. In order to identify potential  
319 tree motifs in the CLT of the differentiating samples, we utilized Developmental cEll Lineage Tree  
320 Alignment (DELTA), an algorithm we previously developed for quantitative comparisons and  
321 alignments between CLTs<sup>43</sup> (**Figure 5A**, see **Methods** and **Text S1**). Using a dynamic programming  
322 scheme analogous to that employed by classical algorithms looking for similarities between  
323 biological sequences (e.g. the Smith-Waterman algorithm), the DELTA algorithm searches for pairs  
324 of homeomorphic sub-CLTs<sup>43</sup> within two given full CLTs. As a result, DELTA identified a large  
325 number of highly similar sub-CLT pairs between and within differentiating samples (**Figure 5B**).  
326 Some of the most frequently occurring sub-CLTs exhibited a consistent structure, comprising  
327 multiple layers of internal cells, a stable composition of terminal cell types, and appeared 20 to 40  
328 times in the three differentiating samples (**Figure 5C**). Groups of such highly similar sub-CLTs  
329 represent strong candidates of tree motifs on the developmental CLT, and strongly supports the  
330 existence of a stereotyped developmental program that contributes to developmental robustness.

## 331 **Discussion**

332 In the current study, we have reconstructed high density developmental CLTs for *in vitro*

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

344 There are a couple limitations of our study that are worth discussing here. First, our study was  
345 based on an *in vitro* directed differentiation model. This choice is a compromise between the  
346 feasibility for reconstruction of high density CLTs and a model that closely reflects the *in vivo*  
347 development. We believe our experiment reasonably recapitulates the *in vivo* situation because clear  
348 morphology of alveolar can be reach on the 20th day of the directed differentiation (**Figure S1A**  
349 and **Video S1**). Ideally, organoid or *in vivo* models should be combined with single-cell  
350 transcriptomes of a larger throughput (in terms of number of cells) in order to assess the question at  
351 a broader scale. Nevertheless, our main conclusion of polyclonal stereotyped development is most  
352 likely NOT an artefact of *in vitro* development, because none of the media components can create  
353 such pattern, and the number of ancestor hESCs seeding the colony is not correlated with the  
354 frequency of recurrence of lineage motifs. Second, we have not inferred detailed molecular  
355 processes and/or trajectories of gene expression changes in the stereotyped sub-CLT, as can be done  
356 for the nematode *Caenorhabditis elegans*<sup>43</sup>, whose temporal changes in gene expression have been  
357 recorded by microscopic image<sup>44,45</sup>. In the near future, this may be possible when the algorithms for  
358 inferring ancestral states based on cell lineage trees become sufficiently accurate<sup>19,46</sup>.

359 As a preliminary assessment on how the stereotyped CLT occurs, we treated the cell type  
360 composition of all descendent tips as a quantitative trait of the ancestral cells (internal nodes of the  
361 CLT) and regressed the difference of this trait between two ancestral nodes (that is not descendent  
362 of each other) onto their relatedness on the cell lineage (see Methods). This method, known in the  
363 genetics literature as a Haseman–Elston Regression<sup>47,48</sup>, is an unbiased estimator of heritability. In  
364 all of our samples, cell type compositions displayed heritability to some degree, with the heritability  
365 in the differentiating samples being significantly greater than that in the non-differentiating sample  
366 (**Figure S5**). Furthermore, similarly estimated heritability of single-cell transcriptome for each  
367 sample were lower than that of cell type composition (**Figure S5**). This result is unlikely to be  
368 explained by the higher measurement accuracy of cell type composition compared to single-cell  
369 transcriptomes for two reasons. First, the cell type itself is inferred based on single-cell  
370 transcriptomes. Second, the heritability of cell type composition in the non-differentiating sample

371 is almost equal to that of the single-cell transcriptome, suggesting similar measurement accuracy  
372 for these two traits. Thus, we concluded that descendent cell type composition is a heritable trait of  
373 ancestral cells. This trait is likely inherited from their earlier common ancestors by a mechanism  
374 independent of transcriptional memory, and is therefore expected to be pervasive in a CLT.

375 Beyond the specific mechanisms underlying developmental robustness, our findings suggest a  
376 novel perspective regarding cell types within the context of stereotyped sub-CLTs. In particular, just  
377 as letters can be better understood within the context of words, and nucleotides/amino acids can be  
378 better understood within the context of sequence motifs, stereotyped sub-CLTs can potentially  
379 bridge our knowledge of the atlas of cell types and their organization into functional tissues. Indeed,  
380 Elowitz and colleagues<sup>49</sup> recently identified statistically overrepresented patterns of cell fates on  
381 lineage trees as indicative of progenitor states or extrinsic interactions. The analysis was done using  
382 their newly proposed Lineage Motif Analysis, which differs from the method presented here that  
383 examined cell type composition and topological structure on incomplete CLTs, as their method uses  
384 the fully resolved CLTs and only analyzes cell type composition. Nevertheless, similar to our  
385 proposition here, they considered lineage motifs as a way of breaking complex developmental  
386 processes down into simpler components<sup>49</sup>.

## 387 **Methods**

### 388 **Design of the lineage tracer hESC cell line**

389 To design the lineage barcode and corresponding sgRNA, we first generated randomized 20-  
390 bp candidate sgRNA sequences with >3 substitutions relative to any human genome fragments.  
391 Among these candidates, the spacer sequence 5'-TATTCGCGACGGTTCGT-ACG-3' was selected  
392 as sgRNA1. A total of 13 protospacer sequences were designed based on sgRNA1 according to the  
393 following criteria: (i) each protospacer contained 2-3 mismatches with sgRNA1, (ii) there was no  
394 recurrence of any sequence of 9 bp or longer, and (iii) consecutive repeats of the same nucleotide  
395 for more than 2 bp were completely absent. The 13 protospacers (along with PAM, or protospacer  
396 adjacent motif) were organized according to decreasing CFD (cutting frequency determination)  
397 scores into the full lineage barcode<sup>50,51</sup>. The next three sgRNAs, sgRNA2, sgRNA3, and sgRNA4,  
398 were designed to perfectly match the 9th, 12th, and 13th protospacers, but with lower CFD scores  
399 (<0.55) for other protospacers, because these three protospacers were rarely edited in preliminary  
400 experiments using only sgRNA1. To facilitate capture by poly-dT reverse transcription primers on  
401 10x gel beads, the full lineage barcode with a 20-nt poly-dA(A20) 3' tail was inserted into the 3'UTR  
402 of an EGFP driven by an EF1 $\alpha$  promoter.

403 We constructed lineage tracer hESC cell lines by genomic integration of the lineage barcode,  
404 doxycycline-inducible Tet-on Cas9 and the sgRNAs. Briefly, the lineage barcode vector (pLV-

405 EF1A>EGFP:T2A:Bsd:V1(Barcode), VectorBuilder, no:VB1709 11-1008qmt) was constructed by  
406 the Gateway system and then transfected into H9 hESCs with MOI=0.15. The EGFP-fused lineage  
407 barcode was confirmed to exist as a single copy in the genome and to be highly expressed after  
408 blasticidin selection (15 µg/ml, InvivoGen, no. ant-bl-1) and flow cytometry sorting. Then the Tet-  
409 on inducible Cas9 vector (PB-Tet-ON-T8>Cas9:T2A:puro-PGK:rtTA, donated by Professor  
410 Jichang Wang, Zhongshan School of Medicine, Sun Yat-sen University) was co-transfected with  
411 hyPBBase (VectorBuilder, no: VB190515-1005npr) in a ratio of 1µg:100ng for 1x10<sup>7</sup>/ml cells by  
412 Neon<sup>TM</sup> transfection system (Life, MPK5000). In order to ensure adequate Cas9 expression for  
413 efficient editing, we applied double reinforced selection of Puromycin (1.0 µg/ml, InvivoGen, no.  
414 ant-pr-1) and Doxycycline (Dox, 1.0µg/ml, sigma, D9891-1G) for 7 days. Lastly, the sgRNA vector  
415 (pLV-U6>sgRNA1>U6>sgRNA2>U6>sgRNA3>U6>sgRNA4-EF1α>Mcherry:T2A:Neo,VB1912  
416 11-3149jwe) was constructed by Golden Gate ligation and transfected at MOI=30. H9 hESC cells  
417 with high expression of sgRNAs (fused with mCherry) were enriched by G418 selection (1000  
418 µg/ml, InvivoGen, ant-gn-1) for 11 days and flow cytometry sorting. Expression levels of Cas9,  
419 lineage barcode and sgRNA1 transcripts were detected by RT-qPCR with primers listed in **Table**  
420 **S8**.

421 The editing efficiency of the lineage tracer hESC cell line was evaluated by inducing Cas9  
422 expression in mTesR media with 1.0 µg/ml Dox for five days. We extracted gDNA from all cells  
423 using DNeasy Blood & Tissue Kits (Qiagen, no.69504). Using primers gDNA-V1-F and gDNA-  
424 V1-R (**Table S8**), we amplified the lineage barcode from gDNA using Phanta Max Super-Fidelity  
425 DNA Polymerase (Vazyme, No. P505), which was then cloned into pCE-Zero vector (Vazyme, No.  
426 C115). The efficiency of editing was then evaluated by colony PCR and Sanger sequencing for 50  
427 recombinant clones.

428 Additionally, we examined editing efficiency in the context of our directed differentiation  
429 experiment, in which only a small number of initial cells were used to form each colony. In 96-well  
430 dishes, matrigel (Corning, No. 354277) was plated and each well was seeded with < 10 log-phased  
431 lineage tracer hESC cells manually by micromanipulation. For 11 days, the cells were cultured in  
432 100 µl of mTesR media, to which 10 µl of cloneR (Stemcell, No.05888) were added on day0 and  
433 day2, and 1.0 µg/ml Dox+ mTesR media was added and refreshed every 48 hours since day2.  
434 Normally surviving colonies after the 11-day culture were harvested by GCDR (Stemcell,  
435 No.07174). Next, 50ng of genomic DNA was extracted from each colony using the QIAamp DNA  
436 Micro Kit (Qiagen, No.56304) and PCR amplified for the lineage barcode. The Cas9-induced  
437 mutations accumulated during colony formation were then identified by Sanger sequencing, TA  
438 cloning-based sequencing or Illumina HiSeq PE250 sequencing. Specifically, the raw HiSeq data  
439 were trimmed by fqtrim (<https://ccb.jhu.edu/software/fqtrim/>) with default parameters. The paired  
440 reads were merged by FLASH<sup>52</sup> using 30 bp of overlapping sequence and 2% mismatches.  
441 Sequences alignable to the human reference genome by Bowtie2 with default parameters<sup>53</sup>, or to  
442 primer sequences of gDNA-V1-F and gDNA-V1-R with two mismatches, were removed as they

443 likely represented nonspecifically amplified sequences. MUSCLE<sup>54</sup> aligned the sequenced lineage  
444 barcode to the wild-type lineage barcode using default parameters. The editing events of each  
445 sequence were identified according to a previous method<sup>50</sup>.

## 446 **Validating directed differentiation from hESC to lung progenitor and** 447 **alveolosphere**

448 Using the BU3 NGST (NKX2-1-GFP; SFTPCtdTomato) iPS cell line (donated by Professor  
449 Darrell N. Kotton, Department of Medicine, Boston University), we tested the protocol of directed  
450 differentiation towards lung progenitor and alveolosphere published by Kotton and colleagues<sup>30</sup>.  
451 Briefly, in six-well dishes pre-coated with Matrigel (Stemcell, No.356230),  $2 \times 10^6$  cells maintained  
452 in mTESR1 media were differentiated into definitive endoderm using the STEMdiff Definitive  
453 Endoderm Kit (StemCell, No.05110), adding supplements A and B on day 0, and supplements B  
454 only on day 1 to day 3. Flow cytometry was used to evaluate the efficiency of differentiation to  
455 definitive endoderm at day 3 using the endoderm markers CXCR4 and c-KIT (Anti-human CXCR4  
456 PE conjugate, Thermo Fisher, MHCXCR404,1:20; Anti-human c-kit APC conjugate, Thermo Fisher,  
457 CD11705, 1:20; PE Mouse IgG2a isotype, Thermo Fisher, MG2A04,1:20; APC Mouse IgG1 isotype,  
458 Thermo Fisher, MG105, 1:20) based on the method of Sahabian and Olmer<sup>55</sup>. After the endoderm-  
459 induction stage, cells were dissociated for 1-2 minutes at room temperature with GCDR and  
460 passaged at a ratio between 1:3 to 1:6 into 6 well plates pre-coated with growth factor reduced  
461 matrigel (Stemcell, No.356230) in “DS/SB” anteriorization media, which consists of complete  
462 serum-free differentiation medium (cSFDM) base, including IMDM (Thermo Fisher, No.12440053)  
463 and Ham’s F12 (Corning, No. 10-080-CV) with B27 Supplement with retinoic acid (Gibco,  
464 No.17504044), N2 Supplement (Gibco, No.17502048), 0.1% bovine serum albumin Fraction V  
465 (Sigma, A1933-5G), monothioglycerol (Sigma, No. M6145), Glutamax (ThermoFisher, No. 35050-  
466 061), ascorbic acid (Sigma,A4544), and primocin with supplements of 10  $\mu$ m SB431542 (“SB”;  
467 Tocris, No.1614) and 2  $\mu$ m Dorsomorphin (“DS”; Sigma, No. P5499). In the first 24 hours  
468 following passage, 10  $\mu$ m Y-27632 was added to the media. After anteriorization in DS/SB media  
469 for three days (72 hrs, from day 3 to day 6, refreshed every 48 hours), cells were cultured in “CBRa”  
470 lung progenitor-induction media for nine days (from day 6 to day 15, refreshed every 48 hours).  
471 This CBRa media consists of cSFDM containing 3  $\mu$ m CHIR99021 (Tocris, No.4423), 10 ng/mL  
472 rhBMP4 (R&D, 314-BP-010), and 100 nM retinoic acid (RA, Sigma, No. R2625). At day 15 of  
473 differentiation, single-cell suspensions were prepared by incubating the cells at 37°C in 0.05%  
474 trypsin-EDTA (Gibco, 25200056) for 7-15 minutes. The cells were then washed in media containing  
475 10% fetal bovine serum (FBS, ThermoFisher), centrifuged at 300 g for 5 minutes, and resuspended  
476 in sort buffer containing Hank’s Balanced Salt Solution (ThermoFisher), 2% FBS, and 10  $\mu$ m Y-  
477 27632. The efficiency of differentiation into NKX2-1<sup>+</sup> lung progenitors was evaluated either by  
478 flow cytometry for NKX2-1-GFP reporter expression, or expression of surrogate cell surface

479 markers CD47<sup>hi</sup>/CD26<sup>lo</sup>. Cells were subsequently stained with CD47-PerCPCy5.5 and CD26-PE  
480 antibodies (Anti-human CD47 PerCP/Cy5.5 conjugate, Biolegend, Cat#323110, 1:200; Anti-human  
481 CD26 PE conjugate, Biolegend, Cat#302705, 1:200; PE mouse IgG1 isotype, Biolegend,  
482 Cat#400113, 1:200, PerCP/Cy5-5 mouse IgG1 isotype, Biolegend, Cat#400149, 1:200) for 30 min  
483 at 4 °C, washed with PBS, and resuspended in sort buffer based on the method of Hawkins and  
484 Kotton<sup>55</sup>. Cells were filtered through a 40 µm strainer (Falcon) prior to sorting. The CD47<sup>hi</sup>/CD26<sup>lo</sup>  
485 cell population was sorted on a high-speed cell sorter (MoFlo Astrios EQs) and resuspended in  
486 undiluted growth factor-reduced 3D matrigel (Corning 356230) at a dilution of 20-50 cells/µl, with  
487 droplets ranging in size from 20 µl (in 96 well plate) to 1 ml (in 10 cm dish). Cells in 3D matrigel  
488 suspension were incubated at 37°C for 20-30 min, followed by the addition of warm media. The  
489 differentiation into distal/alveolar cells after day 15 was performed in “CK+DCI” medium,  
490 consisting of cSFDM base, with 3 µM CHIR (Tocris, No.4423), 10 ng/mL rhKGF(R&D, No.251-  
491 KG-010) (CK), and 50 nM dexamethasone(Sigma, No. D4902), 0.1 mM 8-Bromoadenosine 3',5'-  
492 cyclic monophosphate sodium salt (Sigma, No.B7880) and 0.1 mM 3-Isobutyl-1-methylxanthine  
493 (IBMX; Sigma, No.I5879) (DCI). Immediately after replating cells on day 15, 10 µM Y-27632 was  
494 added to the medium for 24 hours. Upon replating on day 15, alveolospheres developed in 3D  
495 Matrigel culture outgrowth within 3-7 days, and were maintained in CK+DCI media for weeks.  
496 These spheres were analyzed by Z stack live images of alveolospheres taken and processed on the  
497 Leica DMi8 fluorescence microscope.

## 498 **Directed differentiation followed by simultaneous assessment of** 499 **single-cell transcriptomes and cell lineage tree.**

500 Based on the results from the full directed differentiation experiment above, we aimed to  
501 evaluate single-cell transcriptomes and CLTs simultaneously for directed differentiation from  
502 hESCs to PLP, a stage at which the colony has <10,000 cells, allowing us to sample a large  
503 proportion of cells. To prepare suitable ancestor hESCs, the cell colonies outgrowth after 5-7 days,  
504 plated in 96-well dishes with microscopic selection for GFP<sup>+</sup> mCherry<sup>+</sup>, were digested with GCDR  
505 to form ~50 µm aggregates, and cultured in mTesR media until day 5. Combining selection and  
506 induction by dox (1.0 µg/ml) and puro (1.0 µg/ml) from day 5 to day 7, the normally survived GFP<sup>+</sup>  
507 mCherry<sup>+</sup> colonies were capable of Cas9 expression and marked by primary editing events (to  
508 distinguish ancestor cells), as confirmed by DNA extraction and barcode PCR and sanger  
509 sequencing. The cell colonies with primary editing events were digested by GCDR for cell counting  
510 (~ 4000 cells) and resuspended at a density of 10 cells/µl. 1µl cell suspension was added into each  
511 well of 96-well dishes plated with 1:10 diluted Matrigel (Corning, No.354277) for culture in mTesR  
512 media with ClonR (10:1) (Stemcell, No.05888) added in the first 48h to promote the survival of  
513 very few stem cell. Directed differentiation was then initiated by applying both dox (1.0µg/ml, for  
514 editing the lineage barcode) and the STEMdiff Definitive Endoderm Kit to the normally survived

515 colonies. Later stages of directed differentiation followed the differentiation protocols described  
516 above, with the exception that it was stopped on the tenth day after its initiation (**Figure S1B**).  
517 Finally, colonies with intermediate size (~ 5,000 cells as approximated by colony size and cell  
518 counts) and  $\geq 50\%$  GFP<sup>+</sup> Mcherry<sup>+</sup> cells were digested with 0.05% trypsin-EDTA for 1 minute at  
519 37 °C, washed in PBS containing 10% fetal bovine serum (FBS, ThermoFisher), centrifuged at 500  
520 g for five minutes, and resuspended in single cell resuspension buffer containing PBS and 0.04%  
521 BSA. Using the standard 10x Chromium protocol, cDNA libraries were prepared from these single  
522 cell suspensions. Each cDNA library was split into two halves, with the first half subjected to  
523 conventional RNA-seq for single-cell transcriptomes, and the other half subjected to amplification  
524 of the lineage barcode followed by PacBio Sequel-based HiFi sequencing of the lineage barcode  
525 (**Figure 1A**).

## 526 **Analysis of scRNA-seq**

527 Following the 10x Genomics official guidelines, we used the Cell Ranger<sup>56</sup> pipeline to map  
528 raw reads to the human reference genome (GRCh38) by STAR<sup>57</sup> and obtained the read counts for  
529 each gene. Using Seurat v3.2.1<sup>58</sup>, we retained cells with <10% mitochondrial reads and >200  
530 expressing unique features detected. Then highly variable genes were detected by Single-cell  
531 Orientation Tracing (SOT)<sup>59</sup>, which were then subjected to Principle Component Analysis, followed  
532 by batch effect correction by Harmony<sup>60</sup>. We then clustered cells based on the cell-cell distance  
533 calculated by FindNeighbors and FindClusters using the Harmony-normalized matrix of gene  
534 expression. Then, we used runUMAP for visualization and FindAllMarkers to obtain differentially  
535 expressed genes (DEGs) among clusters. To identify cell types, we downloaded microarray data  
536 from Gene Expression Omnibus (GEO)<sup>33,61</sup>, and extracted DEGs (Wilcoxon Rank Sum test,  $P <$   
537 0.01) in different stages of differentiation towards PLP. We scored the clusters base on the average  
538 expression and numbers of expressed stage-specific DEGs. Finally, we named 12 cell cluster based  
539 on the inferred order of appearance in the differentiation progress.

## 540 **Construction of cell lineage trees**

541 Based on the PacBio HiFi sequencing results, we built and assessed the quality of the CLT  
542 from PacBio HiFi reads following our previous pipeline<sup>31</sup>. Briefly, using HiFi-seq raw sequences,  
543 we called consensus sequences separately from positive and negative strand subreads from each  
544 zero-mode waveguide (ZMW). We reserve only consensus sequences with at least three subreads  
545 and identifiable barcode primers (**Table S8**, allowing up to two mismatches). From the consensus  
546 sequences, 10x cell barcodes and UMIs were extracted and matched to those from scRNA-seq, with  
547 one mismatch allowed. Lineage barcode sequences were then extracted from the consensus

548 sequences, grouped by identical cell barcode and UMI, then merged by MUSCLE alignment  
549 followed by selecting the nucleotide with the highest frequency at each site. After MUSCLE  
550 alignment of the merged sequence to the reference lineage barcode, the editing events were called<sup>50</sup>.  
551 Then, for each lineage barcode allele from the same cell, the frequency was calculated as the total  
552 number of UMIs of the allele and its ancestral allele. Here, the ancestral allele of a specific allele  
553 was defined as any allele in which the observed editing events were a subset of the editing events  
554 in the focal allele. Finally, the lineage barcode allele of a cell was defined as the allele with the  
555 highest frequency, prioritizing the alleles with more editing events if the frequencies were equal.  
556 For each sample, all cells with a lineage barcode and a single-cell transcriptome were used to  
557 construct a multifurcating lineage tree based on the lineage barcode using the maximum likelihood  
558 (ML) method implemented by the IQ-TREE LG model<sup>62</sup>.

## 559 **Transcriptome divergence among cell type clusters**

560 To elucidate the transcriptomic divergence among the observed clusters in the context of the  
561 directed differentiation towards PLP, we extracted stage-specific DEGs with the top 10% most  
562 extreme fold-change relative to other stages (**Figure 2A**, using microarray data<sup>33</sup> mentioned above),  
563 and identified the Gene Ontology terms enriched (BH-adjusted  $P < 0.05$ , Fisher's exact test) with  
564 these stage-specific DEGs. After eliminating GO terms that have very few expressed genes, we  
565 focused on 179 GO terms (**Table S8**). For each cell, the activities of the specific cellular functions  
566 represented by these GO terms were estimated by the AddModuleScore function of Seurat, which  
567 basically calculated the average Z-score transformed expression levels of all genes annotated by the  
568 GO term. All cells within a cluster were then combined to determine the average activity of the GO  
569 term for the cluster (**Figure 2B**).

## 570 **Transcriptome divergence among sub-CLTs**

571 As for the divergence among sub-CLTs, estimation of pseudotime was conducted via  
572 Monocle<sup>34</sup> with all cells on differentiating CLTs pooled together. After Principal Component  
573 Analysis of all cells from all samples combined, the transcriptomic divergence ( $D_T$ ) between any  
574 two cells is quantified by one minus Pearson's Correlation Coefficient of the top 100 principal  
575 components. The developmental potential of an ancestor cell (an internal node on the CLT) was then  
576 calculated by the summed squared  $D_T$  of all pairs of its descendant cells. The reduction of  
577 developmental potential ( $\Delta_{DP}$ ) during the growth of an internal node to its daughter nodes was  
578 calculated by the focal internal node's  $\Delta_{DP}$  subtracted by the summed  $\Delta_{DP}$  of all its daughter nodes  
579 (**Figure 2D**). The statistical significance of an observed  $\Delta_{DP}$  was estimated by contrasting the  
580 observation with its null distribution generated by random assignment of single-cell transcriptomes  
581 from hESC samples to the focal CLT (**Figure 2D**). We emphasized here that the null distribution

582 should be estimated by the single-cell transcriptomes from the non-differentiating hESC sample,  
583 since using those from the differentiating CBRAD5 samples would introduce actual divergence into  
584 the null and thus lead to an underestimated statistical significance. It is also worth noting that this  
585 method is very similar to the commonly used nonparametric method of permutational multivariate  
586 analysis of variance (PERMANOVA<sup>63</sup>), except that Pearson's correlation-based divergence replaces  
587 the distance-based divergence used in canonical PERMANOVA, as the correlation-based metric has  
588 consistently been shown to result in superior performance for single-cell transcriptomes<sup>64,65</sup>. We  
589 have also applied this PERMANOVA-based method to subsets of genes within the transcriptome.  
590 For example, only genes annotated with a specific GO term (**Table S8**) were used. A significant  
591 divergence for a specific GO term does not necessarily indicate a significant divergence in the whole  
592 transcriptome, since genes annotated with the GO term may have a small effect on the transcriptome  
593 as a whole. As a result, internal nodes with transcriptomic divergence do not necessarily represent  
594 a larger fraction than nodes with divergence on a specific GO term.

595 In order to perform a retrospective analysis of divergence progression, we need a normalized  
596 temporal scale that is comparable across samples. In theory, this scale could be derived from the  
597 mutation rate of the lineage barcode and/or the topological depth of a node (i.e., the number of nodes  
598 between the root and the focal node). Considering the variability in Cas9 editing efficiency over  
599 barcodes, as well as long inter-site deletions, we discarded the mutation rate-based scale. For the  
600 topological depth scale, due to both biological and experimental stochasticity, the reconstructed  
601 CLTs and their nodes have very different depths, despite the fact that they are supposed to  
602 correspond to the ten-day directed differentiation. Assuming that the internal nodes were evenly  
603 sampled on all root-to-tip paths throughout the CLT, the actual depth of a node should be reflected  
604 equally by its depth from the root and (indirectly) by the depth from the focal node to its descendent  
605 tips. Based on this logic, we defined the normalized depth of a node as  $d = (d_r/d_t +$   
606  $(1 - d_s/d_t))/2$ , where  $d_r$  is the focal node's depth from root,  $d_t$  is the max depth found in the CLT,  
607 and the  $d_s$  is the max depth from the focal node to its descendent tips (**Figure S3A**). Here, via  
608 division by  $d_t$ , all depths were scaled from 0 to 1, with 0 being the root and 1 being the tips with  
609 maximal raw depth within the CLT.

## 610 **Transcriptional memory index**

611 We followed previously proposed methods<sup>29,39</sup> to calculate transcriptional memory index. In  
612 each cell type and for each gene expressed in >10% of cells of this type, the CV of the expression  
613 levels was calculated among all terminal cells of this type within a sub-CLT (containing at least two  
614 cells of this type). The minimal CV among all sub-CLTs, i.e.  $\min(\text{CV})$ , was then used to represent  
615 the expression variability of the focal gene in this cell type. It was also calculated for each of 1,000  
616 randomized CLTs created by reassigning all cells of the same type to a new lineage position that

617 was originally occupied by the same cell type. These 1,000  $\min(CV_{Random})$  from randomized CLTs  
618 were averaged, i.e. mean ( $\min(CV_{Random})$ ), to yield a null expectation for the observed  $\min(CV)$ .  
619 Finally, the memory index was defined as  $M = (\min(CV) - \text{mean}(\min(CV_{Random}))) / \min(CV)$ . Note  
620 that the final division by  $\min(CV)$  is different from the previously defined memory index<sup>29,39</sup>, but  
621 allows comparisons between genes with very different baseline CVs or expression levels.

622 To test the hypothesized role of transcription factors in mediating transcriptional memory, we  
623 obtained lists of gene sets responsive to perturbations of individual transcription factors  
624 (“TF\_Perturbations\_Followed\_By\_Expression” in Enrichr<sup>40</sup>). The genes with highest memory  
625 indices (top 10% across all cell types) were assessed for enrichment in each of these TF-responsive  
626 gene sets using Enrichr<sup>40</sup>. We reported (**Figure 3E**) the “combined score” calculated by Enrichr,  
627 which takes into account both the statistical significance and the magnitude of enrichment  
628 (combined score of enrichment  $c = \log(p) * o$ , where  $p$  is the  $P$  value from Fisher’s exact test and  $o$   
629 is the odds ratio of the enrichment<sup>40</sup>).

## 630 **Composition of terminal cell types compared among sub-CLTs and** 631 **the full CLTs**

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

## 645 **Robustness of random versus stereotyped development**

646 Without loss of generality, we defined a functional unit as consisting of four cell types, namely  
647  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$ , in a 1:1:2:4 ratio. We simulated 1000 binary CLTs, each consisting of 1024 terminal  
648 cells (128  $\alpha$  cells, 128  $\beta$  cells, 256  $\gamma$  cells, 512  $\delta$  cells) generated through ten cell cycles, under two  
649 developmental models. The first “random” model randomly assigns the four types of cells onto the

650 tips of the tree. A second “stereotyped” model strictly assigns  $\alpha$ ,  $\beta$ ,  $\gamma$ , and  $\delta$  cells in a 1:1:2:4 ratio  
 651 onto each sub-CLT consisting of eight tips (three cell cycles). A predefined fraction (0.001, 0.005,  
 652 0.01, 0.05 or 0.1, as on  $x$  axis of **Figure 4F**) of the 2047 (1024 terminal and 1023 internal) cells  
 653 were chosen and removed along with all their descendent cells to mimic random necrosis. Assuming  
 654 sufficient cell migration to allow formation of the functional unit as long as there are enough  
 655 terminal cells of the proper type, the robustness is thus quantified by the number of functional units  
 656 that can be formed by all terminal cells surviving necrosis. A simple example shown in **Figure 4E**.

## 657 **Comparison and alignment of sub-CLTs by mDELTA**

658 Let us denote vectors/nodes as  $V$  and edges connecting nodes as  $E$ . Given a query tree  $Q =$   
 659  $(V, E)$  and a subject tree  $S = (V', E')$ , an isomorphic alignment is a bijection  $A : V \leftrightarrow V'$ , such  
 660 that for every pair of nodes with  $v, u \in V$ , we have  $(v, u) \in E \Leftrightarrow (A(v), A(u)) \in E'$ . Based on  
 661 two types of biologically informed tree editing operations, namely pruning and merging (see  
 662 **Supplementary Text**), a homeomorphic subtree alignment  $A$  between  $Q$  and  $S$  is defined as an  
 663 isomorphic alignment between  $Q'$  and  $S'$ , where  $Q'$  is the result of zero or more pruning and  
 664 merging in  $Q$ , and  $S'$  is the result of zero or more pruning and merging in  $S$ . Here all the pruning in  
 665  $Q$  and  $S$  are collectively denoted as  $\pi(A)$ , and all merging in  $Q$  and  $S$  are collectively denoted as  
 666  $\mu(A)$ . If we further denote the alignment score between two nodes  $v \in V$  and  $v' \in V'$  as  
 667  $a(v, v')$ , the cost for pruning a subtree  $\hat{T}$  as  $p(\hat{T})$ , the cost for merging an internal node  $\hat{v}$  with  
 668 its mother node as  $m(\hat{v})$ . The score of a homeomorphic subtree alignment  $A$  between  $Q$  and  $S$  can  
 669 then be expressed as

$$670 \quad w(Q, S, A) = \sum_{(v, v') \in A} a(v, v') - \sum_{\hat{T} \in \pi(A)} p(\hat{T}) - \sum_{\hat{v} \in \mu(A)} m(\hat{v})$$

671 Our algorithm of mDELTA find the optimal  $A$  (with optimal/highest possible  $w$ ) given  $Q, S, a,$   
 672  $p$  and  $m$  by a dynamic programming procedure. We defined  $a$  based on similarity of single-cell  
 673 transcriptomes,  $p$  based on the number of pruned terminal cells, and  $m$  based on the number of  
 674 merged internal nodes. Detail computational procedures of mDELTA can be found in  
 675 **Supplementary Text**.

## 676 **Heritability of quantitative traits in the CLT**

677 In order to gauge the heritability of quantitative traits on the CLT, we calculated the correlation  
 678 between the relatedness and the phenotypic divergence of a pair of nodes. When the relatedness is  
 679 defined by genomic relatedness like DNA sequence identity, this analysis is the same as the classic  
 680 statistical genetics method called Haseman-Elston Regression<sup>47</sup>. Thus, we consider the correlation  
 681 coefficient from this analysis to be a proxy for phenotypic heritability among nodes on the CLT.

682 However, we would like to emphasize that since the DNA sequences of all cells in our dataset are  
683 presumably nearly identical, the relatedness between nodes is therefore defined by their distance on  
684 the CLT instead (see below), and resulting correlation coefficients cannot be interpreted as  
685 traditional heritability as they are in Haseman-Elston Regression. Specifically, we define the  
686 relatedness between any two nodes on the CLT inversely by the number of cell divisions separating  
687 them, which is then estimated by contrasting the number of their descendent cells with the number  
688 of descendent cells of their latest common ancestor. Following previous Haseman-Elston  
689 Regression applications<sup>48</sup>, the relatedness between nodes was then scaled so that the mean  
690 relatedness between any pair of nodes is 0 and the maximal relatedness is 1. As such scaling is  
691 equivalent to calculating relatedness relative to a different population<sup>48</sup>, comparing the heritability  
692 of one trait relative to that of another trait would not be affected as long as both traits are analyzed  
693 in the same focal population (the focal CLT). On the phenotype side, we examined two quantitative  
694 traits, the single-cell transcriptomes of terminal nodes and the descendent cell type compositions of  
695 internal nodes. Here, the single-cell transcriptomes of terminal nodes were first processed by  
696 Principle Component Analyses, then all principle components of a cell is used to represent its  
697 transcriptome. As for the descendant cell type compositions of an internal node, each internal node  
698 is represented by a vector comprising  $M$  elements, where  $M$  is the total number of cell types  
699 identified in our dataset, and each element represents the percentage of descendent cells of that type.  
700 The phenotypic divergence between two nodes is calculated as the Euclidian distance between the  
701 multidimensional quantitative traits. Lastly, we reported the Spearman's Correlation Coefficient  
702 between the relatedness and the phenotypic divergence between all relevant node pairs in **Figure**  
703 **S5** as a proxy for the heritability of quantitative traits.

704

## 705 **Data availability**

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

708

## 709 **Code availability**

710 Custom R/Python codes that were used in data analysis, are available on GitHub  
711 (<https://github.com/ZhangxyOk/Stereotyped-CLT>). The mDELTA algorithm is deposited on a  
712 separated GitHub repository ([https://github.com/Chenjy0212/mdelta\\_full](https://github.com/Chenjy0212/mdelta_full)).

713

714 **Author contributions**

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

720

721 **Competing interests**

722 The authors declare no conflicts of interest.

723

724 **Funding**

725 This work was supported by the National Key R&D Program of China (grant number  
726 2021YFA1302500 and 2021YFF1200904 to J.-R. Y.), the National Natural Science Foundation of  
727 China (grant numbers 32122022 and 32361133555 to J.-R. Y.), and the Research Grants Council of  
728 the Hong Kong Special Administrative Region, China (grant number 21206223 to X.W.).

729

730 **Figure legends**

731 **Figure 1. Cell lineage tracing for directed differentiation of**  
732 **primordial lung progenitors**

733 (A) Schematic diagram illustrating the overall experimental process. The 10-day directed  
734 differentiation from several Lineage Tracer hESCs to primordial lung progenitors (PLP) was  
735 conducted along out with simultaneous lineage tracing utilizing inducible CRISPR-Cas9 editing of  
736 an expressed lineage barcode (13 editable sites). The resulting colony was assayed for single-cell  
737 transcriptomes by Nova-seq and lineage barcode by PacBio HiFi-seq, which were used to  
738 reconstruct CLTs with single-cell transcriptomes assigned to tips. (B) The variation among single-  
739 cell transcriptomes captured in the four samples (one non-differentiating “HESC” sample and three  
740 differentiating samples) as shown by UMAP. A data point represents a cell, which is colored based  
741 on its source sample on the left panel and the expression level of NKX2-1 (the marker for PLP) on  
742 the right panel. (C) Major clusters of the single-cell transcriptomes are differentially colored and  
743 labeled by their corresponding cell types. (D) In the 12 major cell types ( $y$  axis), differentially  
744 expressed genes (DEGs) found in bulk samples of specific developmental stages preceding PLP ( $x$   
745 axis) were examined for their average expression levels (dot color) and fraction of cells that  
746 expressed the gene (dot size). See also **Figure S2C**. (E) For each of the four samples ( $x$  axis), the  
747 percentage of cells belonging to each type was shown. The cell types are colored identically to those  
748 in panel C. (F) Reconstructed CLTs are visualized as circle packing charts for the four samples.  
749 Circles represent sub-CLTs, whose sizes indicate the number of terminal cells in the sub-CLTs, while  
750 the color (same as panel C) indicates the fraction of terminal cells belonging to each cell type. See  
751 **Figure S2E** for their tree representation. (G) A pair of cells' normalized lineage distance (the  
752 number of internal nodes on the path from one cell to the other, divided by the maximal lineage  
753 distance found in the sample) is highly correlated with the normalized allelic distance of their  
754 lineage barcodes (the total number of target sites that differed from the reference, divided by the  
755 maximum value of 26). All cell pairs were separated into five groups based on their normalized  
756 lineage distance ( $x$  axis), and the distribution of normalized allelic distances ( $y$  axis) within each  
757 group is shown in the form of a standard boxplot, with the mean value indicated by the white point.  
758 On top, Spearman's  $\rho$  and  $P$  value for raw data are indicated. (H) The probability of finding a  
759 common ancestral allele (as yet-to-decay transcripts) between a pair of single-cell tips decreased as  
760 their normalized lineage distance ( $x$  axis) increased. The error bars indicate the standard error  
761 estimated by bootstrapping the cell pairs for 1,000 times.

762 **Figure 2. The transcriptome divergence among cell type clusters and**

763 **among subclones**

764 (A) Heatmap for expression levels of DEGs extracted from microarray-based transcriptomes of  
765 specific developmental stages (color bars on top) of the directed differentiation<sup>33</sup>. (B) Functional  
766 activities of GO terms (*x* axis. Full list in **Table S7**) enriched with stage-specific DEGs were  
767 shown for every cluster (*y* axis) identified in our samples. Here functional activity as indicated by  
768 the color scale was estimated by the average Z-score-transformed expression of all genes  
769 annotated with the GO term. Some important GO terms are boxed and labeled by dashed lines,  
770 and are further analyzed in panel F and **Figure S3B**. (C) A coefficient of variation (CV) was  
771 calculated using pseudotime estimates of single-cell transcriptomes within a sub-CLT. These CVs  
772 were plotted for all real sub-CLTs (*y* axis) and corresponding randomized sub-CLTs generated by  
773 shuffling all tips (*x* axis) in each differentiating sample (name on top). As the dashed line indicates  
774  $x = y$ , sub-CLTs with CVs lower than random expectation (i.e. restricted variation) will appear  
775 below it. Each panel includes the number of CLTs above and below the dashed line, which was  
776 also tested against the binomial expectation (50% below the line) and yielded the *P* values on top.  
777 (D) Schematic diagram for the PERMANOVA-based estimation of transcriptome divergence for  
778 an internal node (see **Methods**). (E) Cumulative fraction (*y* axis) of internal nodes exhibiting  
779 significant transcriptome divergence as the normalized depth (*x* axis) considered increased.  
780 Results from different samples were shown with different colors, as indicated by the color legend.  
781 (F) Same as panel E except that the analyses were limited to specific GO terms indicated on top of  
782 each panel. (G) We calculated the normalized depths (*y* axis) at which the divergence of specific  
783 functions is completed. GO terms enriched of marker genes in representative developmental  
784 stages (*x* axis and colors) were examined. Dots represent GO terms and triangles represent the  
785 average depth within the same-color group. Significant *P* values from between-group Wilcoxon  
786 Rank Sum test are labeled on top.

787 **Figure 3. Limited contribution of transcriptional memory in**  
788 **differentiation**

789 (A) Schematic diagram for the CLT-based estimation of transcriptional memory. (B) Expression  
790 variability in the real CLT (*y* axis) compared to that in the randomized CLT (*y* axis). Each dot  
791 represents a gene in a cell type. Dot color shows the fraction of cells within the cell type that express  
792 the gene, as indicated by the color scale on top. (C) A stacked histogram showing the distribution  
793 of the memory indices calculated. A filled bar represents those estimated from pluripotent cell types  
794 and an empty bar represents those estimated from progenitor cell types. Genes exhibiting strong  
795 transcriptional memory, i.e. those with a memory index ranking among the top 10% (dashed line),  
796 were red, while others were gray. The inset shows a zoomed-in view of the large memory index

797 region. **(D)** Among different cell types, the fraction (height of bar) of genes exhibiting high memory  
798 indices was compared. The bars are colored similarly to those in panel **C**. **(E)** Gene sets responsive  
799 to perturbation of individual transcription factors ( $x$  axis) were tested for the enrichment of genes  
800 exhibiting strong signal of transcriptional memory (see **Methods**). The top ten transcription factors  
801 with the highest combined enrichment score ( $y$  axis) were shown for each cell type. The statistical  
802 significance of enrichment according to Fisher's exact test is indicated as \*: $P<0.05$ ; \*\*: $P<0.01$ ;  
803 \*\*\*: $P<0.001$ .

#### 804 **Figure 4. Stable cell type composition across sub-clones supports** 805 **robust development**

806 **(A-C)** In each panel for each of the CBRAD5 samples (names on top of the panel), the diversity of  
807 compositions of terminal cell types within sub-CLTs were estimated by a summed chi-square value  
808 ( $\chi^2$ ) (see **Methods**) as indicated by the red arrows. The same summed  $\chi^2$  values were calculated  
809 for 1,000 randomized CLTs, whose distribution was shown as a blue histogram. The probability of  
810 a summed  $\chi^2$  value being smaller than the observation (red arrow) is indicated by the  $P$  values in  
811 the panel. **(D)** For 35 sub-CLTs in CBRAD5 samples, the normalized depths of their roots ( $y$  axis)  
812 and the sizes of the sub-CLTs ( $x$  axis) were plotted. These sub-CLTs display highly similar terminal  
813 cell type compositions (less than 10% deviation from 0.13, 0.39, 0.13 and 0.18 respectively for C6,  
814 C7, C9 and C10) **(E)** A schematic diagram showing a simple model of the functional robustness of  
815 the random (left) versus stereotyped (right) development against random necrosis (indicated by "X").  
816 The robustness is quantified by the number of functional units (with cell type compositions indicated  
817 in the triangle) that can be formed by terminal cells surviving necrosis, as exemplified at the bottom.  
818 **(F)** Robustness ( $y$  axis) of the random (blue) *versus* stereotyped (green) development under different  
819 rate of necrosis ( $x$  axis), as estimated by the model in **E**. The statistical significance of student's  $t$ -  
820 test is indicated as \*\*\*: $P<0.001$ .

#### 821 **Figure 5. The heritability and of the stereotyped developmental** 822 **program**

823 **(A)** The input (top) for DELTA includes two CLTs (query and subject) and the expression profiles  
824 of all terminal cells on these CLTs. DELTA uses a dynamic programming procedure (middle) to  
825 compare the two CLTs and identify homeomorphic sub-CLTs. The procedure has three phases,  
826 including (i) a cell pair scoring stage, (ii) a forward stage that maximizes the alignment scores by  
827 finding the best correspondence between terminal cells, and (iii) a backtracking stage for extracting  
828 the alignment behind the maximized scores. The output (lower right) is one or more aligned sub-  
829 CLTs ordered by decreasing alignment scores. See **Methods** and **Supplemental Texts** for more

830 details. **(B)** A circular plot of the top 100 sub-CLT pairs found by mDELTA in each of the six  
831 pairwise comparisons among the CLTs from the three differentiating samples. In the outer circle,  
832 each sub-CLT is represented by a dot, with the color indicating its source sample. Each pair of  
833 homeomorphic sub-CLTs identified by mDELTA is shown by curved links between two  
834 corresponding dots, where inter-sample pairs/links are colored the same as the sample used as the  
835 query CLT, and intra-sample pairs/links are colored purple. A dot's size indicates how many links it  
836 has. Only sub-CLTs with at least one link are included. **(C)** One highly recurrent tree motif found  
837 in all three samples is shown by “densitree” plots. All sub-CLTs homeomorphic to a specific  
838 reference sub-CLT are extracted from mDELTA results in panel **B**. They were separated by their  
839 source sample as indicated on top of each plot. In each plot, the mDELTA-aligned topological  
840 structure of each sub-CLT (including the reference sub-CLT) is drawn with transparency on the left  
841 so that common topologies can be seen as darker lines. Each column of tiles on the right shows the  
842 DELTA-aligned terminal cell types (colored as the label on top) on one of the homeomorphic sub-  
843 CLTs. The left-most column of tiles is always the reference sub-CLT. The number at the bottom  
844 indicates the number of sub-CLTs found as homeomorphic to the reference sub-CLT.

845

846

847 **Supplementary Information**

- 848 Video S1. A typical alveolosphere formed by the directed differentiation procedure
- 849 Figure S1. Reliability of the directed differentiation and experimental lineage tracing
- 850 Figure S2. Quality of the simultaneous directed differentiation and lineage tracing
- 851 Figure S3. Transcriptional divergence among sub-CLTs
- 852 Figure S4. Transcriptional memory in individual cell types
- 853 Figure S5. Heritability of descendent cell type compositions and single-cell transcriptomes
- 854 Table S1. Designed lineage barcode sites and sgRNAs
- 855 Table S2. Summary statistics of single-cell transcriptomes
- 856 Table S3. Number of passes required *versus* sequencing quality of PacBio HiFi-reads
- 857 Table S5. List of unique (cell barcode and UMI) lineage barcode alleles and their editing events
- 858 Table S4. List of the representative lineage barcode of each cell by their editing events
- 859 Table S6. Structure of the constructed cell lineage trees
- 860 Table S7. List of analyzed GO terms enriched with stage-specific DEGs
- 861 Table S8. List of primers used
- 862

863 **Video S1. Alveolospheres developed on day 15 of the *in vitro* directed**  
864 **differentiation**

865 Following the sorting and replating of NKX2-1<sup>+</sup> lung progenitors on day 15, alveolospheres are  
866 developed in 3D Matrigel culture with CK+DCI media within 3-7 days and maintained in CK+DCI  
867 media for weeks. These spheres are examined by Z stack live images on the Leica DMI8  
868 fluorescence microscope.

869 **Figure S1. Reliability of the directed differentiation and experimental**  
870 **lineage tracing**

871 (A) Verification of *in vitro* directed differentiation toward PLP at hallmark steps ranging from day  
872 0 (hESC), day 3 (definitive endoderm, with flow cytometry results below), day 6 (anterior foregut  
873 endoderm), day 15 (primordial lung progenitor, with flow cytometry results below) to day 20 (Lung  
874 alveolar type II epithelial cells, fluorescence imaging) by using the BU3 NGST (NKX2-1-GFP;  
875 SFTPC-tdTomato) iPS cell line. Bars at the bottom right corners indicate 50  $\mu$ m. (B) Key steps of  
876 experimental lineage tracing for *in vitro* directed differentiation from several (~10) lineage tracer  
877 hESCs (sgRNA-mCherry; lineage barcode-GFP) to PLP are shown at the bottom. The process began  
878 with the selection of traceable colonies (GFP<sup>+</sup> and mCherry<sup>+</sup>) by a 7-day culture, during which a  
879 brief Cas9 induction was applied to uniquely label the ancestor cells by the resulting mutations on  
880 the lineage barcode. The selected colonies were digested and plated again at ~10 cells per well for  
881 the subsequent directed differentiation culture, which lasted for 10 days to produce ~5000 cells. On  
882 top, a typical sample is shown with bright field images at several timepoints, with the scale bar  
883 placed at the bottom right corner. (C) Cas9 (left), the lineage barcode (middle), and sgRNAs (right)  
884 are sufficiently expressed/induced in lineage tracer hESCs. The error bars indicate the standard error  
885 of three replicates. (D) The frequency of inter-site (red) and non-inter-site (blue) deletions found in  
886 edited barcode of lineage tracer hESCs. (E) The most frequent editing events are evenly dispersed  
887 within the lineage barcode. Editing events are named (x axis) by length (the number before I/D),  
888 type (I: insertion; D: deletion) and position (the number after the underline). (F) The frequency of  
889 inter-site deletion events of different lengths (in terms of the number of editing sites) among all  
890 inter-site deletion events.

891 **Figure S2. Quality of the simultaneous directed differentiation and**  
892 **lineage tracing.**

893 (A) Morphology and fluorescence imaging of differentiating/CBRAD5 and hESC self-renewal  
894 samples on day 10. (B) Overview of single-cell transcriptomes measured for differentiating and

895 hESC samples. **(C)** Feature plots for average expression level of marker DEGs found in previous  
896 microarray-based transcriptomes of specific developmental stages<sup>33</sup>, based on UMAP visualization  
897 of single-cell transcriptomes as described in **Figure 1B**. The specific marker genes were listed below  
898 the title and above the plot. **(D)** Sequencing quality and accuracy of PacBio HiFi-reads given the  
899 required number of passes. Error bars indicates standard deviation among ZMWs. **(E)** Tree  
900 representation of the CLTs shown in **Figure 1F**. **(F)** Bootstrap support percentages for the internal  
901 nodes of the CLTs in each sample are presented as histograms. The sample names and median  
902 bootstrap support are shown in the plot titles and in-plot texts, respectively, with the median support  
903 further indicated by a red vertical dashed line.

### 904 **Figure S3. Transcriptional divergence among sub-CLTs**

905 **(A)** Schematic diagram for the normalized depth of a node (see **Methods**). **(B)** Same as **Figure 2F**  
906 except that the analyses were limited to specific GO terms indicated on top of each panel.

### 907 **Figure S4. Transcriptional memory in individual cell types**

908 **(A and B)** Similar to **Figure 3B and C**, except that each major cell type was plotted separately.

### 909 **Figure S5. Heritability of descendent cell type compositions and** 910 **single-cell transcriptomes**

911 The Spearman's Correlation Coefficient ( $y$  axis) between relatedness and phenotypic divergence, a  
912 proxy of the phenotypic heritability, is calculated for all pairs of relevant nodes (see **Methods**). For  
913 the differentiating samples ( $x$  axis), the correlation and therefore phenotypic heritability is always  
914 stronger for the phenotype of descendent cell type components (dots) compared to single-cell  
915 transcriptomes (triangles). The correlation is nevertheless indistinguishable between the two  
916 phenotypes in the non-differentiating sample. A filled or empty point is used to indicate whether the  
917 correlation is statistically significant. A slight offset has been applied to the points of the two  
918 phenotypes in order to avoid overplotting.

919

## 920 **References**

- 921 1. Waddington, C.H. CANALIZATION OF DEVELOPMENT AND THE INHERITANCE OF ACQUIRED  
922 CHARACTERS. *Nature* **150**, 563-565 (1942).
- 923 2. Kitano, H. Biological robustness. *Nat Rev Genet* **5**, 826-37 (2004).
- 924 3. Nijhout, H.F. The nature of robustness in development. *Bioessays* **24**, 553-63 (2002).
- 925 4. Siegal, M.L. & Leu, J.Y. On the Nature and Evolutionary Impact of Phenotypic Robustness  
926 Mechanisms. *Annu Rev Ecol Evol Syst* **45**, 496-517 (2014).
- 927 5. Felix, M.A. & Wagner, A. Robustness and evolution: concepts, insights and challenges from a  
928 developmental model system. *Heredity (Edinb)* **100**, 132-40 (2008).
- 929 6. Gibson, G. & Lacey, K.A. Canalization and Robustness in Human Genetics and Disease. *Annu*  
930 *Rev Genet* **54**, 189-211 (2020).
- 931 7. Rutherford, S.L. & Lindquist, S. Hsp90 as a capacitor for morphological evolution. *Nature* **396**,  
932 336-42 (1998).
- 933 8. Hornstein, E. & Shomron, N. Canalization of development by microRNAs. *Nat Genet* **38 Suppl**,  
934 S20-4 (2006).
- 935 9. Wu, C.I., Shen, Y. & Tang, T. Evolution under canalization and the dual roles of microRNAs: a  
936 hypothesis. *Genome Res* **19**, 734-43 (2009).
- 937 10. Siciliano, V. *et al.* MiRNAs confer phenotypic robustness to gene networks by suppressing  
938 biological noise. *Nat Commun* **4**, 2364 (2013).
- 939 11. Levy, S.F. & Siegal, M.L. Network hubs buffer environmental variation in *Saccharomyces*  
940 *cerevisiae*. *PLoS Biol* **6**, e264 (2008).
- 941 12. Mo, N. *et al.* Bidirectional Genetic Control of Phenotypic Heterogeneity and Its Implication  
942 for Cancer Drug Resistance. *Mol Biol Evol* **38**, 1874-1887 (2021).
- 943 13. El-Samad, H. Biological feedback control-Respect the loops. *Cell Syst* **12**, 477-487 (2021).
- 944 14. Kafri, R., Levy, M. & Pilpel, Y. The regulatory utilization of genetic redundancy through  
945 responsive backup circuits. *Proc Natl Acad Sci U S A* **103**, 11653-8 (2006).
- 946 15. Xiao, L., Fan, D., Qi, H., Cong, Y. & Du, Z. Defect-buffering cellular plasticity increases  
947 robustness of metazoan embryogenesis. *Cell Syst* **13**, 615-630 e9 (2022).
- 948 16. Green, R.M. *et al.* Developmental nonlinearity drives phenotypic robustness. *Nat Commun* **8**,  
949 1970 (2017).
- 950 17. Yang, J.R., Ruan, S. & Zhang, J. Determinative developmental cell lineages are robust to cell  
951 deaths. *PLoS Genet* **10**, e1004501 (2014).
- 952 18. Wagner, D.E. & Klein, A.M. Lineage tracing meets single-cell omics: opportunities and  
953 challenges. *Nat Rev Genet* **21**, 410-427 (2020).
- 954 19. Li, Z. *et al.* Reconstructing cell lineage trees with genomic barcoding: approaches and  
955 applications. *J Genet Genomics* (2023).
- 956 20. Waddington, C. *The Strategy of the Genes*. Allen. (Unwin, London, 1957).
- 957 21. Pujadas, E. & Feinberg, A.P. Regulated noise in the epigenetic landscape of development and  
958 disease. *Cell* **148**, 1123-31 (2012).
- 959 22. Sulston, J.E., Schierenberg, E., White, J.G. & Thomson, J.N. The embryonic cell lineage of the  
960 nematode *Caenorhabditis elegans*. *Dev Biol* **100**, 64-119 (1983).
- 961 23. Raj, B., Gagnon, J.A. & Schier, A.F. Large-scale reconstruction of cell lineages using single-cell  
962 readout of transcriptomes and CRISPR-Cas9 barcodes by scGESTALT. *Nat Protoc* **13**, 2685-2713 (2018).
- 963 24. Furman, D.P. & Bukharina, T.A. How *Drosophila melanogaster* Forms its Mechanoreceptors.

964 *Curr Genomics* **9**, 312-23 (2008).

965 25. Koch, U., Lehal, R. & Radtke, F. Stem cells living with a Notch. *Development* **140**, 689-704  
966 (2013).

967 26. Kilimnik, G., Jo, J., Periwai, V., Zielinski, M.C. & Hara, M. Quantification of islet size and  
968 architecture. *Islets* **4**, 167-72 (2012).

969 27. Salipante, S.J., Kas, A., McMonagle, E. & Horwitz, M.S. Phylogenetic analysis of  
970 developmental and postnatal mouse cell lineages. *Evol Dev* **12**, 84-94 (2010).

971 28. Anderson, D.J. *et al.* Simultaneous brain cell type and lineage determined by scRNA-seq  
972 reveals stereotyped cortical development. *Cell Syst* **13**, 438-453 e5 (2022).

973 29. Hughes, N.W. *et al.* Machine-learning-optimized Cas12a barcoding enables the recovery of  
974 single-cell lineages and transcriptional profiles. *Mol Cell* **82**, 3103-3118 e8 (2022).

975 30. Jacob, A. *et al.* Differentiation of Human Pluripotent Stem Cells into Functional Lung Alveolar  
976 Epithelial Cells. *Cell Stem Cell* **21**, 472-488 e10 (2017).

977 31. Chen, F. *et al.* Phylogenetic Comparative Analysis of Single-Cell Transcriptomes Reveals  
978 Constrained Accumulation of Gene Expression Heterogeneity during Clonal Expansion. *Mol Biol Evol*  
979 **40**(2023).

980 32. He, Z. *et al.* Lineage recording in human cerebral organoids. *Nat Methods* **19**, 90-99 (2022).

981 33. Hawkins, F. *et al.* Prospective isolation of NKX2-1-expressing human lung progenitors derived  
982 from pluripotent stem cells.

983 34. Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by  
984 pseudotemporal ordering of single cells.

985 35. Bonasio, R., Tu, S. & Reinberg, D. Molecular signals of epigenetic states. *Science* **330**, 612-6  
986 (2010).

987 36. Shaffer, S.M. *et al.* Memory Sequencing Reveals Heritable Single-Cell Gene Expression  
988 Programs Associated with Distinct Cellular Behaviors. *Cell* **182**, 947-959 e17 (2020).

989 37. Dublanche, Y., Michalodimitrakis, K., Kummerer, N., Foglierini, M. & Serrano, L. Noise in  
990 transcription negative feedback loops: simulation and experimental analysis. *Mol Syst Biol* **2**, 41 (2006).

991 38. Sharon, E. *et al.* Probing the effect of promoters on noise in gene expression using thousands  
992 of designed sequences. *Genome Res* **24**, 1698-706 (2014).

993 39. Eisele, A.S., Tarbier, M., Dormann, A.A., Pelechano, V. & Suter, D.M. Barcode-free prediction  
994 of cell lineages from scRNA-seq datasets. *bioRxiv*, 2022.09.20.508646 (2022).

995 40. Kuleshov, M.V. *et al.* Enrichr: a comprehensive gene set enrichment analysis web server 2016  
996 update. *Nucleic Acids Res* **44**, W90-7 (2016).

997 41. Pan, G. & Thomson, J.A. Nanog and transcriptional networks in embryonic stem cell  
998 pluripotency. *Cell Res* **17**, 42-9 (2007).

999 42. Heslop, J.A., Pournasr, B., Liu, J.T. & Duncan, S.A. GATA6 defines endoderm fate by controlling  
1000 chromatin accessibility during differentiation of human-induced pluripotent stem cells. *Cell Rep* **35**,  
1001 109145 (2021).

1002 43. Yuan, M. *et al.* Alignment of Cell Lineage Trees Elucidates Genetic Programs for the  
1003 Development and Evolution of Cell Types. *iScience* **23**, 101273 (2020).

1004 44. Li, Y. *et al.* A full-body transcription factor expression atlas with completely resolved cell  
1005 identities in *C. elegans*. *Nat Commun* **15**, 358 (2024).

1006 45. Murray, J.I. *et al.* Multidimensional regulation of gene expression in the *C. elegans* embryo.  
1007 *Genome Res* **22**, 1282-94 (2012).

1008 46. Khalil, O. *et al.* Reconstructing unobserved cellular states from paired single-cell lineage  
1009 tracing and transcriptomics data. *bioRxiv*, 2021.05.28.446021 (2021).

1010 47. Haseman, J.K. & Elston, R.C. The investigation of linkage between a quantitative trait and a  
1011 marker locus. *Behav Genet* **2**, 3-19 (1972).

1012 48. Yang, J. *et al.* Common SNPs explain a large proportion of the heritability for human height.  
1013 *Nat Genet* **42**, 565-9 (2010).

1014 49. Tran, M., Askary, A. & Elowitz, M.B. Lineage motifs as developmental modules for control of  
1015 cell type proportions. *Dev Cell* (2024).

1016 50. McKenna, A. *et al.* Whole-organism lineage tracing by combinatorial and cumulative genome  
1017 editing. *Science* **353**, aaf7907 (2016).

1018 51. Doench, J.G. *et al.* Optimized sgRNA design to maximize activity and minimize off-target  
1019 effects of CRISPR-Cas9. *Nat Biotechnol* **34**, 184-191 (2016).

1020 52. Magoc, T. & Salzberg, S.L. FLASH: fast length adjustment of short reads to improve genome  
1021 assemblies. *Bioinformatics* **27**, 2957-63 (2011).

1022 53. Langmead, B. & Salzberg, S.L. Fast gapped-read alignment with Bowtie 2. *Nat Methods* **9**,  
1023 357-9 (2012).

1024 54. Edgar, R.C. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
1025 *Nucleic Acids Res* **32**, 1792-7 (2004).

1026 55. Sahabian, A., Dahlmann, J., Martin, U. & Olmer, R. Production and cryopreservation of  
1027 definitive endoderm from human pluripotent stem cells under defined and scalable culture conditions.  
1028 *Nat Protoc* **16**, 1581-1599 (2021).

1029 56. Zheng, G.X. *et al.* Massively parallel digital transcriptional profiling of single cells. *Nat*  
1030 *Commun* **8**, 14049 (2017).

1031 57. Dobin, A. *et al.* STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21 (2013).

1032 58. Hao, Y. *et al.* Integrated analysis of multimodal single-cell data. *Cell* **184**, 3573-3587.e29  
1033 (2021).

1034 59. Guo, L. *et al.* Resolving Cell Fate Decisions during Somatic Cell Reprogramming by Single-Cell  
1035 RNA-Seq. *Mol Cell* **73**, 815-829 e7 (2019).

1036 60. Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony.  
1037 *Nat Methods* **16**, 1289-1296 (2019).

1038 61. Barrett, T. *et al.* NCBI GEO: archive for functional genomics data sets--update. *Nucleic Acids*  
1039 *Res* **41**, D991-5 (2013).

1040 62. Umkehrer, C. *et al.* Isolating live cell clones from barcoded populations using CRISPRa-  
1041 inducible reporters. *Nature Biotechnology* (2020).

1042 63. Anderson, M.J. A new method for non-parametric multivariate analysis of variance. *Austral*  
1043 *ecology* **26**, 32-46 (2001).

1044 64. Ozgode Yigin, B. & Saygili, G. Effect of distance measures on confidences of t-SNE embeddings  
1045 and its implications on clustering for scRNA-seq data. *Sci Rep* **13**, 6567 (2023).

1046 65. Kim, T. *et al.* Impact of similarity metrics on single-cell RNA-seq data clustering. *Brief*  
1047 *Bioinform* **20**, 2316-2326 (2019).

1048