

# PhyloVelo enhances transcriptomic velocity field mapping using monotonically expressed genes

Zheng Hu (✉ [zheng.hu@siat.ac.cn](mailto:zheng.hu@siat.ac.cn))

Shenzhen Institutes of Advanced Technology, Chinese Academy of Sciences <https://orcid.org/0000-0003-1552-0060>

**Kun Wang**

Xiamen University <https://orcid.org/0009-0009-6503-4678>

**Liangzhen Hou**

University of Macau

**Xin Wang**

Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences <https://orcid.org/0000-0001-7440-1075>

**Xiangwei Zhai**

Sun Yat-Sen University <https://orcid.org/0000-0003-2967-5573>

**Zhaolian Lu**

Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences

**Zhike Zi**

Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences <https://orcid.org/0000-0002-7601-915X>

**Weiwei Zhai**

Institute of Zoology, Chinese Academy of Sciences <https://orcid.org/0000-0001-7938-0226>

**Xionglei He**

Sun Yat-Sen University <https://orcid.org/0000-0003-1050-802X>

**Christina Curtis**

Stanford University School of Medicine <https://orcid.org/0000-0003-0166-3802>

**Da Zhou**

Xiamen University <https://orcid.org/0000-0002-0272-6644>

---

**Article**

**Keywords:**

**Posted Date:** August 1st, 2023

**DOI:** <https://doi.org/10.21203/rs.3.rs-2197712/v1>

**License:** © ⓘ This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

**Additional Declarations:** **Yes** there is potential Competing Interest. C.C. is an advisor and stockholder in Grail, Ravel, and DeepCell and an advisor to Genentech, Bristol Myers Squibb, 3T Biosciences, and NanoString. Other authors declare no competing interests.

---

**Version of Record:** A version of this preprint was published at Nature Biotechnology on July 31st, 2023. See the published version at <https://doi.org/10.1038/s41587-023-01887-5>.

1 **Ed summary: A new velocity model improves cell-fate mapping with**  
2 **lineage-traced scRNA-seq data.**

3

4 **PhyloVelo enhances transcriptomic velocity field mapping using**  
5 **monotonically expressed genes**

6 Kun Wang<sup>1,2</sup>, Liangzhen Hou<sup>1,3</sup>, Xin Wang<sup>1</sup>, Xiangwei Zhai<sup>4</sup>, Zhaolian Lu<sup>1</sup>, Zhike Zi<sup>1</sup>, Weiwei  
7 Zhai<sup>5,6</sup>, Xionglei He<sup>4</sup>, Christina Curtis<sup>7,8,9</sup>, Da Zhou<sup>2,10\*</sup>, Zheng Hu<sup>1\*</sup>

8 <sup>1</sup>CAS Key Laboratory of Quantitative Engineering Biology, Shenzhen Institute of Synthetic  
9 Biology, Shenzhen Institute of Advanced Technology, Chinese Academy of Sciences,  
10 Shenzhen, China

11 <sup>2</sup>School of Mathematical Sciences, Xiamen University, Xiamen, China

12 <sup>3</sup>Faculty of Health Sciences, University of Macau, Taipa, Macau, China

13 <sup>4</sup>MOE Key Laboratory of Gene Function and Regulation, State Key Laboratory of Biocontrol,  
14 School of Life Sciences, Sun Yat-Sen University, Guangzhou, China

15 <sup>5</sup>CAS Key Laboratory of Zoological Systematics and Evolution, Institute of Zoology, Chinese  
16 Academy of Sciences, Beijing, China

17 <sup>6</sup>Center for Excellence in Animal Evolution and Genetics, Chinese Academy of Sciences,  
18 Kunming, China

19 <sup>7</sup>Department of Medicine, Division of Oncology, Stanford University School of Medicine,  
20 Stanford, CA, USA

21 <sup>8</sup>Department of Genetics, Stanford University School of Medicine, Stanford, CA, USA

22 <sup>9</sup>Stanford Cancer Institute, Stanford University School of Medicine, Stanford, CA, USA

23 <sup>10</sup>National Institute for Data Science in Health and Medicine, Xiamen University, Xiamen,  
24 China

25

26 \* Correspondence to: [zhouda@xmu.edu.cn](mailto:zhouda@xmu.edu.cn) (D.Z.); [zheng.hu@siat.ac.cn](mailto:zheng.hu@siat.ac.cn) (Z.H.)

27

## 28 **Abstract**

29 Single-cell RNA-sequencing (scRNA-seq) is a powerful approach for studying cellular  
30 differentiation, but accurately tracking cell-fate transitions can be challenging, especially in  
31 disease conditions. Here, we introduce PhyloVelo, a computational framework that  
32 estimates the velocity of transcriptomic dynamics by using monotonically expressed genes  
33 (MEGs), or genes with expression patterns that either increase or decrease, but don't cycle,  
34 through phylogenetic time. Through integration of scRNA-seq data with lineage information,  
35 PhyloVelo identifies MEGs and reconstructs a transcriptomic velocity field. We validate  
36 PhyloVelo using simulated data and *C. elegans* ground-truth data, successfully recovering  
37 linear, bifurcated, and convergent differentiations. Applying PhyloVelo to seven lineage-  
38 traced scRNA-seq datasets, generated via CRISPR/Cas9 editing, lentiviral barcoding or  
39 immune repertoire profiling, demonstrates its high accuracy and robustness in inferring  
40 complex lineage trajectories, while outperforming RNA velocity. Additionally, we discover  
41 that MEGs across tissues and organisms share similar functions in translation and ribosome  
42 biogenesis.

## 43 **Main**

44 Organism development and disease progression both involve serial cell-fate transitions  
45 upon repeated cell divisions. Essentially, all cells in an organism are related by a  
46 phylogenetic tree where the root represents the zygote, the branches represent cell  
47 divisions, and the leaves represent the terminal cells at various phenotypic states (e.g. cell  
48 types)<sup>1-4</sup>. To understand how cell fate is determined, it is important to identify the order of

49 cell-state transitions acting in the lineage tree and the underlying gene regulatory  
50 mechanisms that precipitate these transitions<sup>5</sup>.  
51  
52 Single-cell RNA sequencing (scRNA-seq) has been a powerful approach to study cellular  
53 differentiations<sup>6-10</sup>. However, the transcriptomic trajectories may or may not be equivalent to  
54 the true lineage paths of a progenitor population<sup>11-14</sup>. One example is convergent  
55 differentiation where distinct progenitors can converge on the same terminal state<sup>15-19</sup>. In  
56 this case, similar cellular states do not reflect a closer lineage relationship<sup>13</sup>. Moreover,  
57 predicting the fate directions often requires prior knowledge of the initial/terminal cell types<sup>9</sup>  
58 or relies on the information of gene expression diversity during development<sup>10</sup>, thus limiting  
59 their applications to normally differentiating systems<sup>20</sup>. Abnormal development or disease  
60 progression often involves noncanonical cell-fate transitions such as dedifferentiation and  
61 transdifferentiation<sup>21</sup>, while tackling these processes is still challenging with current  
62 approaches. RNA velocity<sup>22, 23</sup> provides a powerful framework to predict cellular state  
63 transitions by leveraging the internal kinetics of spliced/unspliced RNAs, and can be readily  
64 applied to diseased or perturbed conditions. However, the intrinsic high dynamics of RNA  
65 kinetics including transcription, splicing and degradation often violates the constant rate  
66 assumptions in the model, which can lead to uncertain estimates<sup>23-25</sup>. Taken together, cell-  
67 state transitions are challenging to distinguish using transcriptomic data alone<sup>11-14</sup>.  
68  
69 The recent use of CRISPR/Cas9 editing to record cell lineages offers an opportunity to  
70 reconstruct the cell lineage tree at whole-organism or whole-organ level<sup>26-29</sup>. Importantly,  
71 simultaneous analysis of single-cell transcriptomes and lineage tree makes it possible to  
72 uncover complex developmental dynamics, as well as the molecular mechanisms, of cell  
73 fate commitment<sup>13, 30-35</sup>. For instance, CRISPR lineage tracing has enabled the identification  
74 of transcriptional convergence of endodermal cells from both extra-embryonic and  
75 embryonic origins during mouse embryogenesis<sup>32</sup>. Although the significance of CRISPR

76 lineage tracing coupled with single-cell transcriptomics has been widely acknowledged in  
77 developmental biology and somatic evolution<sup>13, 29, 33</sup>, computationally integrating dual  
78 information for reconstructing cellular trajectories is challenging, partially due to the distinct  
79 data modalities. Previous effort such as the CoSpar algorithm<sup>36</sup> has been made to use the  
80 paired scRNA-seq and lineage information to infer transition maps and predict fate bias of  
81 progenitor cells, which is better fit for static barcoding information (e.g. LARRY system<sup>37</sup>).  
82 Another algorithm lineageOT has taken advantage of lineage tree for trajectory inference<sup>38</sup>,  
83 however, this relies on time-course scRNA-seq data and an invariant cell lineage tree. This  
84 type of data is currently only available in *C. elegans*<sup>14</sup>, thus preventing its application to  
85 more common datasets such as CRISPR-based lineage tracing data.

86

87 In this study, we described a method to systematically map cell-fate transitions by using  
88 both single-cell transcriptomic and lineage information. Our method, called PhyloVelo,  
89 leverages monotonically expressed genes (MEGs) along cell divisions to quantify the  
90 transcriptomic velocity fields from lineage-resolved single-cell RNA-seq data (**Fig. 1**). We  
91 verified the capacity and robustness of PhyloVelo to resolve complex lineage structures in  
92 comprehensive simulations and real lineage tracing data in embryo development (*C.*  
93 *elegans* and mouse embryos), tumor evolution (both initiation and metastasis), *in vitro*  
94 hematopoiesis and intratumoral T cell dynamics. We further demonstrated that the velocities  
95 estimated from one scRNA-seq dataset were sufficiently robust to infer the lineage  
96 trajectory with independent datasets in similar biological conditions, even in the absence of  
97 lineage information. Finally, we found MEGs were strongly enriched in ribosome-mediated  
98 processes across tissues and organisms, thus exposing an internal clock-like gene  
99 expression program during cell proliferation and differentiation.

## 100 **Results**

### 101 ***A transcriptomic velocity field reconstructed by MEGs***

102 RNA velocity<sup>22, 23</sup> exploits the kinetics of spliced/unspliced RNAs to estimate the time  
103 derivative (or velocity) of single-cell gene expression states ( $ds/dt$  with  $s$  representing the  
104 high-dimensional expression state and  $t$  representing time). This enables the prediction of  
105 future gene expression states and reconstruction of a velocity vector field of cellular state  
106 transitions on low-dimensional embedding. We anticipated that other measurements of  
107  $ds/dt$  can be similarly employed to establish the velocity field. In particular, lineage tracing  
108 by endogenous mutations or evolving barcodes (e.g. CRISPR/Cas9 editing) reconstructs a  
109 cell phylogenetic tree that records the cell division history from a common progenitor (**Fig.**  
110 **1a**). Because more cell divisions usually indicate more advanced differentiation stages in a  
111 stem cell hierarchy, the phylogenetic time potentially associates with the differentiation time  
112 of individual cells. To bridge differentiation and phylogenetic time, we focused on a group of  
113 genes (namely MEGs) whose expressions increase or decrease monotonically over  
114 phylogenetic time (**Fig. 1a-b**). The monotonic feature of MEGs enables them to serve a  
115 “clock” of cell differentiations (**Fig. 1a**).

116

117 To identify MEGs and estimate their expression velocity  $ds/dt$ , we first sought to estimate  
118 the latent gene expression of each gene at a phylogenetic time, which was akin to the latent  
119 variables in single-cell RNA-seq denoising<sup>39</sup>. Inspired by the classic models of trait evolution  
120 in phylogenetics<sup>40</sup>, we modeled the continuously varying gene expressions by a diffusion  
121 process (also called stochastic differential equation) (**Methods**). Each gene has a specific  
122 rate of expression dynamics, namely the drift coefficient  $v(t, z_t)$  per unit of time, where  $t$  is  
123 the time started from the root of tree and  $z_t$  is the latent expression for a gene at time  $t$ .  
124 MEGs were identified by their significant association (Pearson’s correlation,  $q < 0.05$ )  
125 between the latent expressions,  $\mathbf{Z} = (\hat{z}_1, \hat{z}_2, \dots, \hat{z}_n)$ , and the time from terminal cells to the  
126 root,  $\mathbf{B} = (b_1, b_2, \dots, b_n)$ , where  $n$  was the cell number (**Methods**). It’s worth noting that  
127 MEGs are defined with respect to the observed phylogenetic time range  $\mathbf{B}$ , which means  
128 the monotonic expression might not retain out of  $\mathbf{B}$ . Moreover, although we focused on

129 linear MEGs with phylogenetic time, non-linear MEGs can be also identified if the linear  
130 regression is statistically significant. We will show later by simulations that linear  
131 approximation is technically sound way to accurately map cell trajectories. The drift  
132 coefficients of all  $G$  MEGs in a dataset,  $\mathbf{v} = (v_1, v_2, \dots, v_G)$ , were thus referred to as  
133 phylogenetic velocity (or PhyloVelo). As shown in **Fig. 1b-c**, phylogenetic velocity can be  
134 used to predict the past expression state of each cell before a unit of time  $\Delta t$  (one cell  
135 division or mutation),  $\mathbf{s}^* = \mathbf{s} - \mathbf{v}\Delta t$ . Similar to RNA velocity<sup>22, 23</sup>, phylogenetic velocity  $\mathbf{v}$  can  
136 also be projected into low dimensional embedding such as t-distributed Stochastic Neighbor  
137 Embedding (tSNE) or Uniform Manifold Approximation and Projection (UMAP), which  
138 reconstructs the velocity vector fields (**Fig. 1d**). Unlike RNA velocity where the velocity fields  
139 point to future extrapolated states, phylogenetic velocity fields point to the instantaneously  
140 past states, thus reconstructing fate-transition map in backward directions (**Fig. 1d**).

#### 141 ***PhyloVelo recovers complex lineages in simulations and *C. elegans****

142 We next sought to test PhyloVelo with simulation data where various lineage structures  
143 were considered, including linear, bifurcated and convergent differentiations (**Fig. 2a-c**). A  
144 lineage-imbedded scRNA-seq simulator PROSSTT<sup>41</sup> was modified to record individual cell  
145 divisions and generate single-cell UMI counts simultaneously (**Methods**). To model cell  
146 differentiations, different cell types each showing a characteristic gene expression program  
147 were simulated in the three lineage structures, respectively (**Fig 2a-c, Supplementary Fig.**  
148 **1**). We also simulated random mutations that occur during cell divisions, which allows to  
149 build mutation-based cell lineage tree (**Methods**). Of note, simulations showed that different  
150 cell types were highly intermixed on the lineage tree (**Fig. 2a-c**), a phenomenon that  
151 appears to be common for organ development across diverse species, such as flies<sup>42</sup>,  
152 zebrafish<sup>43</sup> and mice<sup>44</sup>. Nevertheless, the dimensionality reduction embedding (tSNE) of  
153 simulated scRNA-seq data recapitulated the actual lineage structures (**Fig. 2a-c**).

154



155 By applying PhyloVelo to the simulation data, we first found that MEGs following either  
156 increasing or decreasing dynamics can be robustly detected with our algorithm (**Extended**  
157 **Data Fig. 1**). With the estimated phylogenetic velocities of MEGs, PhyloVelo mapped the  
158 state transitions in backward directions which point to the extrapolated past states of  
159 individual cells on low dimensional embedding (**Fig. 2d-f**). We used two quantitative metrics  
160 to systematically evaluate the performance of PhyloVelo with simulation datasets generated  
161 under a variety of parameters and conditions: (1) precision rate of MEG identification,  $M_1$ ;  
162 (2) Pearson's correlation of estimated velocity directions using the identified MEGs vs using  
163 the genuine MEGs,  $M_2$ . As shown in **Extended Data Fig. 1** and **Supplementary Fig. 2-5**,  
164 we found PhyloVelo inferences were highly robust to the cell number (mean  $M_1 > 80\%$  and  
165 mean  $M_2 > 90\%$  even at low cell number of 100 cells), non-linear dynamics of MEGs (mean  
166  $M_1 > 75\%$ ,  $M_2$  was not available here), the methods of dimensionality reduction embedding  
167 (mean  $M_2 > 90\%$  for both tSNE and UMAP), and the sparsity level of single-cell data (mean  
168  $M_1 > 90\%$  and  $M_2 > 90\%$  even for  $\sim 0.2$  UMIs per cell per gene). Interestingly, both  $M_1$  and  $M_2$   
169 remained high when the number of genuine MEGs exceeded 50 ( $M_1 > 85\%$  and  $M_2 > 90\%$ ,  
170 **Extended Data Fig. 1, Supplementary Fig. 6**). Mathematical analysis verified a small  
171 angle (upper bound  $< 37^\circ$ ) between the estimated and true velocity vectors at 50 MEGs and  
172 a precision rate of  $M_1 = 80\%$  (**Supplementary Note**). Using a stringent threshold ( $M_2 = 95\%$ )  
173 for good performance, our simulations revealed at least 35, 82 and 43 MEGs were required  
174 for linear, bifurcated and convergent differentiation model, respectively. In summary, our  
175 comprehensive benchmarking and mathematical analysis demonstrated the high accuracy  
176 and robustness of PhyloVelo to systematically map cell-state trajectories.

177 In addition, we found the estimated phylogenetic velocities based on the mutation-based  
178 phylogenies and the ground-truth division history were highly concordant, although  
179 inaccurate velocity estimations in local lineages were noted when the mutation rate was  
180 rather low (mean mutation rate  $u = 0.1$  per cell division) (**Supplementary Fig. 7-8**).

181 Importantly, phylogenetic velocity estimates were robust to different phylogenetic methods  
182 (e.g. maximum likelihood, neighbor joining or maximum parsimony), which was because  
183 these methods gave highly consistent inferences on the phylogenetic distances  
184 (**Supplementary Fig. 9**). Finally, while classic trajectory inference algorithms such as  
185 monocle3<sup>7</sup>, slingshot<sup>45</sup>, and PAGA<sup>46</sup> can accurately identify the backbones of linear and  
186 bifurcated lineage structures, only PAGA was able to identify the circular structure in  
187 convergent differentiation (**Supplementary Fig. 10**). In fact, additional information on initial  
188 or terminal cell types is needed to define the directions using the aforementioned three  
189 algorithms. This is expected because most trajectory inference methods are inadequate for  
190 single-cell datasets containing a convergent trajectory, and also rely on prior information of  
191 initial/terminal cell types<sup>9</sup>.

192 We next applied PhyloVelo to *C. elegans* given that the embryonic lineage tree of this  
193 organism is entirely known<sup>2</sup>. The scRNA-seq data from temporal *C. elegans* embryos are  
194 also available and have been mapped to the invariant lineage tree, as described by Packer  
195 *et al.*<sup>14</sup>. Thus, *C. elegans* is an ideal system to benchmark our method. We focused on the  
196 AB lineage with mostly ectoderm accounting for ~70% of the terminal cells in the embryo  
197 (**Fig. 3a**), which also had the densest single-cell annotations in Packer *et al.* dataset<sup>14</sup>  
198 spanning from generation 5 (32-cell stage) to 12 (threefold stage of development). Since  
199 many nodes on the lineage tree have been sampled multiple times through pooled  
200 sequencing of multiple embryos, one cell was randomly chosen to represent the  
201 corresponding lineage node. This resulted in 298 non-repetitive cells for the AB lineage,  
202 denoted as a single pseudo-embryo. By analyzing the correlation between latent gene  
203 expressions and cell generation times, we identified 326 significant ( $q < 0.05$ ) MEGs with 22  
204 and 304 increasing and decreasing in expressions, respectively (**Fig. 3b, Supplementary**  
205 **Table 1**). This was consistent with the observed global decline in gene expressions during  
206 *C. elegans* embryogenesis<sup>14</sup>.

207 To generate a ground-truth velocity field, each cell was assigned a vector on the UMAP plot  
208 that points to its immediate parental cell in the ground-truth lineage tree (**Fig. 3c**). The  
209 vector fields together tracked the cell lineages back to the earliest cells in development.  
210 Therefore, by comparing the quantitative directions of PhyloVelo velocity fields with the  
211 ground-truth and also RNA velocity fields, we were able to evaluate the accuracy of our  
212 method. The UMAP embedding clearly reflected the differentiation trajectories along cell  
213 divisions (**Fig. 3c-e**) or embryo time (**Fig. 3f**). Surprisingly, we found that RNA velocity  
214 (scVelo - dynamical mode, **Fig. 3d**) failed to recover the expected trajectories, where the  
215 directions were even reversed from the ground truth (**Supplementary Fig. 11a**). In fact, the  
216 scVelo latent time (**Fig. 3g**) was negatively correlated with the real embryo time in early  
217 development before ~300 minutes (**Fig. 3d, Supplementary Fig. 11b**). In contrast, the  
218 directions of phylogenetic velocities recapitulated the actual development orders (**Fig. 3e,**  
219 **Supplementary Fig. 11c-d**). Other single pseudo-embryo data also showed similar results  
220 (**Supplementary Fig. 12**).

221 The RNA velocity fields estimated by pooling all 29,600 AB lineage cells from multiple  
222 embryos were improved (**Fig. 3i-j**), suggesting that RNA velocity estimates had been  
223 hindered by a small cell number. Remarkably, the phylogenetic velocities of 326 MEGs  
224 estimated from single pseudo-embryo data (~300 cells) can be used to accurately infer the  
225 velocity fields for all 29,600 AB lineage cells, even though their lineage trees were not  
226 utilized (**Fig. 3k-l**). In fact, these MEGs were even applicable to non-AB lineage cells such  
227 as hypodermis, body wall muscle (BWM) and pharynx (**Extended Data Fig. 2**).

228 Interestingly, a convergent trajectory for the first row of head body wall muscle (BWM) and  
229 all other BWMs (including C, D and MS lineages) can be identified (**Extended Data Fig.**  
230 **2e**). These results demonstrated a general transcriptomic clock during *C. elegans*  
231 embryogenesis, and also suggest that compiling a reference panel of MEGs will greatly  
232 facilitate the applications of PhyloVelo to conventional scRNA-seq data where lineage data

233 is unavailable. In summary, the benchmarking on comprehensive simulations and *C.*  
234 *elegans* embryo lineages demonstrated the high robustness of PhyloVelo to recover  
235 complex developmental trajectories with phylogeny-resolved scRNA-seq data even with  
236 relatively limited cell numbers.

### 237 ***PhyloVelo resolves multiple-rate kinetics in mouse embryos***

238 We next applied PhyloVelo to a CRISPR/Cas9-based lineage tracing dataset from mouse  
239 early embryos (E8.0 or E8.5), described by Chan *et al.*<sup>32</sup>. This study provided both cell  
240 lineage tree and scRNA-seq data via CRISPR lineage tracing of mouse fertilization through  
241 gastrulation. By analyzing four embryos (embryo 1, 2, 3 and 6, each with 6,328 to 19,071  
242 cells and more than 500 unique barcode alleles), we have identified 426, 460, 420 and 418  
243 MEGs ( $q < 10^{-5}$ ), respectively at whole embryo level (**Extended Data Fig. 3, Supplementary**  
244 **Table 1**). Notably, about 50% (n=212) of MEGs were overlapped by all four embryos and  
245 the phylogenetic velocities of these overlapped MEGs were strongly correlated (Pearson's  
246  $r=0.65-0.95$ , **Extended Data Fig. 3**). Given the generally low barcode diversity in CRISPR  
247 lineage tracing<sup>47</sup> and also high noise in scRNA-seq data, these data actually indicated the  
248 high robustness of PhyloVelo for identifying MEGs. Because of the rapid cell replication in  
249 early embryogenesis, the MEGs identified from one snapshot sample from Chan *et al.*  
250 dataset<sup>32</sup> might only represent a short-term monotonic effect. Nevertheless, we found about  
251 a half (104 out of 212) of overlapped MEGs identified from Chan *et al.* dataset<sup>32</sup> also  
252 showed significant correlation ( $p < 0.05$ ) with the capture time (E6.5-8.5) of temporal mouse  
253 embryos from Pijuan-Sala *et al.*<sup>19</sup> (**Supplementary Fig. 13**). We thus called these 104  
254 genes long-term MEGs, or LT-MEGs. As expected, these 104 LT-MEGs identified from  
255 Chan *et al.* dataset<sup>32</sup> enabled accurate prediction of the entire differentiation trajectories of  
256 mouse embryogenesis with the temporal scRNA-seq data (**Extended Data Fig. 4a-c**).  
257 Remarkably, these LT-MEGs were also highly robust to infer the velocity fields when

258 transferred to mouse brain tissues across broader developmental stages (E7-18) and over  
259 18 cell types<sup>48</sup> (**Extended Data Fig. 4d-f**).

260 To directly compare PhyloVelo with RNA velocity, and also quantify the state-transition  
261 probabilities between cell types, we next focused on the erythroid lineage given its well-  
262 defined differentiation trajectory during mouse gastrulation<sup>49</sup>. Embryo 3 (E8.5) had the  
263 largest cell number and more diverse cell types in erythroid developmental lineages  
264 (n=2,419 cells), thus being selected for a representative case while other embryos (1, 2 and  
265 6) were also analyzed (**Fig. 4a, Extended Data Fig. 5**). RNA velocity failed to identify  
266 hematopoietic/endothelial progenitors as the earliest cell types (**Fig. 4b-c**). In addition, the  
267 fractions of varying cell types only changed slightly along the scVelo latent time (**Fig. 4d**). In  
268 contrast, PhyloVelo correctly predicted the expected trajectory from  
269 hematopoietic/endothelial/primitive blood progenitors to primitive blood early/late based on  
270 the velocity fields and pseudotime (**Fig. 4e-g**). Transferring the MEGs identified from  
271 erythroid cells of embryo 3 to other three embryos also robustly recovered their erythroid  
272 differentiation orders (**Extended Data Fig. 5b-d**). Dynamo<sup>50</sup> was further used to incorporate  
273 PhyloVelo velocity fields, which can quantify the transition probabilities between any two cell  
274 types (**Extended Data Fig. 5e-h**). Dynamo successfully placed hematopoietic endothelial  
275 progenitors and primitive blood late as the starting and ending states, respectively  
276 (**Extended Data Fig. 5i-l**). Importantly, the possible ancestral states of a particular cell type  
277 (non-zero transition probabilities) recapitulated well where the cell type was differentiated.  
278 Finally, by applying CellRank<sup>20</sup> with the input of the PhyloVelo pseudotime, we were also  
279 able to identify the known driver genes underlying erythroid maturation (e.g. *Alas2*, *Bpgm*,  
280 *Car2*, *Slc4a1*, *Hemgn*, **Supplementary Fig. 14**).

281 Studies have shown that multiple-rate kinetics (MURK) of RNA violates the constant  
282 assumptions in RNA velocity model, which might lead to erroneous estimates of velocities<sup>24</sup>,  
283 <sup>25</sup>. Erythroid development is a salient example, where due to MURK, the directions of RNA

284 velocity were even reversed from the expected trajectory<sup>19, 24</sup> (**Fig. 4h**). Remarkably, using  
285 the phylogenetic velocities of MEGs in erythroid development from a single embryo (embryo  
286 3) of Chan *et al.* dataset<sup>32</sup> (**Supplementary Fig. 15**), PhyloVelo accurately predicted the  
287 expected erythroid trajectory with the scRNA-seq data of temporal mouse embryos (E6.5-  
288 8.5) from Pijuan-Sala *et al.*<sup>19</sup>, despite the lineage tree was not being available (**Fig. 4i**). The  
289 PhyloVelo pseudotime was also strongly correlated with mouse embryo time (**Fig. 4j-k**).  
290 Together, these data demonstrated that PhyloVelo can circumvent the MURK issue of RNA  
291 velocity and the MEGs identified from one dataset can be also applied to independent  
292 datasets, even when phylogenetic information is not available.

### 293 ***PhyloVelo identifies lung tumor dedifferentiation***

294 We next applied PhyloVelo to a CRISPR/Cas9-based lineage tracing dataset in a  
295 genetically-engineered mouse model (GEMM) of lung adenocarcinoma (*Kras*<sup>LSL-G12D/+</sup>;  
296 *Trp53*<sup>fl/fl</sup>, or KP model), described by Yang *et al.*<sup>51</sup>. Cancer GEMMs allow one to study  
297 tumor evolutionary trajectory in its native microenvironment. Two primary tumors from KP  
298 mice (3726\_NT\_T1 and 3435\_NT\_T1) were selected because of their relatively high  
299 resolution of the lineage trees and composition of diverse cell types (including AT2-like,  
300 AT1-like, Gastric-like, High plasticity, Lung-mixed, Endoderm-like, Early EMT (epithelial-  
301 mesenchymal transition)-1, etc.) (**Fig. 5a, Extended Data Fig. 6a**). In total, 337 and 344  
302 MEGs ( $q < 0.05$ ) were identified from these two tumors, respectively (**Supplementary Fig.**  
303 **16, Supplementary Table 1**). RNA velocity by scVelo performed reasonably well in  
304 3435\_NT\_T1 to recapitulate the expected trajectory from AT2-like to High plasticity, and to  
305 Lung-mixed cells (**Extended Data Fig. 6b**), whereas no clear trajectory was inferred in  
306 3726\_NT\_T1 by scVelo (**Fig. 5b**). In contrast, in both tumors PhyloVelo identified AT2-like  
307 cells as the cell-of-origin of KP lung adenocarcinoma and also recovered the trajectory from  
308 AT2-like to lung-mixed or Early EMT (**Fig. 5c, Extended Data Fig. 6c**). In 3726\_NT\_T1, two  
309 trajectories appeared to coexist, namely 1) AT2-like > lung-mixed > Early EMT and 2) AT2-

310 like > Endoderm-like > Early EMT (**Fig. 5c**), thus recapitulating the findings in the original  
311 study<sup>51</sup>. As previous reports<sup>52, 53</sup>, lung tumor development was accompanied by the loss of  
312 AT2 identity and gain of highly plastic phenotypes such as lung-mixed and EMT.

313 Yang *et al.*<sup>51</sup> defined a single-cell fitness signature (**Fig. 5d, Extended Data Fig. 6d**) where  
314 the expression of a specific gene module is associated with the cell proliferating fitness  
315 estimated from the phylogenetic tree. While no overt association between scVelo latent time  
316 and the fitness signatures was found (**Fig. 5d, Extended Data Fig. 6e**), the PhyloVelo  
317 pseudotime showed a strong correlation with the fitness signature in both tumors  
318 (3726\_NT\_T1, Spearman's  $\rho=0.86$ ,  $P=1.2\times 10^{-218}$ ; 3435\_NT\_T1, Spearman's  $\rho=0.83$ ,  
319  $P=4.0\times 10^{-280}$ , **Fig. 5d, Extended Data Fig. 6f**). This indicates an intrinsic link between our  
320 measure of phylogenetic velocity and the cell fitness. Interestingly, CytoTRACE<sup>10</sup>, a  
321 computational algorithm to predict the cellular differentiation states with scRNA-seq data,  
322 revealed a drastic increase of the expressed gene number during the tumor evolution (**Fig.**  
323 **5e-f, Extended Data Fig. 6g-i**), which was in line with a dedifferentiation model. Reanalysis  
324 of a scRNA-seq dataset from human non-small cell lung cancers<sup>54</sup> verified that CytoTRACE  
325 scores increased as tumor evolved from normal lung tissue, early-stage cancer, advanced-  
326 stage cancer to pleural fluids and lymph node metastasis (**Fig. 5g**). This suggested that  
327 dedifferentiation might be a general phenomenon during tumor evolution. Also, this  
328 indicated that although gene expression diversity is a key feature of the developmental  
329 potential<sup>10</sup>, the directions of cell-state transitions are highly context-dependent and the  
330 cellular trajectories in normal differentiation and disease progression can be completely  
331 reversed.

332 Again, we showed that the phylogenetic velocity of MEGs identified from only one pilot  
333 tumor 3726\_NT\_T1 (754 cells, **Supplementary Fig. 16**) enabled the robust inference of  
334 velocity fields for other independent KP tumors even the lineage trees were not utilized  
335 (n=58,022 cell, **Fig. 5h-i**). In fact, the inferred and expected trajectories were highly

336 consistent for the cell types that existed in 3726\_NT\_T1 but not for the cell types such as  
337 Mesenchymal 1 and 2 (**Fig. 5h-i**). Interestingly, the phenomenon that two trajectories  
338 coexist as in 3726\_NT\_T1 (**Fig. 5c**) was more evident on the pooled PhyloVelo velocity  
339 fields (**Fig. 5h**) and the quantitative state-transition map computed by Dynamo based on the  
340 PhyloVelo velocity fields (**Fig. 5j-k**). These results demonstrated the generality of  
341 transcriptomic clock across the KP tumors, but also implied that a large single-cell lineage  
342 tree spanning numerous cell types must be reconstructed in order to identify more  
343 ubiquitously clock-like MEGs.

#### 344 ***PhyloVelo for clonal lineage tracing data using static barcodes***

345 Clonal lineage tracing by static barcoding has been paired with single-cell transcriptomics,  
346 such as LARRY<sup>37</sup>, CellTagging<sup>55</sup> and immune repertoire profiling<sup>56</sup> (e.g. scVDJ-seq), which  
347 provides both clonality and gene expression profiles of individual cells. Clonal lineage  
348 tracing identifies cells of common ancestry but can't resolve phylogenetic relationship within  
349 each clonal subpopulation. Nevertheless, similar to mutation number, the clone sizes  
350 (number of cells sharing a unique static barcode) also indicate the relative proliferative  
351 activity of the cells in the past division history. Hence, we considered using clone size as a  
352 surrogate of phylogenetic time in PhyloVelo (**Fig. 6a**). According to a simple exponential  
353 growth model ( $c_t = c_0 e^{rt}$ ), the logarithm of clone size ( $\log(c)$ ) has a linear relationship with  
354 cell proliferation rate ( $r$ ). Therefore, here "MEGs" can be identified by the significant  
355 association of latent expressions  $\mathbf{Z} = (\hat{z}_1, \hat{z}_2, \dots, \hat{z}_n)$  with the logarithm of clone size  $\mathbf{B} =$   
356  $(\log(c_1), \log(c_2), \dots, \log(c_n))$  (**Fig. 6b**). The velocities were estimated the same way as  
357 using phylogeny-resolved scRNA-seq data (e.g. CRISPR lineage tracing) (**Methods**).

358

359 We first applied this extended model of PhyloVelo to a lentiviral barcoding dataset from *in*  
360 *vitro* hematopoiesis<sup>37</sup>. This dataset sampled hematopoietic differentiation over the culture of  
361 2, 4, and 6 days and contained 29,242 cells where each could be traced by one unique



362 barcode (**Fig. 6c**). In total, 419 MEGs ( $q < 0.05$ ) were identified, with 297 positively and 122  
363 negatively associated with the logarithm of clone sizes at day 6 (**Supplementary Fig. 17a**,  
364 **Supplementary Table 1**). PhyloVelo velocity fields accurately traced differentiated cells  
365 (erythrocytes, megakaryocytes, mast cells, neutrophils, monocytes, *etc*) backward to  
366 undifferentiated progenitor cells (**Fig. 6d-e**). In fact, PhyloVelo pseudotime was strongly  
367 correlated with the clonal fate potency inferred by the Cospar algorithm<sup>36</sup> (**Supplementary**  
368 **Fig. 18**), indicating clone size-based PhyloVelo has successfully recovered the  
369 hematopoietic differentiation trajectories.

370

371 We also showcased the application of PhyloVelo to immune repertoire profiling data, where  
372 simultaneous lineage receptor sequences and gene expression profiles of individual T cells  
373 are available. With a lineage tracing dataset of intratumoral CD8+ T cells in basal cell  
374 carcinoma<sup>57</sup> (**Supplementary Fig. 17b**), PhyloVelo combined with Dynamo<sup>50</sup> quantified the  
375 T cell state-transition rates pre and post PD-1 blockade treatment (**Fig. 6f-k**). The  
376 quantitative transition map revealed that the enriched CD8+ activated T cells post treatment  
377 had few origin (2.4%) from the infiltrated naïve or memory CD8+ T cells (**Fig. 6j**), in line with  
378 the clonal replacement model of T lymphocytes after PD-1 treatment<sup>57</sup>. Interestingly, the  
379 hybrid activated/exhausted CD8+ T cells appeared to be mainly (81%) derived from  
380 exhausted CD8+ T cells before PD-1 treatment (**Fig. 6h-i**), while they were instead almost  
381 all (99%) from activated T cells after treatment (**Fig. 6j-k**). Therefore, our quantitative  
382 analyses revealed the drastic fate plasticity of intratumoral CD8+ T cells during checkpoint  
383 blockade immunotherapy.

#### 384 ***Comparison of PhyloVelo with different RNA velocity methods***

385 We noticed several methods for estimating RNA velocity have been developed, which  
386 model cell-specific and/or gene-specific RNA kinetics with deep learning framework such as  
387 VeloVAE<sup>58</sup>, DeepVelo<sup>59</sup> and cellDancer<sup>60</sup>, or by the radial basis function such as

388 UniTVelo<sup>61</sup>. These methods highlighted their improved performance as compared to  
389 scVelo. We therefore sought to compare PhyloVelo with these RNA velocity estimators  
390 using the scRNA-seq datasets of *C. elegans*, mouse erythroid cells and KP mouse lung  
391 tumor in this study. For *C. elegans*, only cellDancer seemed to improve the RNA velocity  
392 estimates relative to scVelo, where all others still gave reversed directions against the *C.*  
393 *elegans* embryo time (**Supplementary Fig. 19**). For mouse erythroid development (E8.5),  
394 UniTVelo showed the best performance amongst the five RNA velocity methods with  
395 competitively accurate estimations as PhyloVelo (**Extended Data Fig. 7**). Finally, for KP  
396 mouse lung tumor (3726\_NT\_T1), DeepVelo and UnitVelo performed reasonably well to  
397 recapitulate the expected trajectory from AT2-like to Lung-mixed cells. PhyloVelo  
398 pseudotime still showed the best correlation with cell fitness signatures (**Supplementary**  
399 **Fig. 20**). Overall, these preliminary comparison analyses highlighted the superior  
400 performance of PhyloVelo relative to RNA velocity methods.

#### 401 ***MEGs are enriched in ribosome-mediated processes***

402 To systematically investigate the potential functions of MEGs across tissues and organisms,  
403 we analyzed three additional CRISPR-based lineage tracing datasets that were derived  
404 from mouse or human cell lines, including pancreatic cancer KPCY<sup>62</sup>, lung cancer A549<sup>63</sup>  
405 and normal epithelial cells HEK293T<sup>64</sup>. The KPCY and A549 cells were sampled from *in*  
406 *vivo* mouse xenograft model, while the HEK293T cells were from a single-cell derived clone  
407 of *in vitro* culture. Interestingly, although these cell lines are known to be non-differentiating,  
408 continuous cell-state transitions were evident according to the PhyloVelo velocity fields  
409 (**Extended Data Fig. 8, Supplementary Figs. 21-22**). For instance, PhyloVelo recovered a  
410 dynamic EMT trajectory during the metastatic progression of KPCY and A549 cells in  
411 mouse xenografts (**Extended Data Fig. 8, Supplementary Figs. 21**). Even for *in vitro*  
412 culture of HEK293T cells, PhyloVelo and scVelo consistently showed continuous state  
413 transitions (**Supplementary Fig. 22**), and this phenomenon was not caused by cell-cycle

414 heterogeneity (**Extended Data Fig. 9**). Interestingly, the CytoTRACE “stemness” scores  
415 were strongly associated with PhyloVelo pseudotime in A549 mouse xenografts  
416 (**Supplementary Fig. 23**), also in line with a dedifferentiation process during the *in vivo*  
417 tumor progression.

418 We found the MEGs identified across organisms (mouse and human) and tissue or cell  
419 types (embryo, tumor tissues, cell lines and intratumoral T cells) were significantly  
420 overlapped (**Extended Data Fig. 10a-b**). Interestingly, the ribosome machinery was  
421 strongly enriched across the tissues and organisms, including translation, ribonucleoprotein  
422 complex biogenesis, ribosome biogenesis and assembly (**Fig. 6I**). For instance, the gene  
423 scores of ribosomal protein (RP) genes were significantly associated with the phylogenetic  
424 time based on tree distance in KP lung tumors or based on clone size in *in vitro*  
425 hematopoiesis and Intratumoral CD8+ T cells (**Supplementary Fig. 24**). To rule out the  
426 possibility that ribosomal genes were identified because of their high expression  
427 heterogeneity amongst the cells, we further performed permutation analysis where the  
428 phylogenetic distances were randomly shuffled and assigned to the cells. Here, although  
429 some “pseudo-MEGs” can still be identified (**Extended Data Fig. 10c**), they only showed  
430 weak associations (most  $q$  values were around 0.05) with the phylogenetic distances.  
431 Importantly, no significant enrichment in ribosome-mediated processes was found  
432 (**Extended Data Fig. 10d**). These results strongly suggest that many ribosomal genes  
433 genuinely follow the clock-like expression dynamics during cell proliferation and  
434 differentiation.

## 435 **Discussion**

436 Defining the correct directions of cell-fate transitions is crucial for unraveling the (epi)genetic  
437 regulators that drivers of lineage specification in diverse biological contexts<sup>20</sup>. Although RNA  
438 velocity and its improvements<sup>22, 23, 50, 61</sup> are powerful approaches for quantifying cellular

439 transitions from single-cell transcriptomic data, an accurate estimation of the velocity fields  
440 is still challenging because of the highly dynamic RNA kinetics (transcription, splicing and  
441 degradation)<sup>23-25</sup> and the biased capture of intron regions by droplet-based scRNA-seq<sup>65</sup>.  
442 The fundamental objective of our PhyloVelo algorithm is the same as RNA velocity – to  
443 extrapolate the gene expression of single cells to their near future or past states. However,  
444 unlike RNA velocity, PhyloVelo quantifies the transcriptomic velocity by measuring the rate  
445 of expression changes along a cell division history. Using various single-cell datasets where  
446 the coupled lineage information was available, PhyloVelo not only recovered the expected  
447 trajectories more accurately, but also gave more consistent estimates of the velocities  
448 relative to RNA velocity, across diverse biological contexts.

449

450 Analysis of lineage-resolved scRNA-seq datasets with PhyloVelo across mouse embryo  
451 development, hematopoietic differentiation, tumor evolution and immune cell dynamics  
452 yields insights into cell state dynamics. First, in each of the lineage tracing datasets, we  
453 have identified 100-500 MEGs, suggesting a considerable number of genes follow  
454 directional expression trajectories along cell divisions, at least within the phylogenetic time  
455 range of sampled cells. Interestingly, the MEGs across tissues and organisms had highly  
456 similar functions in translation and ribosome biogenesis, in line with their crucial role in  
457 regulating cell proliferation. Previous studies have also shown that ribosomal protein (RP)  
458 genes are commonly downregulated during differentiation<sup>66, 67</sup>, which represent robust  
459 markers of differentiation potency<sup>68, 69</sup>. Our study provides an explanation on why they serve  
460 as markers of differentiation potency - that is probably through regulating cell cycle and  
461 proliferation. In other words, the downregulation of some RP genes might suppress cell  
462 proliferation and thus promote differentiation. Second, we showcased that the phylogenetic  
463 velocities of MEGs estimated from one lineage-resolved scRNA-seq dataset can be reused  
464 in independent scRNA-seq datasets in similar biological conditions, in the absence of  
465 lineage information. Because obtaining a coupled lineage tree for every scRNA-seq dataset

466 is rather laborious, the transferability of MEGs facilitates the application of PhyloVelo to  
467 conventional scRNA-seq datasets. It is important to note that the transfer of MEGs is limited  
468 to similar biological conditions. However, it is not recommended to transfer MEGs between  
469 different conditions, such as normal and disease as the MEGs can differ significantly. To  
470 study normal development, it would be beneficial for future efforts to compile a  
471 comprehensive set of MEGs encompassing whole-organism development by employing  
472 whole-organism lineage tracing. These MEGs can then be used to assess their  
473 phylogenetic velocities in various organs or tissues. On the other hand, in the context of  
474 diseased conditions, it is essential to identify specific MEGs for each dataset, utilizing the  
475 corresponding lineage tracing data. Third, by combining PhyloVelo and Dynamo<sup>50</sup>, we were  
476 able to estimate the transition probability between any two cell states. For instance, in  
477 mouse lung tumors, we found two competing trajectories of cell-state evolution through  
478 dedifferentiation. In another case of intratumor CD8<sup>+</sup> T cells, we found distinct origin of  
479 activated T cells pre and post anti-PD-1 treatment. These quantitative analyses revealed  
480 high cell plasticity for both tumor cells and immune microenvironment during tumor  
481 progression and treatment. Importantly, as RNA velocity, PhyloVelo velocities fields are  
482 useful for identifying cell-fate drivers<sup>20</sup> or core gene regulatory networks<sup>70</sup>.

483

484 Despite the rapid development of CRISPR lineage tracing methods<sup>4, 29, 71</sup>, building high-  
485 precision lineage trees with single-cell resolution is still challenging because of the small  
486 number of Cas9 target sites (typically <50), rapid saturation, frequent inter-site deletions,  
487 and other factors<sup>47, 72, 73</sup>. This requires a more reliable lineage tracing method that has a  
488 larger lineage-labeling space and more stable mutagenesis strategy. We recently developed  
489 a base editor-based lineage tracing method, called SMALT<sup>42</sup>, which leverages a genetically-  
490 evolved activation-induced cytidine deaminase (AID) to specifically target a 3k synthetic  
491 DNA barcode and induce C to T mutations on it with high efficacy. The lineage tree  
492 reconstructed by SMALT achieved nearly single-cell resolution and over 80% statistically

493 bootstrapping support<sup>42</sup>. We envision the combination of SMALT lineage tracing and single-  
494 cell transcriptomics will greatly empower PhyloVelo to resolve complex lineage dynamics in  
495 more diverse biological contexts, such as genetic perturbation or disease progression.

496

497 In summary, we provide a theoretical framework to quantify cell-fate transitions by  
498 leveraging both single-cell lineage and transcriptomic information. With the rapid  
499 development of single-cell lineage tracing technologies and emergence of lineage-traced  
500 multi-omic data<sup>74</sup>, we envision our method will facilitate the lineage analysis for complex  
501 cellular processes and the discovery of the cell-fate determinants in diverse organisms,  
502 tissues, and diseases.

## 503 **Acknowledgments**

504 We thank Yuanhua Huang, Jiguang Wang, Jin Xu, Liang Ma, Wanze Chen and members of  
505 Hu laboratory for constructive discussions. This work was supported by National Key R&D  
506 Program of China (2021YFA1302500 to Z.H.), National Natural Science Foundation of  
507 China (11971405 to D.Z., 32270693 to Z.H.), Guangdong Basic and Applied Basic  
508 Research Foundation (2021B1515020042 to Z.H.), Fundamental Research Funds for the  
509 Central Universities (20720230023 to D.Z.) and China Postdoctoral Science Foundation  
510 (2021M693303 to Z.L, 2022M723301 to X.W.).

## 511 **Author contributions**

512 Z.H. and K.W. conceived the concept of phylogenetic velocity. Z.H., K.W., and D.Z.  
513 designed the study. K.W. developed the mathematical framework and implemented the  
514 software. K.W., Z.H., L.H., Z.L., X.W., X.Z., Z.Y. analyzed the data. W.Z. and Z.Z. provided  
515 constructive suggestions on the model. K.W., Z.H., D.Z., C.C., X.H., interpreted results. Z.H.  
516 and K.W. wrote the manuscript with contributions from all co-authors. Z.H., and D.Z.  
517 supervised the study.

## 518 **Competing interests**

519 C.C. is an advisor and stockholder in Grail, Ravel, and DeepCell and an advisor to  
520 Genentech, Bristol Myers Squibb, 3T Biosciences, and NanoString. Other authors declare  
521 no competing interests.

## 522 **Figure Legends**

523 **Fig. 1. Schematic of the PhyloVelo framework.** (a) Schematic of monotonically expressed  
524 genes (MEGs) over phylogenetic time on a cell phylogenetic tree. (b) Two examples of  
525 MEGs whose latent expressions are associated with the phylogenetic time (cell divisions or  
526 mutation number). A diffusion process of gene expressions was used to model the changes  
527 of latent expressions over phylogenetic time. This enables the estimation of the  
528 phylogenetic velocity,  $v = (v_1, v_2, \dots, v_G)$ , which corresponds to the drift coefficients of  $G$   
529 MEGs in the diffusion process (approximate to the slope of linear regression between latent  
530 expression and phylogenetic time). Whiskers: minimum and maximum. (c) Phylogenetic  
531 velocity predicts the past transcriptional state of a cell before a unit of phylogenetic time  
532 (one cell division or mutation). (d) Projection of the phylogenetic velocity into low  
533 dimensional embedding enables the mapping of cell-state trajectory in backward directions.

534  
535 **Fig. 2. PhyloVelo recovers complex cell lineages in simulations.** Simulation of single-  
536 cell RNA-seq data and paired cell-division history under linear (a), bifurcated (b), and  
537 convergent (c) differentiation models, respectively. Colors are labeled by cell types. Each  
538 simulation consists of 1,000 cells randomly sampled from a growing cell population at  
539 10,000 cells. Each cell has 2,000 expressed genes, including 200-300 MEGs. (d-f)  
540 Phylogenetic velocity fields reconstructed by PhyloVelo for the corresponding differentiation  
541 scenarios. The left panel shows the single-cell level of velocity fields, while the right panel

542 shows the same velocity fields visualized as streamlines in scVelo. PhyloVelo velocity fields  
543 are at backward directions.

544

545 **Fig. 3. PhyloVelo reconstructs the embryonic differentiation trajectories of *C.***

546 ***elegans*.** (a) Phylogenetic tree of the *C. elegans* AB lineage. (b) Heatmap showing the  
547 expressions (z-score normalized) of MEGs along *C. elegans* embryo time. (c) The ground-  
548 truth velocity fields represent vectors superimposed on the cells that point to their immediate  
549 parental cells on the Uniform Manifold Approximation and Projection (UMAP) plot. (d-e) The  
550 velocity fields estimated by scVelo (dynamical mode) (d) or PhyloVelo (e). Dash square  
551 indicates the early embryonic lineages where RNA velocity gave erroneous estimations on  
552 the fate directions. (f) *C. elegans* embryo time as Packer *et al.*<sup>14</sup>. (g) scVelo latent time. (h)  
553 PhyloVelo pseudotime. (i) RNA velocity fields for all 29,600 AB lineage cells. Colors are  
554 labeled by scVelo latent time. (j) The correlation between scVelo latent time and embryo  
555 time for all AB lineage cells. (k) PhyloVelo velocity fields for all 29,600 AB lineage cells,  
556 estimated by the phylogenetic velocity of MEGs in a single embryo (n=298 cells). Cell colors  
557 are labelled by PhyloVelo pseudotime. (l) The correlation between PhyloVelo pseudotime  
558 and embryo time for all AB lineage cells. The Spearman correlation coefficients and *P*  
559 values are shown.

560

561 **Fig. 4. PhyloVelo reconstructs the cellular trajectory of mouse erythroid maturation.**

562 (a) Phylogenetic tree of the 2,419 erythroid lineage cells (embryo 3, E8.5) in Chan *et al.*  
563 dataset<sup>32</sup>. (b-c) RNA velocity fields (scVelo - dynamical mode) and the latent time of mouse  
564 erythroid development. (d) Muller plot showing the fractions of four cell types that change  
565 over scVelo latent time. (e-f) PhyloVelo velocity fields and the pseudotime of mouse  
566 erythroid development. (g) Muller plot showing the fractions of four cell types that change  
567 over PhyloVelo pseudotime. (h) Erroneous estimations of RNA velocity fields on erythroid  
568 maturation because of multiple rate kinetics (MURK). Data were from Pijuan-Sala *et al.*<sup>19</sup>. (i)



569 PhyloVelo velocity fields of erythroid maturation for Pijuan-Sala *et al.* dataset while using the  
570 MEGs identified from Chan *et al.* dataset. (j) PhyloVelo pseudotime of erythroid maturation  
571 in Pijuan-Sala *et al.* dataset. (k) The correlation between PhyloVelo pseudotime and mouse  
572 embryo time (n=12,324 cells). The Spearman correlation coefficient and *P* value are shown  
573 here. Whiskers: minimum and maximum; center lines: median.

574

575 **Fig. 5. PhyloVelo identifies a dedifferentiation trajectory in lung tumor evolution.** (a)  
576 Phylogenetic tree of 754 cells from a KP-mouse primary lung tumor, 3726\_NT\_T1, in Yang  
577 *et al.* dataset<sup>51</sup>. The scRNA-seq data, cell type annotations, and lineage trees were  
578 obtained from the original study. (b) RNA velocity fields (scVelo - dynamical mode). (c)  
579 PhyloVelo velocity fields. (d) Fitness signatures of individual cells, as defined by Yang *et al.*  
580 (e) CytoTRACE score of individual cells. (f) The correlation between PhyloVelo pseudotime  
581 and CytoTRACE scores. The Spearman correlation coefficient and *P* value are shown here.  
582 (g) CytoTRACE score of single tumor cells from human lung primary sites (tLung and tL/B),  
583 pleural fluids (PE), lymph node metastases (mLN), and brain metastases (mBrain), as well  
584 as normal tissues from lungs (nLung), as described in Kim *et al.*<sup>54</sup>. Bar, median; box, 25th to  
585 75th percentile (IQR); vertical line, data within 1.5 times the IQR. (h) PhyloVelo velocity  
586 fields for all 58,022 single cells from pooled KP primary lung tumors, estimated by the  
587 MEGs identified from 3726\_NT\_T1. (i) PhyloVelo velocity fields for the cell types that  
588 existed in 3726\_NT\_T1. (j) Cell-type transition graph (backward) based on the transition  
589 rate matrix between any two cell types (k), estimated by Dynamo using PhyloVelo velocity  
590 fields as input. The arrows point from the current states to the past states.

591

592 **Fig. 6. PhyloVelo inference with clonal lineage tracing data and MEGs are enriched in**  
593 **ribosome-mediated processes.** (a) Schematic of clonal lineage tracing data where static  
594 barcodes identify cells of common ancestry. Clone size, denoted by  $c_k$  for  $k$  clones,  
595 represents the number of cells carrying the same unique barcode. (b) Two examples of

596 clonal size-based MEGs whose latent expressions are positively or negatively associated  
597 with the logarithm of clone sizes, respectively. Whiskers: minimum and maximum. (c)  
598 scRNA-seq data of in vitro hematopoietic differentiation from Weinreb *et al.*<sup>37</sup>, where each  
599 cell over the course of 2, 4, and 6 days culture could be traced by one unique barcode. (d)  
600 The velocity fields estimated by PhyloVelo. (e) Cell type transition graph (backward) of in  
601 vitro hematopoietic differentiation. (f) UMAP of tumor-infiltrating CD8+ T cells in BCC  
602 samples pre- and post-PD-1 blockade, colored by anti-PD-1 treatment status. Data were  
603 from Yost *et al.*<sup>57</sup> (g) The velocity fields estimated by PhyloVelo. (h-i) Cell-type transition  
604 graph and transition matrix (backward) at pre-treatment. (j-k) Cell-type transition graph and  
605 transition matrix (backward) at post-treatment. CD8\_act: CD8+ activated T cells; CD8\_ex:  
606 CD8+ exhausted T cells; CD8\_ex\_act: CD8+ exhausted/activated T cells; CD8\_eff: CD8+  
607 effector T cells; CD8\_mem: CD8+ memory T cells. (l) Gene ontology (GO) enrichment of  
608 MEGs identified across tissues and organisms. The top and most commonly shared 20  
609 biological processes are shown. Ribosome-mediated processes are highlighted.

610

611 **Extended Data Fig. 1. Quantitative metrics for evaluating PhyloVelo's performance on**  
612 **simulation data.** Two quantitative metrics with varied cell numbers (a), non-linear MEGs  
613 (b), different dimensionality reduction methods (c), varied data sparsity (d) and varied  
614 numbers of MEGs (e). All benchmarks are simulated 50 times independently. Bar, median;  
615 box, 25th to 75th percentile (IQR); vertical line, data within 1.5 times the IQR.

616

617 **Extended Data Fig. 2. PhyloVelo velocity fields in three additional lineages of *C.***  
618 ***elegans*.** (a-c) Hypodermis, body wall muscle (BWM) and Pharynx lineage cells,  
619 respectively. Colors are labeled by the estimated embryo time (minutes). (d-f) PhyloVelo  
620 velocity fields of the three lineages respectively each consisting of 2,000 randomly sampled  
621 cells from multiple embryos, which were reconstructed using the MEGs identified from 298  
622 AB lineage cells. Colors are labeled by the PhyloVelo pseudotime. (g-i) The correlation

623 between PhyloVelo pseudotime and embryo time for the cells in the three lineages. The  
624 Spearman correlation coefficients and  $P$  values are shown.

625

626 **Extended Data Fig. 3. High concordance of MEGs identified from 4 mouse embryos**  
627 **(E8.0/8.5) in Chan *et al.*<sup>32</sup>** (a) Venn diagram showing the overlap of MEGs identified from  
628 four mouse embryos in the dataset of Chan *et al.*  $P$  value, one-sided SuperExactTest multi-  
629 set intersection test. (b-g) The correlation of phylogenetic velocities  $v$  for the overlapped  
630 MEGs between any two embryos. The Pearson correlation coefficients and  $P$  values are  
631 shown.

632

633 **Extended Data Fig. 4. The global differentiation trajectories of whole mouse embryos**  
634 **and brain tissues predicted by LT-MEGs.** (a) PhyloVelo velocity fields of mouse embryos  
635 (E6.5-8.5) mapped by 104 LT-MEGs with the temporal scRNA-seq dataset from Pijuan-Sala  
636 *et al.*<sup>19</sup> (b-c) UMAP plot colored by PhyloVelo pseudotime (b) or sample capture time (c).  
637 (d) PhyloVelo velocity fields of mouse brain (E7-18) mapped by LT-MEGs with the temporal  
638 scRNA-seq dataset from La Manno *et al.*<sup>48</sup> (e-f) tSNE plot colored by PhyloVelo pseudotime  
639 (e) or sample capture time (f). UMAP or tSNE coordinates were as the original studies.

640

641 **Extended Data Fig. 5. PhyloVelo velocity fields and quantitative state transitions of**  
642 **mouse erythroid development for four embryos from Chan *et al.*<sup>32</sup>** (a-d) PhyloVelo  
643 velocity fields. (e-h) The transition rate (backward) between any two cell types. (i-l) cell-type  
644 transition graph (backward) visualized based on the cell-type transition rates. PhyloVelo  
645 velocity fields were used as the input of Dynamo.

646

647 **Extended Data Fig. 6. PhyloVelo reconstructs the cellular trajectory of lung cancer**  
648 **evolution in 3435\_NT\_T1.** (a) Single-cell phylogenetic tree of primary lung tumor  
649 3435\_NT\_T1 (n=1,109 cells) from KP (Kras<sup>LSL-G12D/+</sup>;Trp53<sup>fl/fl</sup>) mouse model. The single-cell

650 RNA data, cell type annotations and lineage tree were obtained from Yang *et al.*<sup>51</sup> (b) RNA  
651 velocity fields (scVelo - dynamical mode). (c) PhyloVelo velocity fields. (d) The fitness  
652 signatures of single cells as defined by Yang *et al.* (e) The correlation between scVelo latent  
653 time and fitness signatures. (f) The correlation between PhyloVelo pseudotime and fitness  
654 signatures. (g) CytoTRACE score of individual cells. (h) The correlation between scVelo  
655 latent time and CytoTRACE scores. (i) The correlation between PhyloVelo pseudotime and  
656 CytoTRACE scores. (j) The correlation of phylogenetic velocities for the overlapped MEGs  
657 between KP primary tumor 3435\_NT\_T1 and 3726\_NT\_T1. The Pearson correlation  
658 coefficients and *P* values are shown here.

659

660 **Extended Data Fig. 7. Comparison of PhyloVelo with scVelo, VeloVAE, DeepVelo,**  
661 **CellDancer and UniTVelo respectively on mouse erythroid data.** scVelo - RNA velocity  
662 fields (a), latent time (b) and the fractions of different cell types along latent time (c).  
663 VeloVAE - RNA velocity fields (d), latent time (e) and the fractions of different cell types  
664 along latent time (f). DeepVelo - RNA velocity fields (g), latent time (h) and the fractions of  
665 different cell types along latent time (i). cellDancer - RNA velocity fields (j), pseudotime (k)  
666 and the fractions of different cell types along pseudotime (l). UniTVelo - RNA velocity fields  
667 (m), latent time (n) and the fractions of different cell types along latent time (o). PhyloVelo -  
668 velocity fields (p), pseudotime (q) and the fractions of different cell types along pseudotime  
669 (r). PhyloVelo velocity fields are in backward directions.

670

671 **Extended Data Fig. 8. The dynamic EMT trajectory in metastatic progression of**  
672 **pancreatic cancer KPCY cells.** (a) Phylogenetic tree of 601 non-repetitive terminal cells in  
673 tumor subclone M1.1 from Simeonov *et al.*<sup>62</sup> Cell colors are labeled by EMT pseudotime as  
674 defined in the original study. (b) The total UMI count (normalized) of MEGs changing with  
675 the phylogenetic distance from the root. (c) Heatmap of MEG expressions (z-score  
676 normalized) with EMT pseudotime. (d) RNA velocity fields (scVelo - dynamical mode). Cell

677 colors are labeled by EMT pseudotime. (e) scVelo latent time. (f) The correlation between  
678 scVelo latent time and EMT pseudotime. (g) PhyloVelo velocity fields. Cell colors are  
679 labeled by EMT pseudotime. (h) PhyloVelo pseudotime. (i) The correlation between  
680 PhyloVelo pseudotime and EMT pseudotime. The Spearman correlation coefficients and  $P$   
681 values are shown here.

682

683 **Extended Data Fig. 9. Continuous state transitions inferred by PhyloVelo after**

684 **regressing out cell-cycle effect. (a-d)** PhyloVelo velocity fields after regressing out cell-

685 cycle dynamics in KPCY, A549 Ig1, A549 Ig2 and HEK293T, respectively. (e-h) The

686 correlation of PhyloVelo pseudotime between original analysis and post regressing out of

687 cell-cycle effect in KPCY, A549 Ig1, A549 Ig2 and HEK293T, respectively. The Pearson

688 correlation coefficients and  $P$  values are shown here.

689

690 **Extended Data Fig. 10. Overlap of MEGs across organisms and tissue/cell types and**

691 **the permutation analysis of MEG identification. (a)** The overlap of MEGs identified in

692 different datasets as stratified by mouse vs human. (b) The overlap of MEGs identified in

693 different datasets as stratified by normal vs tumor cells.  $P$  values are by one-sided

694 hypergeometric test. (c) The  $q$  values of MEGs in standard and permutation analysis.

695 Permutation analysis was done by randomly shuffling the phylogenetic distances of the

696 cells, followed by the PhyloVelo inference procedure. The number of detected MEGs in

697 standard and permutation analysis respectively are:  $n=1,724$  and  $n=941$  genes in Embryo

698 E8/E8.5;  $n=681$  and  $n=445$  genes in KP lung tumor;  $n=424$  and  $n=141$  genes in KPCY;

699  $n=629$  and  $n=50$  genes in A549;  $n=243$  and  $n=90$  genes in HEK293T;  $n=419$  and  $n=112$

700 genes in *in vitro* hematopoiesis;  $n=368$  and  $n=270$  genes in CD8+T cells. Bar, median; box,

701 25th to 75th percentile (IQR); vertical line, data within 1.5 times the IQR. (d) The GO

702 enrichment of pseudo-MEGs across the seven lineage tracing datasets.

703

## 704 References

- 705 1. Salipante, S.J. & Horwitz, M.S. Phylogenetic fate mapping. *Proc Natl Acad Sci U S A* **103**, 5448-  
706 5453 (2006).
- 707 2. Sulston, J.E., Schierenberg, E., White, J.G. & Thomson, J.N. The embryonic cell lineage of the  
708 nematode *Caenorhabditis elegans*. *Dev Biol* **100**, 64-119 (1983).
- 709 3. Stadler, T., Pybus, O.G. & Stumpf, M.P. Phylodynamics for cell biologists. *Science* **371**,  
710 eaah6266 (2021).
- 711 4. Baron, C.S. & van Oudenaarden, A. Unravelling cellular relationships during development and  
712 regeneration using genetic lineage tracing. *Nat Rev Mol Cell Biol* **20**, 753-765 (2019).
- 713 5. Moris, N., Pina, C. & Arias, A.M. Transition states and cell fate decisions in epigenetic  
714 landscapes. *Nat Rev Genet* **17**, 693-703 (2016).
- 715 6. Bendall, S.C. et al. Single-cell trajectory detection uncovers progression and regulatory  
716 coordination in human B cell development. *Cell* **157**, 714-725 (2014).
- 717 7. Trapnell, C. et al. The dynamics and regulators of cell fate decisions are revealed by  
718 pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386 (2014).
- 719 8. Haghverdi, L., Buttnner, M., Wolf, F.A., Buettner, F. & Theis, F.J. Diffusion pseudotime robustly  
720 reconstructs lineage branching. *Nat Methods* **13**, 845-848 (2016).
- 721 9. Saelens, W., Cannoodt, R., Todorov, H. & Saeys, Y. A comparison of single-cell trajectory  
722 inference methods. *Nat Biotechnol* **37**, 547-554 (2019).
- 723 10. Gulati, G.S. et al. Single-cell transcriptional diversity is a hallmark of developmental potential.  
724 *Science* **367**, 405-411 (2020).
- 725 11. Weinreb, C., Wolock, S., Tusi, B.K., Socolovsky, M. & Klein, A.M. Fundamental limits on  
726 dynamic inference from single-cell snapshots. *Proc Natl Acad Sci U S A* **115**, E2467-E2476  
727 (2018).
- 728 12. Tritschler, S. et al. Concepts and limitations for learning developmental trajectories from  
729 single cell genomics. *Development* **146** (2019).
- 730 13. Wagner, D.E. & Klein, A.M. Lineage tracing meets single-cell omics: opportunities and  
731 challenges. *Nat Rev Genet* **21**, 410-427 (2020).
- 732 14. Packer, J.S. et al. A lineage-resolved molecular atlas of *C. elegans* embryogenesis at single-cell  
733 resolution. *Science* **365** (2019).
- 734 15. Mulas, C., Chaigne, A., Smith, A. & Chalut, K.J. Cell state transitions: definitions and challenges.  
735 *Development* **148** (2021).
- 736 16. Gerber, T. et al. Single-cell analysis uncovers convergence of cell identities during axolotl limb  
737 regeneration. *Science* **362** (2018).

- 738 17. Liu, X. et al. Single-Cell RNA-Seq of the Developing Cardiac Outflow Tract Reveals Convergent  
739 Development of the Vascular Smooth Muscle Cells. *Cell Rep* **28**, 1346-1361 e1344 (2019).
- 740 18. Nowotschin, S. et al. The emergent landscape of the mouse gut endoderm at single-cell  
741 resolution. *Nature* **569**, 361-367 (2019).
- 742 19. Pijuan-Sala, B. et al. A single-cell molecular map of mouse gastrulation and early  
743 organogenesis. *Nature* **566**, 490-495 (2019).
- 744 20. Lange, M. et al. CellRank for directed single-cell fate mapping. *Nat Methods* **19**, 159-170  
745 (2022).
- 746 21. Gupta, P.B., Pastushenko, I., Skibinski, A., Blanpain, C. & Kuperwasser, C. Phenotypic Plasticity:  
747 Driver of Cancer Initiation, Progression, and Therapy Resistance. *Cell Stem Cell* **24**, 65-78  
748 (2019).
- 749 22. La Manno, G. et al. RNA velocity of single cells. *Nature* **560**, 494-498 (2018).
- 750 23. Bergen, V., Lange, M., Peidli, S., Wolf, F.A. & Theis, F.J. Generalizing RNA velocity to transient  
751 cell states through dynamical modeling. *Nat Biotechnol* **38**, 1408-1414 (2020).
- 752 24. Barile, M. et al. Coordinated changes in gene expression kinetics underlie both mouse and  
753 human erythroid maturation. *Genome Biol* **22**, 197 (2021).
- 754 25. Bergen, V., Soldatov, R.A., Kharchenko, P.V. & Theis, F.J. RNA velocity-current challenges and  
755 future perspectives. *Mol Syst Biol* **17**, e10282 (2021).
- 756 26. McKenna, A. et al. Whole-organism lineage tracing by combinatorial and cumulative genome  
757 editing. *Science* **353**, aaf7907 (2016).
- 758 27. Frieda, K.L. et al. Synthetic recording and in situ readout of lineage information in single cells.  
759 *Nature* **541**, 107-111 (2017).
- 760 28. Kalhor, R. et al. Developmental barcoding of whole mouse via homing CRISPR. *Science* **361**  
761 (2018).
- 762 29. VanHorn, S. & Morris, S.A. Next-Generation Lineage Tracing and Fate Mapping to Interrogate  
763 Development. *Dev Cell* **56**, 7-21 (2021).
- 764 30. Alemany, A., Florescu, M., Baron, C.S., Peterson-Maduro, J. & van Oudenaarden, A. Whole-  
765 organism clone tracing using single-cell sequencing. *Nature* **556**, 108-112 (2018).
- 766 31. Raj, B. et al. Simultaneous single-cell profiling of lineages and cell types in the vertebrate brain.  
767 *Nat Biotechnol* **36**, 442-450 (2018).
- 768 32. Chan, M.M. et al. Molecular recording of mammalian embryogenesis. *Nature* **570**, 77-82  
769 (2019).
- 770 33. Kester, L. & van Oudenaarden, A. Single-Cell Transcriptomics Meets Lineage Tracing. *Cell Stem*  
771 *Cell* **23**, 166-179 (2018).

- 772 34. Spanjaard, B. et al. Simultaneous lineage tracing and cell-type identification using CRISPR-  
773 Cas9-induced genetic scars. *Nat Biotechnol* **36**, 469-473 (2018).
- 774 35. Bowling, S. et al. An Engineered CRISPR-Cas9 Mouse Line for Simultaneous Readout of Lineage  
775 Histories and Gene Expression Profiles in Single Cells. *Cell* **181**, 1693-1694 (2020).
- 776 36. Wang, S.W., Herriges, M.J., Hurley, K., Kotton, D.N. & Klein, A.M. CoSpar identifies early cell  
777 fate biases from single-cell transcriptomic and lineage information. *Nat Biotechnol* **40**, 1066-  
778 1074 (2022).
- 779 37. Weinreb, C., Rodriguez-Fraticelli, A., Camargo, F.D. & Klein, A.M. Lineage tracing on  
780 transcriptional landscapes links state to fate during differentiation. *Science* **367** (2020).
- 781 38. Forrow, A. & Schiebinger, G. LineageOT is a unified framework for lineage tracing and  
782 trajectory inference. *Nat Commun* **12**, 4940 (2021).
- 783 39. Lopez, R., Regier, J., Cole, M.B., Jordan, M.I. & Yosef, N. Deep generative modeling for single-  
784 cell transcriptomics. *Nat Methods* **15**, 1053-1058 (2018).
- 785 40. Butler, M.A. & King, A.A. Phylogenetic Comparative Analysis: A Modeling Approach for  
786 Adaptive Evolution. *Am Nat* **164**, 683-695 (2004).
- 787 41. Papadopoulos, N., Gonzalo, P.R. & Soding, J. PROSSTT: probabilistic simulation of single-cell  
788 RNA-seq data for complex differentiation processes. *Bioinformatics* **35**, 3517-3519 (2019).
- 789 42. Liu, K. et al. Mapping single-cell-resolution cell phylogeny reveals cell population dynamics  
790 during organ development. *Nat Methods* **18**, 1506-1514 (2021).
- 791 43. Wagner, D.E. et al. Single-cell mapping of gene expression landscapes and lineage in the  
792 zebrafish embryo. *Science* **360**, 981-987 (2018).
- 793 44. Salipante, S.J., Kas, A., McMonagle, E. & Horwitz, M.S. Phylogenetic analysis of developmental  
794 and postnatal mouse cell lineages. *Evol Dev* **12**, 84-94 (2010).
- 795 45. Street, K. et al. Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics.  
796 *BMC Genomics* **19**, 477 (2018).
- 797 46. Wolf, F.A. et al. PAGA: graph abstraction reconciles clustering with trajectory inference  
798 through a topology preserving map of single cells. *Genome Biol* **20**, 59 (2019).
- 799 47. Salvador-Martinez, I., Grillo, M., Averof, M. & Telford, M.J. Is it possible to reconstruct an  
800 accurate cell lineage using CRISPR recorders? *Elife* **8** (2019).
- 801 48. La Manno, G. et al. Molecular architecture of the developing mouse brain. *Nature* **596**, 92-96  
802 (2021).
- 803 49. Baron, M.H., Isern, J. & Fraser, S.T. The embryonic origins of erythropoiesis in mammals. *Blood*  
804 **119**, 4828-4837 (2012).
- 805 50. Qiu, X. et al. Mapping transcriptomic vector fields of single cells. *Cell* **185**, 690-711 e645 (2022).



- 806 51. Yang, D. et al. Lineage tracing reveals the phylodynamics, plasticity, and paths of tumor  
807 evolution. *Cell* **185**, 1905-1923 e1925 (2022).
- 808 52. Marjanovic, N.D. et al. Emergence of a High-Plasticity Cell State during Lung Cancer Evolution.  
809 *Cancer Cell* **38**, 229-246 e213 (2020).
- 810 53. LaFave, L.M. et al. Epigenomic State Transitions Characterize Tumor Progression in Mouse  
811 Lung Adenocarcinoma. *Cancer Cell* **38**, 212-228 e213 (2020).
- 812 54. Kim, N. et al. Single-cell RNA sequencing demonstrates the molecular and cellular  
813 reprogramming of metastatic lung adenocarcinoma. *Nat Commun* **11**, 2285 (2020).
- 814 55. Bidy, B.A. et al. Single-cell mapping of lineage and identity in direct reprogramming. *Nature*  
815 **564**, 219-224 (2018).
- 816 56. Penter, L., Gohil, S.H. & Wu, C.J. Natural barcodes for longitudinal single cell tracking of  
817 leukemic and immune cell dynamics. *Frontiers in Immunology* **12**, 788891 (2022).
- 818 57. Yost, K.E. et al. Clonal replacement of tumor-specific T cells following PD-1 blockade. *Nat Med*  
819 **25**, 1251-1259 (2019).
- 820 58. Gu, Y., Blaauw, D. & Welch, J.D. Bayesian Inference of RNA Velocity from Multi-Lineage Single-  
821 Cell Data. *bioRxiv*, 2022.2007.2008.499381 (2022).
- 822 59. Cui, H., Maan, H., Taylor, M.D. & Wang, B. DeepVelo: Deep Learning extends RNA velocity to  
823 multi-lineage systems with cell-specific kinetics. *bioRxiv*, 2022.2004.2003.486877 (2022).
- 824 60. Li, S. et al. A relay velocity model infers cell-dependent RNA velocity. *Nat Biotechnol* (2023).
- 825 61. Gao, M., Qiao, C. & Huang, Y. UniTVelo: temporally unified RNA velocity reinforces single-cell  
826 trajectory inference. *Nat Commun* **13**, 6586 (2022).
- 827 62. Simeonov, K.P. et al. Single-cell lineage tracing of metastatic cancer reveals selection of hybrid  
828 EMT states. *Cancer Cell* **39**, 1150-1162 e1159 (2021).
- 829 63. Quinn, J.J. et al. Single-cell lineages reveal the rates, routes, and drivers of metastasis in cancer  
830 xenografts. *Science* **371** (2021).
- 831 64. Choi, J. et al. A time-resolved, multi-symbol molecular recorder via sequential genome editing.  
832 *Nature* (2022).
- 833 65. Qiu, Q. et al. Massively parallel and time-resolved RNA sequencing in single cells with scNT-  
834 seq. *Nat Methods* **17**, 991-1001 (2020).
- 835 66. Athanasiadis, E.I. et al. Single-cell RNA-sequencing uncovers transcriptional states and fate  
836 decisions in haematopoiesis. *Nat Commun* **8**, 2045 (2017).
- 837 67. Fei, L. et al. Systematic identification of cell-fate regulatory programs using a single-cell atlas  
838 of mouse development. *Nat Genet* **54**, 1051-1061 (2022).

- 839 68. Shi, J., Teschendorff, A.E., Chen, W., Chen, L. & Li, T. Quantifying Waddington's epigenetic  
840 landscape: a comparison of single-cell potency measures. *Brief Bioinform* (2018).
- 841 69. Teschendorff, A.E. & Feinberg, A.P. Statistical mechanics meets single-cell biology. *Nat Rev*  
842 *Genet* **22**, 459-476 (2021).
- 843 70. Singh, R., Wu, A.P., Mudide, A. & Berger, B. Unraveling causal gene regulation from the RNA  
844 velocity graph using Velorama. *bioRxiv*, 2022.2010.2018.512766 (2023).
- 845 71. Hughes, N.W. et al. Machine-learning-optimized Cas12a barcoding enables the recovery of  
846 single-cell lineages and transcriptional profiles. *Mol Cell* (2022).
- 847 72. Gong, W. et al. Benchmarked approaches for reconstruction of in vitro cell lineages and in  
848 silico models of *C. elegans* and *M. musculus* developmental trees. *Cell Syst* **12**, 810-826 e814  
849 (2021).
- 850 73. Espinosa-Medina, I., Garcia-Marques, J., Cepko, C. & Lee, T. High-throughput dense  
851 reconstruction of cell lineages. *Open Biol* **9**, 190229 (2019).
- 852 74. Jindal, K. et al. Multiomic single-cell lineage tracing to dissect fate-specific gene regulatory  
853 programs. *bioRxiv*, 2022.2010.2023.512790 (2022).
- 854 75. Gillespie, D.T. Exact stochastic simulation of coupled chemical reactions. *The journal of*  
855 *physical chemistry* **81**, 2340-2361 (1977).
- 856 76. Minh, B.Q. et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in  
857 the Genomic Era. *Mol Biol Evol* **37**, 1530-1534 (2020).
- 858 77. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R  
859 language. *Bioinformatics* **20**, 289-290 (2004).
- 860 78. Chen, W. et al. UMI-count modeling and differential expression analysis for single-cell RNA  
861 sequencing. *Genome Biol* **19**, 70 (2018).
- 862 79. Jia, C. Kinetic foundation of the zero-inflated negative binomial model for single-cell RNA  
863 sequencing data. *SIAM Journal on Applied Mathematics* **80**, 1336-1355 (2020).
- 864 80. Prim, R.C. Shortest connection networks and some generalizations. *The Bell System Technical*  
865 *Journal* **36**, 1389-1401 (1957).
- 866 81. Cock, P.J. et al. Biopython: freely available Python tools for computational molecular biology  
867 and bioinformatics. *Bioinformatics* **25**, 1422-1423 (2009).
- 868 82. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree  
869 display and annotation. *Nucleic Acids Res* **49**, W293-W296 (2021).
- 870 83. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data  
871 analysis. *Genome Biol* **19**, 15 (2018).
- 872 84. Glanville, J. et al. Identifying specificity groups in the T cell receptor repertoire. *Nature* **547**,  
873 94-98 (2017).

- 874 85. Wu, T. et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.  
875 *Innovation (Camb)* **2**, 100141 (2021).
- 876 86. Chan, M.M. et al. Molecular recording of mammalian embryogenesis. *Nature* **570**, 77-82  
877 (2019).

878

## 879 **Methods**

### 880 **The mathematical framework of PhyloVelo**

881 The dynamics of the latent expression  $z$  for each gene on a phylogeny  $\mathcal{T}$  was assumed to  
882 follow a diffusion process (also known as the stochastic differential equation, SDE), which  
883 varies along cell divisions:

$$884 \quad dz_t = v(t, z_t)dt + \sigma(t, z_t)dW_t \quad (1)$$

885 Here,  $W_t$  is a standard Brownian motion. In our model, we hypothesized that there is a  
886 group of genes  $G_m$  whose drift coefficient  $v(t, z_t)$  and diffusion coefficient  $\sigma(t, z_t)$  are  
887 independent of both  $t$  and  $z_t$ , thus  $v(t, z_t) = v$  and  $\sigma(t, z_t) = \sigma$ . We called them  
888 monotonically expressed genes (MEGs). For this type of genes, the dynamics of its latent  
889 expression  $z$  is thus formulated as:

$$890 \quad dz_t = vdt + \sigma dW_t \quad (2)$$

891 and its expectation is given by:

$$892 \quad \mathbb{E}(z_t) = \mathbb{E}(z_{t_0})v(t - t_0) \quad (3)$$

893 For the observed scRNA-seq measurement  $x$  (read or UMI count), we assumed that it is  
894 sampled from the negative binomial (NB) distribution or the zero-inflated negative binomial  
895 (ZINB) distribution:

$$896 \quad \mathbb{P}(x|z', \alpha, \psi) = \begin{cases} (1 - \psi) + \psi \left(\frac{\alpha}{\alpha + z'}\right)^\alpha, & x = 0 \\ \psi \frac{\Gamma(x + \alpha)}{x! \Gamma(\alpha)} \left(\frac{\alpha}{z' + \alpha}\right)^\alpha \left(\frac{z'}{z' + \alpha}\right)^x, & x \geq 1 \end{cases} \quad (4a)$$

897 where  $z'$  is the exponential function of latent expression  $z$ ,  $\alpha$  is the scale parameter, and  $\psi$   
 898 is the zero-inflation parameter. The expectation of the distribution is  $z'\psi$ . For the negative  
 899 binomial distribution,  $\psi = 1$ . We used the likelihood ratio test to verify zero inflation for each  
 900 gene (**Supplementary Note**).

901

902 For the scRNA-seq data after normalization (e.g. using *scanpy.pp.normalize\_per\_cell* and  
 903 *scanpy.pp.log1p*), we also provided a Gaussian model of latent expression for normalized  
 904 data:

905

$$906 \quad \mathbb{P}(x|z, \alpha, \psi) = \begin{cases} (1 - \psi), & x = 0 \\ \psi \frac{1}{\sqrt{2\pi\alpha}} \exp\left(-\frac{(x - z)^2}{2\alpha}\right), & \text{otherwise} \end{cases} \quad (4b)$$

907

908 To estimate the latent expression  $z$ , we used the maximum a posteriori probability (MAP)  
 909 estimate. For the ZINB model using raw UMI count data, we estimated  $z'$  and then took  
 910 logarithm to get the estimated latent expression  $z$ :

$$911 \quad \hat{z}_{\text{MAP}}(x) = \log\left(\underset{z'}{\operatorname{argmax}}(\mathbb{P}(x|z, \alpha, \psi)\mathbb{P}(z'))\right) \quad (5a)$$

912 For the Gaussian model using normalized UMI count data, we directly performed the MAP  
 913 estimate of the latent expression  $z$ :

$$914 \quad \hat{z}_{\text{MAP}}(x) = \underset{z}{\operatorname{argmax}}(\mathbb{P}(x|z, \alpha, \psi)\mathbb{P}(z)). \quad (5b)$$

915 For a MEG  $g$ , its drift coefficient can be estimated as:

$$916 \quad v_g = \frac{\mathbf{Z}^T \mathbf{B} - n\bar{\mathbf{Z}}\bar{\mathbf{B}}}{\mathbf{B}^T \mathbf{B} - n\bar{\mathbf{B}}^2} \quad (6)$$

917 Here  $\mathbf{z} = (\hat{z}_1, \hat{z}_2, \dots, \hat{z}_n)$  represents the estimated latent expressions and  $\mathbf{B} = (b_1, b_2, \dots, b_n)$   
 918 the phylogenetic distances from terminal cells to the root ( $n$  is the cell number in a dataset).  
 919 The drift coefficients of all  $G$  MEGs in a dataset  $\mathbf{v} = (v_1, v_2, \dots, v_G)$  were thus referred to be  
 920 as phylogenetic velocity.

921

922 For clonal lineage-resolved scRNA-seq data by lentiviral barcoding or immune cell receptor  
923 sequences,  $\mathbf{B}$  represent the logarithm of clonal sizes of individual cells at the time of  
924 sampling, namely  $\mathbf{B} = (\log(c_1), \log(c_2), \dots, \log(c_n))$ , where  $c$  is the size of corresponding  
925 clone that a cell belongs to and  $n$  the cell number. The estimation of  $\nu$  is the same as using  
926 phylogeny-resolved scRNA-seq data.

## 927 **Simulation of phylogeny-resolved scRNA-seq data**

928 To generate simultaneous single-cell phylogenetic and transcriptomic data *in silico*, a  
929 lineage-embedded scRNA-seq data simulator, PROSSTT<sup>41</sup>, was modified to account for  
930 dividing cell populations, so that the whole cell division history initiated from a single cell can  
931 be recorded. The simulation consisted of three parts:

- 932 1) Simulate a cell division and differentiation process using the Gillespie algorithm to obtain  
933 the cell division history;
- 934 2) Given a cell differentiation model (linear, bifurcated or convergent), use the diffusion  
935 process to generate gene expression programs;
- 936 3) Assign the gene expression programs onto the cell division history in order to obtain the  
937 read/UMI count data for each gene in each cell.

938 ***Simulating cell division history and mutation-based phylogeny.*** We used a continuous-  
939 time Markov process to simulate cell division and differentiation. In particular, each cell type  
940  $i$  has a specific division rate  $p_i(t)$  and differentiation rate  $q_{ij}(t)$ , given as follows:

$$\begin{aligned} 941 \quad p_i(t) &= r_i \left( 1 - \frac{1}{1 + e^{-k_i(t-t_{0i})}} \right) \\ q_{ij}(t) &= p_{ij} r_i \left( \frac{1}{1 + e^{-k_i(t-t_{0i})}} \right) \end{aligned} \quad (7)$$

942 where  $r_i$ ,  $k_i$  and  $t_{0i}$  are the cell-type specific parameters and  $p_{ij}$  is the probability of cell type  
943  $i$  differentiating into cell type  $j$ ,  $\sum_{i \neq j} p_{ij} = 1$ ,  $p_{ii} = 0$ .

944

945 Now, we can simulate the cell growth process using the Gillespie algorithm<sup>75</sup>. Each  
946 simulation ended when the population size reached 10,000 cells. Then, 1,000 cells were  
947 randomly sampled to obtain their division history.

948

949 To simulate the mutation-based cell phylogeny, we assumed that mutations randomly occur  
950 during each cell division following a Poisson distribution:

951 
$$P(m = i) = \frac{u^i e^{-u}}{i!} \quad (8)$$

952 where  $u$  is the mean mutation rate per cell division. Different mutation rates ( $u = 0.1, 0.3,$  or  
953  $1$ ) were used. After obtaining the cell mutational information, we used three different  
954 algorithms to reconstruct the phylogenetic tree, respectively, namely *Maximum Likelihood*  
955 (using IQ-TREE 2<sup>76</sup>), *Neighbor-Joining* (using R package ape 5.6-2<sup>77</sup>) and *Maximum*  
956 *Parsimony* (using R package phangorn 2.11.1).

957

958 ***Simulating scRNA-seq data.*** We first simulated the latent expression process of genes.

959 For each gene, we randomly generated its initial expression  $\mu_0$ , drift coefficient  $\nu$ , and  
960 variance  $\sigma^2$ , and then simulated gene-specific diffusion process as follows:

961 
$$z_{t+dt} \sim Normal(loc = z_t + \nu dt, scale = \sigma^2) \quad (9)$$

962 When cells differentiated at time  $t_d$ , for MEGs, their gene expressions remained unchanged  
963 as the same with the values in the previous process. For a non-MEG, its drift coefficient and  
964 variance were regenerated randomly and the value of expression was reset to  $z_{t_d}$ . We

965 called the diffusion process  $z_t$  as the gene expression program. For each gene, the initial  
966 value of the expression program  $z_0$  was randomly drawn from a gamma distribution  $z_0 \sim$

967  $\Gamma(0.5, 20)$ , the drift coefficient  $\nu$  was drawn from a normal distribution  $\nu \sim Normal(0, 1)$ , and

968 the diffusion coefficient  $\sigma$  was drawn from a truncated normal distribution  $\sigma \sim$

969  $TruncatedNormal\left(kz_0, \frac{|kz_0|}{3}, \min = 0.00001\right)$ . In all simulations, we set the drift coefficient of

970 each gene expression program to change with probability 0.4 when cell differentiation  
971 happens, ultimately resulting in 10-15% of genes with unchanged expression programs  
972 upon cell differentiations, thus behaving as MEGs. The other 85-90% of genes will change  
973 dynamically with cell differentiations and thus behave as non-MEGs. Each cell was  
974 assumed to have 2,000 expressed genes, thus including 200-300 MEGs in total in each  
975 simulated dataset.

976

977 After generating the latent expression process of all genes, in order to simulate the  
978 variations introduced in real scRNA-seq experiments, the NB<sup>78</sup> or ZINB distribution<sup>79</sup> was  
979 used to obtain read/UMI count  $x$ :

$$980 \quad x \sim \text{ZINB}(\psi, z, \alpha), \quad (10)$$

981 where  $\psi$  is the zero-inflation parameter ( $\psi = 1$  for negative binomial model),  $\alpha$  is the scale  
982 parameter and the expectation of the distribution is  $\psi z$ .

983

984 **Assigning the gene expression programs to the cell division history.** Having each  
985 gene expression program  $(z_1, z_2, \dots, z_G)$  and phylogeny  $\mathcal{T}$ , we traversed all nodes  $V \in \mathcal{T}$  and  
986 assigned the latent expression programs  $z_1(d), \dots, z_G(d)$  to the nodes with branch length of  
987  $d$ . For each gene  $g$ , a random number obeying the distribution  $z_g(d)$  was drawn as the  
988 latent expression of that gene. By traversing all the genes and cells, the latent expression  
989 matrix can be obtained and denoted as  $Z$ . We also simulated a Gaussian noise  $\varepsilon$  to be  
990 imposed on  $Z$ , thus the latent expression matrix would be updated as  $Z + \varepsilon$ . Finally, given  
991 the zero-inflation factor  $\psi$  and the scale parameter  $\alpha$ , the expression matrix  $X$  of the  
992 simulated data can be obtained by random sampling according to Equation (10).

### 993 **Inference of PhyloVelo pseudotime**

994 To infer the PhyloVelo pseudotime (forward) of each cell, we first constructed a minimum  
995 spanning tree based on the distance of cellular states on tSNE/UMAP embedding using

996 Prim's algorithm<sup>80</sup>. Thus, we can obtain a subset of the edges  $\mathcal{E}$  connecting all cells together  
997 with a minimum possible total distance. We then chose any cell  $c_0$  as the starting point and  
998 set its pseudotime  $pt_{c_0} = 0$ . For any other cells  $c \in U\{e \in \mathcal{E}: c_0 \in e\}$ , we calculated its  
999 pseudotime using the following equation:

$$1000 \quad pt_c = pt_{c_0} + \int_{x_{c_0}}^{x_c} \frac{1}{v_{emb}(x)} dx \quad (11)$$

1001 where  $x_c$  is the coordinate of cell  $c$  in the embedding space and  $v_{emb}(x)$  is the phylogenetic  
1002 velocity in the embedding space and varies with its coordinates.

1003

1004 To simplify the calculation, we replaced the velocity in this path with the average velocity of  
1005  $v(c)$  and  $v(c_0)$ , denoted by  $v_a$ , and used the line segments  $\mathbf{l}_{c,c_0}$  to approximate the path.

1006 Hence, we have:

$$1007 \quad pt_c = pt_{c_0} + \frac{\|\mathbf{l}_{c,c_0}\|_2^2}{\mathbf{v}_a^T \mathbf{l}_{c,c_0}} \quad (12)$$

1008 Following the path generated from the minimum spanning tree, we can estimate the  
1009 pseudotime of all cells and finally normalize to  $[0,1]$ . It should be noted that although  
1010 PhyloVelo velocity fields are in backward directions, PhyloVelo pseudotime is still set to be  
1011 forward as scVelo latent time.

## 1012 **Analysis of phylogeny-resolved scRNA-seq datasets**

1013 **Datasets and pre-processing.** We have applied PhyloVelo to six real phylogeny-resolved  
1014 scRNA-seq datasets that are publicly available through online sources (see **Data**  
1015 **availability**). These included *C. elegans*<sup>14</sup>, mouse embryos<sup>32</sup>, GEMM of lung  
1016 adenocarcinoma<sup>51</sup>, mouse xenograft models using pancreatic cancer cell line KPCY<sup>62</sup> and  
1017 lung cancer cell line A549<sup>63</sup>, and *in vitro* culture of human kidney cell line HEK293T<sup>64</sup>. The  
1018 embryonic lineage tree of *C. elegans* is entirely known and was obtained from  
1019 <http://dulab.genetics.ac.cn/TF-atlas/Cell.html>, while the CRISPR-based lineage trees in



1020 other five datasets were obtained from the original studies which were reconstructed by the  
1021 mutational scars on CRISPR lineage barcodes. In the *C. elegans* dataset, because multiple  
1022 synchronous embryos were pooled for the scRNA-seq experiment, many nodes in the  
1023 lineage tree have been sampled multiple times. Thus, only one random cell was chosen to  
1024 represent the corresponding node, while these non-repetitive cells (~300 cells) from one  
1025 lineage tree constituted a “pseudo-embryo”. For the mouse embryos (E8.0/8.5)<sup>32</sup>, four  
1026 (embryos 1,2,3 and 6) out of seven embryos were analyzed for their higher barcode  
1027 diversity where the number of unique barcode alleles was > 500 in each embryo. For the  
1028 scRNA-seq data of *C. elegans*<sup>14</sup> and mouse lung adenocarcinoma<sup>51</sup>, the coordinates of  
1029 tSNE or UMAP from the original studies were used. All phylogenetic trees were read and  
1030 branch lengths were calculated using biopython<sup>81</sup> and visualized using iTOL<sup>82</sup>. For the  
1031 scRNA-seq data of mouse embryos<sup>32</sup>, cell lines KPCY<sup>62</sup>, A549<sup>63</sup> and HEK293T<sup>64</sup>, the  
1032 dimensionality reduction and tSNE or UMAP visualization were performed using Scanpy<sup>83</sup>  
1033 following the recommended data processing procedures and parameters as [https://scanpy-](https://scanpy-tutorials.readthedocs.io/en/latest/)  
1034 [tutorials.readthedocs.io/en/latest/](https://scanpy-tutorials.readthedocs.io/en/latest/). In each dataset, the genes with total count < 20 were  
1035 filtered out.

1036

1037 **Applying PhyloVelo.** For *C. elegans*, whose embryonic cell division history is entirely  
1038 known, the cell generation time was used to denote the phylogenetic distance. For the other  
1039 five CRISPR/Cas9 lineage tracing datasets<sup>32, 51, 62-64</sup>, the phylogenetic distance on a lineage  
1040 tree corresponds to the number of Cas9 cutting scars on the evolving barcodes. To estimate  
1041 the latent gene expressions, for *C. elegans*, the ZINB model was used to analyze the raw  
1042 UMI count data because of the high-quality lineage tree. For the CRISPR/Cas9-based  
1043 lineage tracing datasets, the Gaussian model was used on the post-normalized data where  
1044 `normalize_per_cell()` and `log1p()` by Scanpy<sup>83</sup> were applied to the raw UMI counts. To  
1045 prioritize the high-confident candidates of MEGs and speed up the computation, rather than  
1046 estimating the latent expression for all genes, we firstly searched for candidate MEGs by

1047 directly analyzing the correlation between each gene's normalized UMI counts and the  
1048 phylogenetic distances to root of single cells. The top 5% of genes with the highest  
1049 Spearman's correlations were first selected and then proceeded for follow-up latent  
1050 expression estimations. Final MEGs were identified by the significant association (Pearson's  
1051 correlation,  $q < 0.05$  after Benjamini-Hochberg correction; a stringent threshold  $q < 10^{-5}$  was  
1052 used for Chan *et al.* dataset<sup>32</sup> given the large number of cells in individual embryos each  
1053 with 6,328-19,071 cells) between the latent expressions and the phylogenetic distances  
1054 from terminal nodes to the root of tree. The phylogenetic velocity was computed  
1055 independently for each MEG. To project the phylogenetic velocity into the dimensionality  
1056 reduction embedding, we built a k-nearest neighbor (kNN) graph (k=15 for *C. elegans*  
1057 dataset while it was chosen by approximate to one third of total number of cells for the  
1058 CRISPR lineage tracing datasets). The kNN graph was based on the Euclidean distance as  
1059 the base vector and was used to estimate the coordinates of velocity embedding, as the  
1060 projection of RNA velocity<sup>22, 23</sup>.

1061

1062 **Applying scVelo.** The spliced and unspliced read counts were obtained by running  
1063 velocity (v0.6)<sup>22</sup> on the bam files from the output of CellRanger (6.0.2) using the raw  
1064 sequence reads. To estimate RNA velocity, scVelo (version 0.2.4)<sup>23</sup> and the dynamical  
1065 mode were used following the recommended data processing procedures as  
1066 <https://scvelo.readthedocs.io/VelocityBasics/>. Spliced/unspliced read counts were pre-  
1067 processed using the following default setting:

```
1068 scv.pp.filter_and_normalize(adata, min_shared_counts=20, n_top_genes=2000)
```

```
1069 scv.pp.moments(adata, n_neighbors=30, n_pcs=30)
```

1070

1071 **Applying VeloVAE.** VeloVAE<sup>58</sup> applies the same data preprocessing steps as scVelo.

1072 There are three data training models including Basic Model (assuming fixed transcription  
1073 rates), Full Model (assuming variable transcription rates) and Full VB Model (treating the

1074 rate parameters as random variables). Full Model was used as recommended in the paper.  
1075 The model training parameters were used following the example in its GitHub repository  
1076 ([https://github.com/welch-lab/VeloVAE/blob/master/notebooks/velovae\\_example.ipynb](https://github.com/welch-lab/VeloVAE/blob/master/notebooks/velovae_example.ipynb)).

```
1077 vae = vv.VAE(adata, tmax=20, dim_z=5)  
1078 vae.train(adata, gene_plot=gene_plot, plot=True, figure_path=figure_path)
```

1079

1080 **Applying DeepVelo.** DeepVelo<sup>59</sup> also applies the same data preprocessing steps as  
1081 scVelo. Model configurations were the same as the default setting

1082 (<https://github.com/bowang-lab/DeepVelo/blob/main/examples/figure2.ipynb>) except some  
1083 updates as following:

```
1084 configs = dict(  
1085     "name": "DeepVelo", # name of the experiment  
1086     "loss": dict("args": dict("coeff_s": autoset_coeff_s(adata))),  
1087     "trainer": dict("verbosity": 0), # increase verbosity to show training progress  
1088     "n_gpu":0  
1089 )
```

```
1090 configs = update_dict(Constants.default_configs, configs)
```

1091

1092 **Applying cellDancer.** cellDancer<sup>60</sup> applies the same data preprocessing steps as scVelo.

1093 The format conversion of the data is according to its tutorial

1094 ([https://guangyuwanglab2021.github.io/cellDancer\\_website/index.html](https://guangyuwanglab2021.github.io/cellDancer_website/index.html)), and the velocity

1095 inference uses all genes and proceeds according to the default parameters. The velocity

1096 field is visualized using the Dynamo<sup>50</sup>.

1097

1098 **Applying UniTVelo.** UniTVelo<sup>61</sup> also applies the same data preprocessing steps as scVelo.

1099 Model configurations were the same as the default setting

1100 ([https://unitvelo.readthedocs.io/en/latest/Figure2\\_ErythroidMouse.html](https://unitvelo.readthedocs.io/en/latest/Figure2_ErythroidMouse.html)):

1101 `velo_config = utv.config.Configuration()`

1102 `velo_config.R2_ADJUST = True`

1103 `velo_config.IROOT = None`

1104 `velo_config.FIT_OPTION = '1'`

1105 `velo_config.AGENES_R2 = 1`

1106

1107 ***Applying Dynamo.*** `Dynamo`<sup>50</sup> was used to infer the quantitative cell-state transition matrix  
1108 and visualize cell state transition graph. We use the velocity field inferred by `PhyloVelo` as  
1109 input and calculate the transition matrix as follows:

1110 `dyn.vf.VectorField(adata, basis='umap', M=1000, pot_curl_div=True)`

1111 `dyn.vf.topography(adata, basis='umap')`

1112 `dyn.ext.ddhodge(adata, basis='umap')`

1113 `dyn.pd.state_graph(adata, group='cell_states', basis='umap', method='vf', approx=False)`

#### 1114 **Analysis of static barcoding-based lineage tracing datasets**

1115 ***Datasets and pre-processing.*** We have applied the extended model of `PhyloVelo` to two  
1116 static barcoding datasets including LARRY hematopoietic differentiation<sup>37</sup> and intratumoral  
1117 CD8+ T cells in BCC<sup>57</sup> that are publicly available through online sources (see **Data**  
1118 **availability**). For the LARRY dataset, lentiviral barcoding data at day 6 was used to obtain  
1119 the clone size information for each cell. For the CD8+ T cells data, the TCR specificity  
1120 clones were identified by `GLIPH`<sup>84</sup> which defines clones based on the following two criteria:  
1121 1) global similarity, TCR sequences within the same T cell clone have at most one amino  
1122 acid difference; 2) local similarity, two TCRs in same clone contain an identical CDR3 motif,  
1123 which is 2-4 k-mer amino acids in length and is significantly enriched from random sub-  
1124 sampling of unselected repertoires. To avoid batch effect, patient 9 with the largest cell  
1125 number (4,659 cells) was selected for identification of MEGs and inference of phylogenetic  
1126 velocities. The phylogenetic velocities were then transferred to all CD8+ T cells (12,788

1127 cells) from all 12 BCC patients. Cells whose clonal barcodes were not determined were  
1128 filtered out. The coordinates of dimensionality reduction embedding, SPRING (LARRY  
1129 dataset) or UMAP (T cell dataset), from the original studies were used for visualization. In  
1130 each dataset, the genes with total count < 20 were filtered out.

1131

1132 **Applying PhyloVelo.** For both datasets, the phylogenetic time of a cell corresponded to the  
1133 logarithm of the clone size. To estimate the latent gene expressions, Gaussian process  
1134 model was used on the post-normalized data. Same as the scRNA-seq data analysis  
1135 above, *normalize\_per\_cell()* and *log1p()* by *Scanpy*<sup>83</sup> were applied to the raw UMI counts.  
1136 The top 5% genes with the highest Spearman's correlation between normalized gene  
1137 expression and phylogenetic time were first selected, then proceeded for follow-up latent  
1138 expression estimations. Final MEGs were identified by the significant association (Pearson's  
1139 correlation,  $q < 0.05$ ) between the latent expressions and the logarithm of clone size.  
1140 Projecting phylogenetic velocities into the embedding followed the same procedure as  
1141 CRISPR lineage tracing data analysis.

#### 1142 **Transferring the phylogenetic velocities of MEGs to independent datasets**

1143 To evaluate whether the phylogenetic velocities of MEGs estimated from one phylogeny-  
1144 resolved scRNA-seq dataset are sufficiently robust to infer the velocity fields in independent  
1145 datasets in the absence of phylogenetic information, three datasets were analyzed including  
1146 *C. elegans*<sup>14</sup>, mouse erythroid development<sup>19, 32</sup>, and the GEMM of lung adenocarcinoma<sup>51</sup>.  
1147 Here, the MEGs and corresponding phylogenetic velocity estimates were directly applied to  
1148 another scRNA-seq datasets in similar biological conditions. For *C. elegans*, we applied the  
1149 phylogenetic velocities from AB lineage in a single pseudo-embryo (n=298 cells) to all AB  
1150 lineage cells (n=29,600) in multiple embryos. We also applied them to non-AB lineages that  
1151 differentiate to hypodermis, body wall muscles (BWM) and pharynx, respectively. For  
1152 mouse erythroid differentiation, we applied the phylogenetic velocity estimates in the

1153 erythroid lineage cells from a single embryo (E8.5, n=2,419 cells) of the Chan *et al.*  
1154 dataset<sup>32</sup> to the other three embryo and the temporally-sequenced mouse embryos (E6.5-  
1155 E8.5, n=12,324 cells) of the Pijuan-Sala *et al.* dataset<sup>19</sup>. We also applied the phylogenetic  
1156 velocities of LT-MEGs identified from the Chan *et al.* dataset<sup>32</sup> to predict the entire embryo  
1157 development with the Pijuan-Sala *et al.* dataset<sup>19</sup> (E6.5-E8.5, n=10,000 out of 116,312 cells  
1158 were randomly sampled), and predict mouse brain development with Manno *et al.* dataset<sup>48</sup>  
1159 (E7-18, n=10,000 out of 292,495 cells were randomly sampled). For lung  
1160 adenocarcinoma<sup>51</sup>, the phylogenetic velocity estimates in one KP primary lung tumor  
1161 (3726\_NT\_T1, n=754 cells) were applied to all 58,022 single cells from all pooled KP  
1162 primary lung tumors. Finally, for intratumoral CD8+ T cells, in order to avoid the batch effect,  
1163 the phylogenetic velocity estimates in 4,659 cells from patient 9 were applied to all 12,788  
1164 CD8+ T cells from 12 BCC patients.

### 1165 **Gene ontology (GO) enrichment analysis**

1166 GO enrichment analysis was performed using clusterProfiler v4.4.4<sup>85</sup>. The cutoff for *p* value  
1167 and *q* value were set to 0.05 and 0.25, respectively. After excluding Cellular Components  
1168 (CC) terms, all significant terms were retained for downstream analyses. Subsequently, top  
1169 20 GO terms of each sample were merged and these terms were sorted by their total  
1170 occurrence and mean *q* value across samples. Finally, the top 20 GO terms enriched were  
1171 visualized using ggplot2 v3.4.0.

### 1172 **Data availability**

1173 All data analyzed in this article are publicly available through online sources. The annotated  
1174 data, lineage trees, results and Python implementation are available at  
1175 <https://phylovelo.readthedocs.io/>. The raw data for the *C. elegans* dataset<sup>14</sup> can be  
1176 accessed with [GSE126954](https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE126954) and the lineage tree can be accessed from  
1177 <http://dulab.genetics.ac.cn/TF-atlas/Cell.html>. The CRISPR lineage tracing datasets from

1178 the mouse embryos<sup>32, 86</sup> can be accessed with [GSE117542](#). The single cell RNA-seq data  
1179 of mouse brain development<sup>48</sup> can be accessed with [PRJNA637987](#). The time-course  
1180 single-cell RNA-seq data of whole mouse embryos (E6.5-8.5)<sup>19</sup> can be accessed with [E-](#)  
1181 [MTAB-6967](#). The dataset of mouse primary lung tumors<sup>51</sup> can be accessed with  
1182 [PRJNA803321](#) and from Zenodo (<https://zenodo.org/record/5847462#.Yt4-PewRXUI>). The  
1183 dataset of mouse pancreatic cancer cell line KPCY<sup>62</sup> can be accessed with [GSE173958](#) and  
1184 from Mendeley (<https://doi.org/10.17632/t98picd7t6.1>). The dataset of human lung cancer  
1185 cell line A549<sup>63</sup> can be accessed with [GSE161363](#). The dataset of human kidney cell line  
1186 HEK293T<sup>64</sup> can be accessed with [PRJNA757179](#). The LARRY lentiviral barcoding dataset  
1187 of hematopoiesis<sup>37</sup> can be accessed with [GSE140802](#). The single-cell TCR and RNA  
1188 sequencing data of T cells in BCC<sup>57</sup> can be accessed with [GSE123813](#).

## 1189 Code availability

1190 PhyloVelo<sup>87</sup> is freely available as Python package at  
1191 <https://github.com/kunwang34/PhyloVelo>. Detailed workflows to reproduce figures and  
1192 results in this paper are written as Jupyter notebook in the repository. The annotated data,  
1193 lineage trees, results and Python implementation are available at  
1194 <https://phylovelo.readthedocs.io/>.

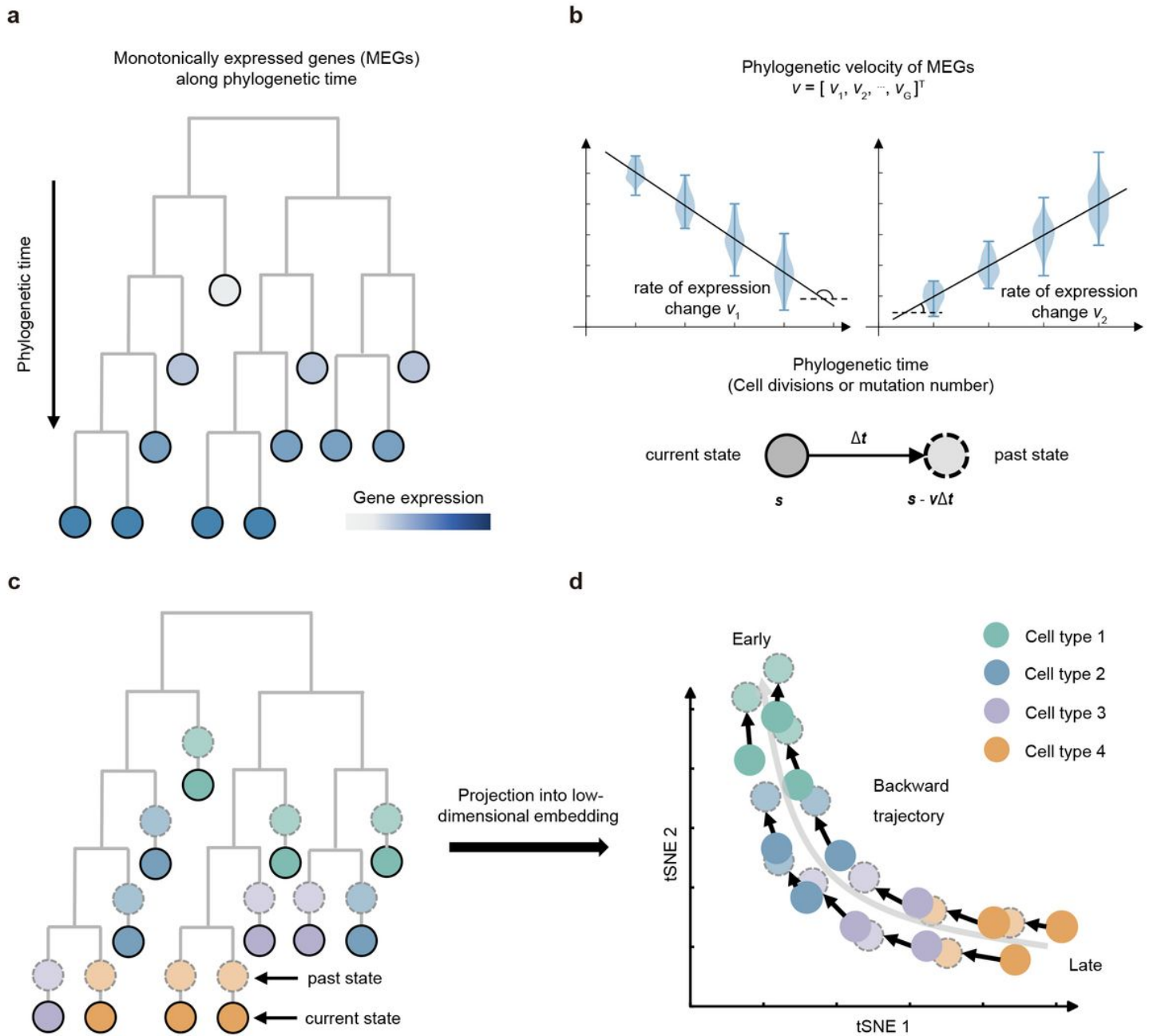
## 1195 Methods-only references

- 1196 75. Gillespie, D.T. Exact stochastic simulation of coupled chemical reactions. *The journal of*  
1197 *physical chemistry* **81**, 2340-2361 (1977).
- 1198 76. Minh, B.Q. et al. IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in  
1199 the Genomic Era. *Mol Biol Evol* **37**, 1530-1534 (2020).
- 1200 77. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R  
1201 language. *Bioinformatics* **20**, 289-290 (2004).
- 1202 78. Chen, W. et al. UMI-count modeling and differential expression analysis for single-cell RNA  
1203 sequencing. *Genome Biol* **19**, 70 (2018).

- 1204 79. Jia, C. Kinetic foundation of the zero-inflated negative binomial model for single-cell RNA  
1205 sequencing data. *SIAM Journal on Applied Mathematics* **80**, 1336-1355 (2020).
- 1206 80. Prim, R.C. Shortest connection networks and some generalizations. *The Bell System Technical*  
1207 *Journal* **36**, 1389-1401 (1957).
- 1208 81. Cock, P.J. et al. Biopython: freely available Python tools for computational molecular biology  
1209 and bioinformatics. *Bioinformatics* **25**, 1422-1423 (2009).
- 1210 82. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v5: an online tool for phylogenetic tree  
1211 display and annotation. *Nucleic Acids Res* **49**, W293-W296 (2021).
- 1212 83. Wolf, F.A., Angerer, P. & Theis, F.J. SCANPY: large-scale single-cell gene expression data  
1213 analysis. *Genome Biol* **19**, 15 (2018).
- 1214 84. Glanville, J. et al. Identifying specificity groups in the T cell receptor repertoire. *Nature* **547**,  
1215 94-98 (2017).
- 1216 85. Wu, T. et al. clusterProfiler 4.0: A universal enrichment tool for interpreting omics data.  
1217 *Innovation (Camb)* **2**, 100141 (2021).
- 1218 86. Chan, M.M. et al. Molecular recording of mammalian embryogenesis. *Nature* **570**, 77-82  
1219 (2019).
- 1220 87. Wang, K., Hou, L., Wang, X., Zhai, X., Lu, Z., Zi, Z., Zhai, W., He, X., Curtis, C., Zhou, D., & Hu, Z.  
1221 PhyloVelo, Phylogeny-based transcriptomic velocity of single cells. GitHub.  
1222 <https://github.com/kunwang34/PhyloVelo> (2023).
- 1223
- 1224
- 1225



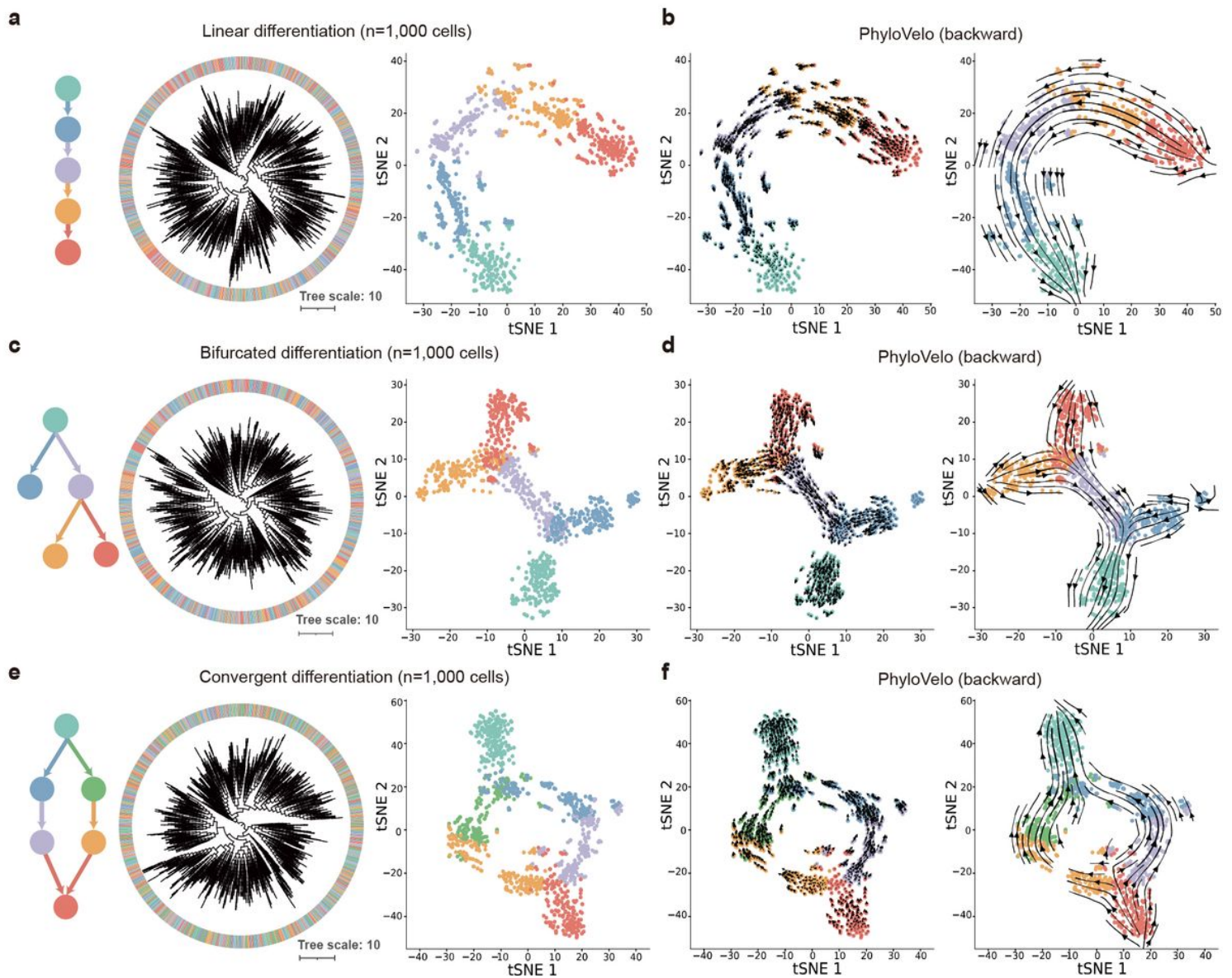
# Figures



**Figure 1**

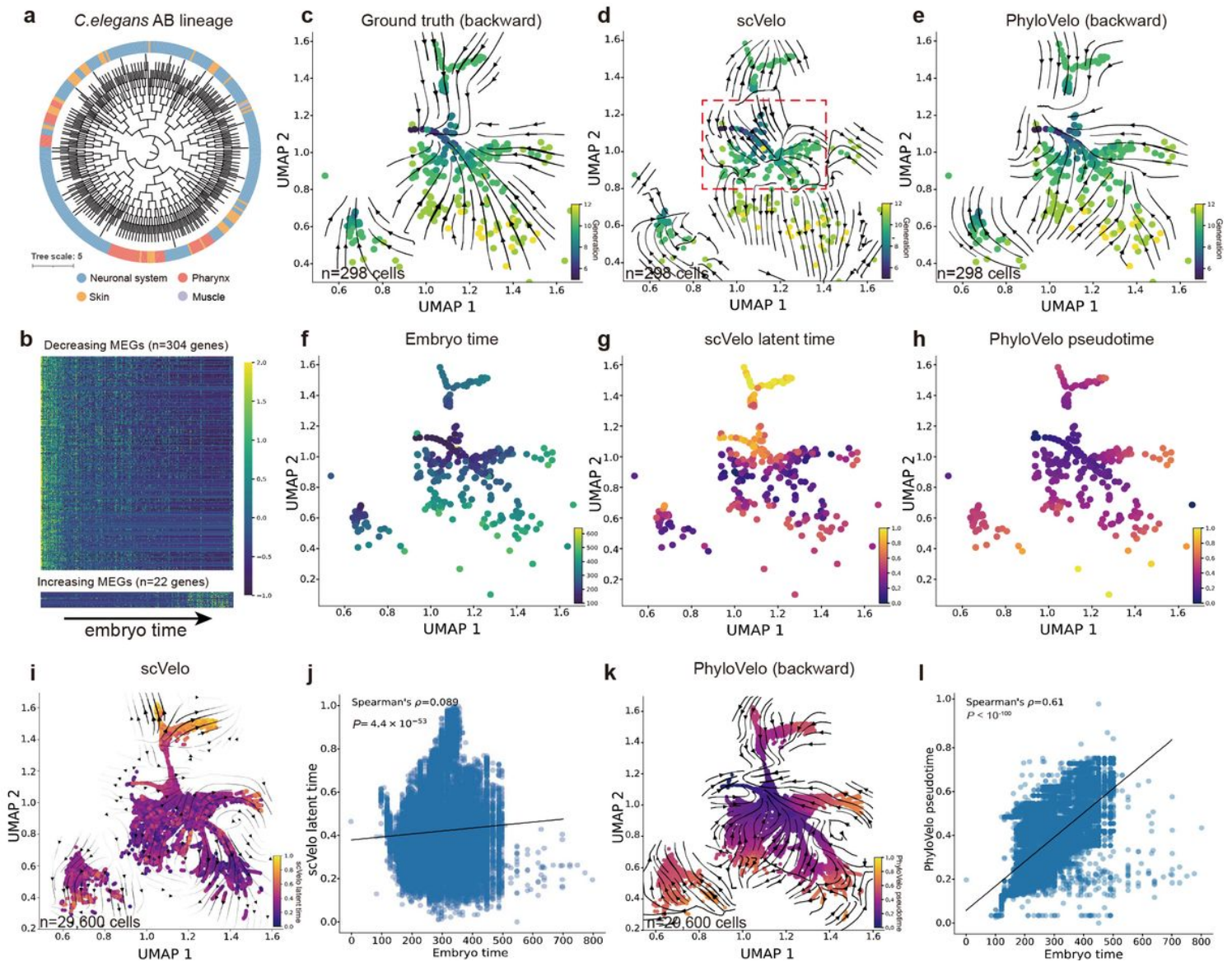
**Schematic of the PhyloVelo framework.** (a) Schematic of monotonically expressed genes (MEGs) over phylogenetic time on a cell phylogenetic tree. (b) Two examples of MEGs whose latent expressions are associated with the phylogenetic time (cell divisions or mutation number). A diffusion process of gene expressions was used to model the changes of latent expressions over phylogenetic time. This enables the estimation of the phylogenetic velocity,  $v = (v_1, v_2, \dots, v_G)$ , which corresponds to the drift coefficients of  $G$  MEGs in the diffusion process (approximate to the slope of linear regression between latent expression and phylogenetic time). Whiskers: minimum and maximum. (c) Phylogenetic velocity predicts the past transcriptional state of a cell before a unit of phylogenetic time (one cell division or mutation). (d)

Projection of the phylogenetic velocity into low dimensional embedding enables the mapping of cell-state trajectory in backward directions.



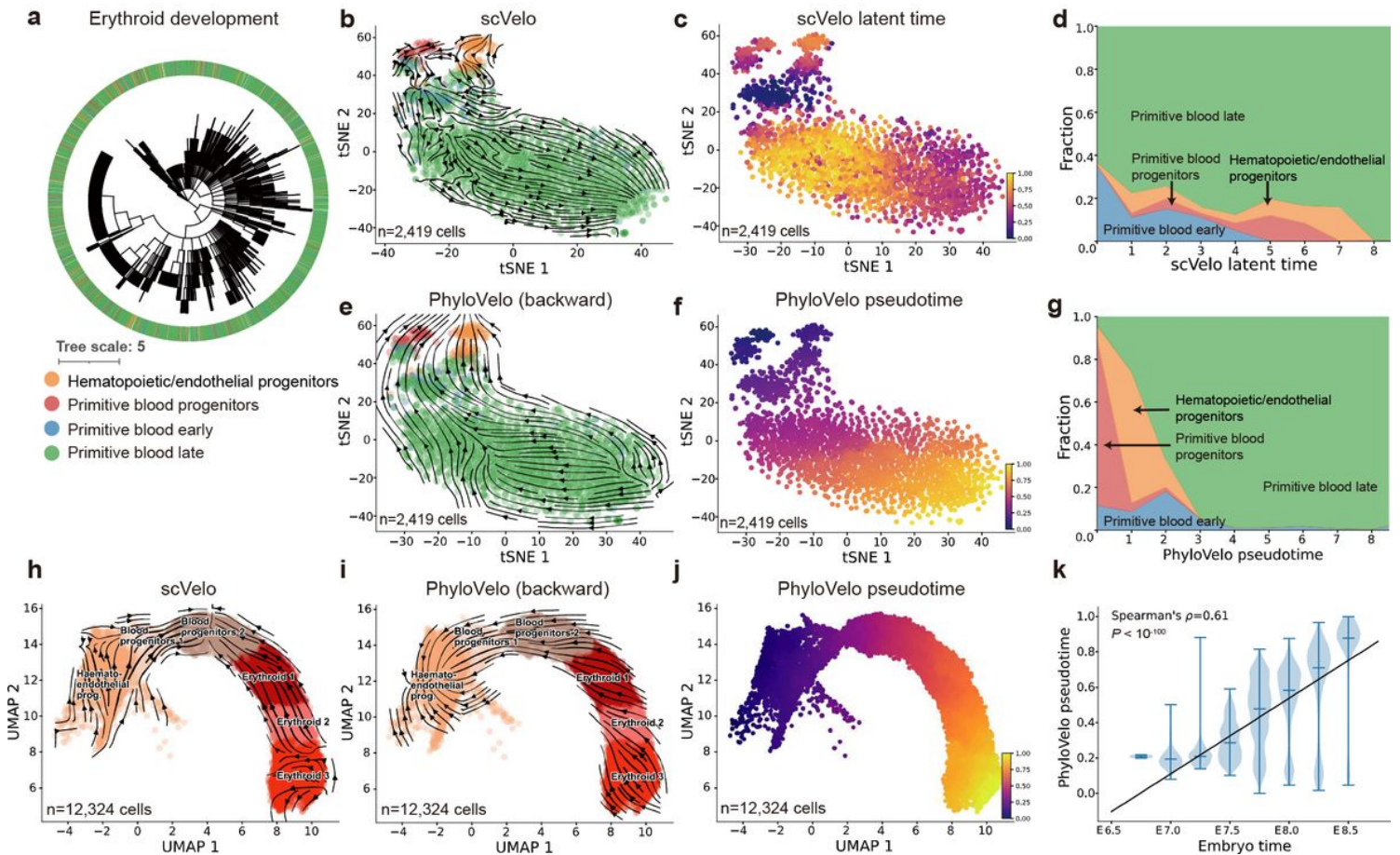
**Figure 2**

**PhyloVelo recovers complex cell lineages in simulations.** Simulation of single-cell RNA-seq data and paired cell-division history under linear (a), bifurcated (b), and convergent (c) differentiation models, respectively. Colors are labeled by cell types. Each simulation consists of 1,000 cells randomly sampled from a growing cell population at 10,000 cells. Each cell has 2,000 expressed genes, including 200-300 MEGs. (d-f) Phylogenetic velocity fields reconstructed by PhyloVelo for the corresponding differentiation scenarios. The left panel shows the single-cell level of velocity fields, while the right panel shows the same velocity fields visualized as streamlines in scVelo. PhyloVelo velocity fields are at backward directions.



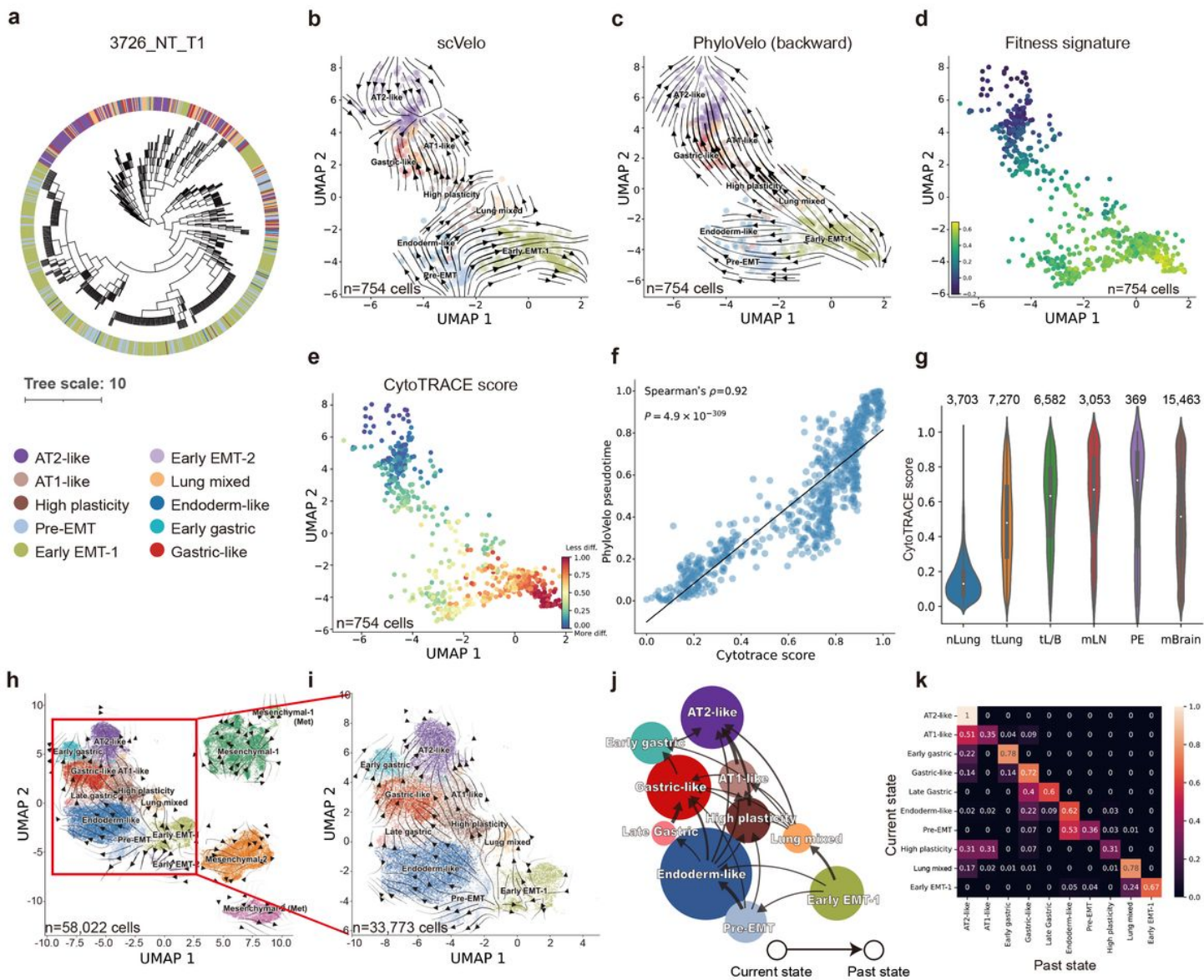
**Figure 3**

**PhyloVelo reconstructs the embryonic differentiation trajectories of *C. elegans*.** (a) Phylogenetic tree of the *C. elegans* AB lineage. (b) Heatmap showing the expressions (z-score normalized) of MEGs along *C. elegans* embryo time. (c) The ground-truth velocity fields represent vectors superimposed on the cells that point to their immediate parental cells on the Uniform Manifold Approximation and Projection (UMAP) plot. (d-e) The velocity fields estimated by scVelo (dynamical mode) (d) or PhyloVelo (e). Dash square indicates the early embryonic lineages where RNA velocity gave erroneous estimations on the fate directions. (f) *C. elegans* embryo time as Packer *et al.*<sup>14</sup>. (g) scVelo latent time. (h) PhyloVelo pseudotime. (i) RNA velocity fields for all 29,600 AB lineage cells. Colors are labeled by scVelo latent time. (j) The correlation between scVelo latent time and embryo time for all AB lineage cells. (k) PhyloVelo velocity fields for all 29,600 AB lineage cells, estimated by the phylogenetic velocity of MEGs in a single embryo (n=298 cells). Cell colors are labelled by PhyloVelo pseudotime. (l) The correlation between PhyloVelo pseudotime and embryo time for all AB lineage cells. The Spearman correlation coefficients and *P* values are shown.



**Figure 4**

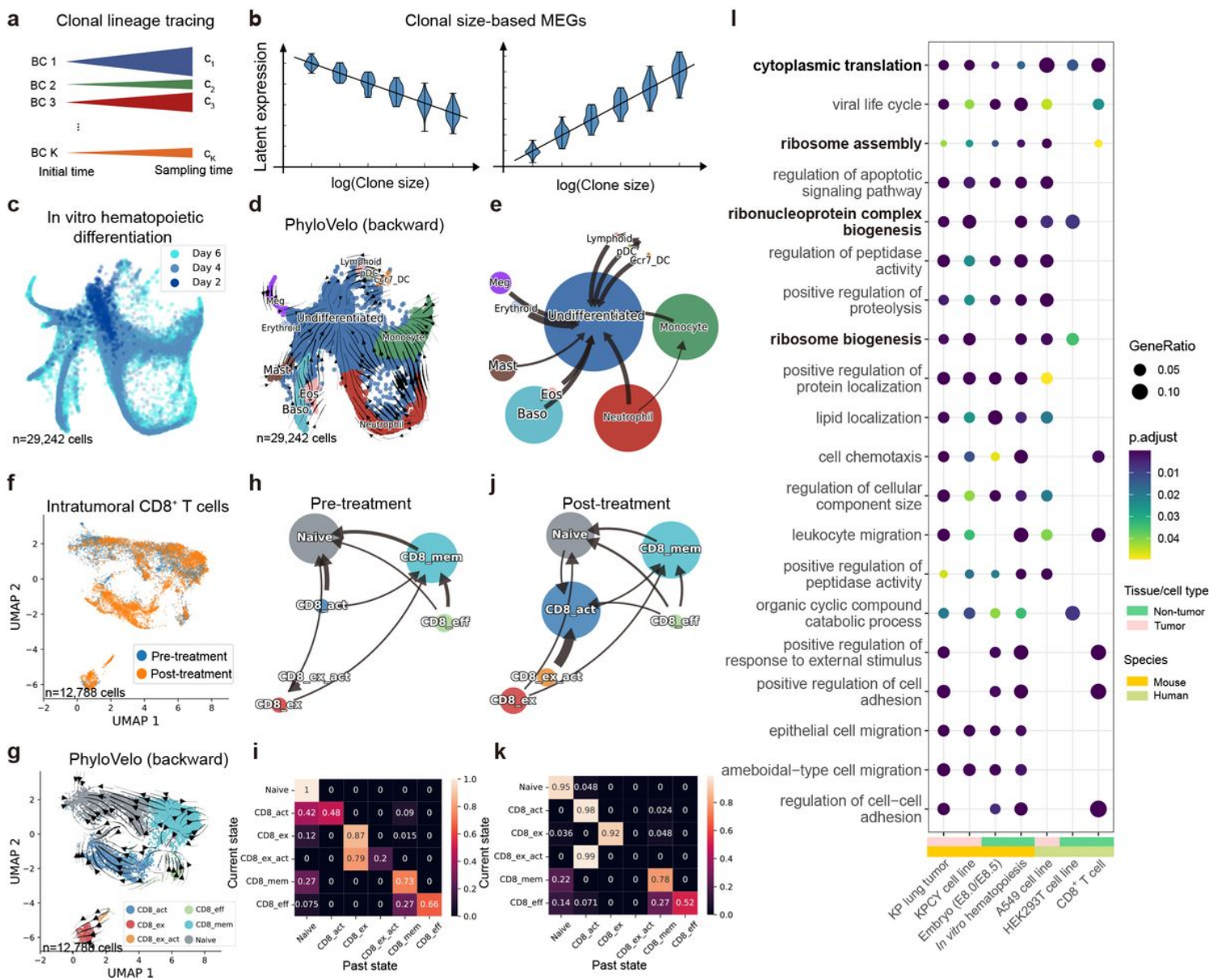
**PhyloVelo reconstructs the cellular trajectory of mouse erythroid maturation.** (a) Phylogenetic tree of the 2,419 erythroid lineage cells (embryo 3, E8.5) in Chan *et al.* dataset<sup>32</sup>. (b-c) RNA velocity fields (scVelo - dynamical mode) and the latent time of mouse erythroid development. (d) Muller plot showing the fractions of four cell types that change over scVelo latent time. (e-f) PhyloVelo velocity fields and the pseudotime of mouse erythroid development. (g) Muller plot showing the fractions of four cell types that change over PhyloVelo pseudotime. (h) Erroneous estimations of RNA velocity fields on erythroid maturation because of multiple rate kinetics (MURK). Data were from Pijuan-Sala *et al.*<sup>19</sup>. (i) PhyloVelo velocity fields of erythroid maturation for Pijuan-Sala *et al.* dataset while using the MEGs identified from Chan *et al.* dataset. (j) PhyloVelo pseudotime of erythroid maturation in Pijuan-Sala *et al.* dataset. (k) The correlation between PhyloVelo pseudotime and mouse embryo time ( $n=12,324$  cells). The Spearman correlation coefficient and  $P$  value are shown here. Whiskers: minimum and maximum; center lines: median.



**Figure 5**

**PhyloVelo identifies a dedifferentiation trajectory in lung tumor evolution.** (a) Phylogenetic tree of 754 cells from a KP-mouse primary lung tumor, 3726\_NT\_T1, in Yang *et al.* dataset<sup>51</sup>. The scRNA-seq data, cell type annotations, and lineage trees were obtained from the original study. (b) RNA velocity fields (scVelo - dynamical mode). (c) PhyloVelo velocity fields. (d) Fitness signatures of individual cells, as defined by Yang *et al.* (e) CytoTRACE score of individual cells. (f) The correlation between PhyloVelo pseudotime and CytoTRACE scores. The Spearman correlation coefficient and  $P$  value are shown here. (g) CytoTRACE score of single tumor cells from human lung primary sites (tLung and tL/B), pleural fluids (PE), lymph node metastases (mLN), and brain metastases (mBrain), as well as normal tissues from lungs (nLung), as described in Kim *et al.*<sup>54</sup>. Bar, median; box, 25th to 75th percentile (IQR); vertical line, data within 1.5 times the IQR. (h) PhyloVelo velocity fields for all 58,022 single cells from pooled KP primary lung tumors, estimated by the MEGs identified from 3726\_NT\_T1. (i) PhyloVelo velocity fields for the cell types that existed in 3726\_NT\_T1. (j) Cell-type transition graph (backward) based on the transition

rate matrix between any two cell types ( $\mathbf{k}$ ), estimated by Dynamo using PhyloVelo velocity fields as input. The arrows point from the current states to the past states.



**Figure 6**

**PhyloVelo inference with clonal lineage tracing data and MEGs are enriched in ribosome-mediated processes.** (a) Schematic of clonal lineage tracing data where static barcodes identify cells of common ancestry. Clone size, denoted by  $c_k$  for  $k$  clones, represents the number of cells carrying the same unique barcode. (b) Two examples of clonal size-based MEGs whose latent expressions are positively or negatively associated with the logarithm of clone sizes, respectively. Whiskers: minimum and maximum. (c) scRNA-seq data of in vitro hematopoietic differentiation from Weinreb *et al.*<sup>37</sup>, where each cell over the course of 2, 4, and 6 days culture could be traced by one unique barcode. (d) The velocity fields estimated by PhyloVelo. (e) Cell type transition graph (backward) of in vitro hematopoietic differentiation. (f) UMAP of tumor-infiltrating CD8<sup>+</sup> T cells in BCC samples pre- and post-PD-1 blockade, colored by anti-PD-1

treatment status. Data were from Yost *et al.*<sup>57</sup> (g) The velocity fields estimated by PhyloVelo. (h-i) Cell-type transition graph and transition matrix (backward) at pre-treatment. (j-k) Cell-type transition graph and transition matrix (backward) at post-treatment. CD8\_act: CD8+ activated T cells; CD8\_ex: CD8+ exhausted T cells; CD8\_ex\_act: CD8+ exhausted/activated T cells; CD8\_eff: CD8+ effector T cells; CD8\_mem: CD8+ memory T cells. (l) Gene ontology (GO) enrichment of MEGs identified across tissues and organisms. The top and most commonly shared 20 biological processes are shown. Ribosome-mediated processes are highlighted.

## Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [ReportingSummary0622.pdf](#)
- [SupplementaryInformation0622.pdf](#)
- [SupplementaryTable1.xlsx](#)
- [ExtendedDataFig.1.jpg](#)
- [ExtendedDataFig.2.jpg](#)
- [ExtendedDataFig.3.jpg](#)
- [ExtendedDataFig.4.jpg](#)
- [ExtendedDataFig.5.jpg](#)
- [ExtendedDataFig.6.jpg](#)
- [ExtendedDataFig.7.jpg](#)
- [ExtendedDataFig.8.jpg](#)
- [ExtendedDataFig.9.jpg](#)
- [ExtendedDataFig.10.jpg](#)