

1 VIA: Generalized and scalable trajectory inference in single-cell omics data

2
3 Shobana V. Stassen¹, Gwinky G. K. Yip¹, Kenneth K. Y. Wong^{1,3}, Joshua W. K. Ho^{2,4} and Kevin K. Tsia^{1,3}

4 ¹Department of Electrical & Electronic Engineering, The University of Hong Kong, Pokfulam Road, Hong Kong

5 ²School of Biomedical Sciences, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Pokfulam, Hong Kong

6 ³Advanced Biomedical Instrumentation Centre, Hong Kong Science Park, Shatin, New Territories, Hong Kong

7 ⁴Laboratory of Data Discovery for Health, Hong Kong Science Park, Shatin, New Territories, Hong Kong

8 Abstract

9 Inferring cellular trajectories using a variety of omic data is a critical task in single-cell data science.
10 However, accurate prediction of cell fates, and thereby biologically meaningful discovery, is challenged
11 by the sheer size of single-cell data, the diversity of omic data types, and the complexity of their
12 topologies. We present VIA, a scalable trajectory inference algorithm that overcomes these limitations by
13 using lazy-teleporting random walks to accurately reconstruct complex cellular trajectories beyond
14 tree-like pathways (e.g. cyclic or disconnected structures). We show that VIA robustly and efficiently
15 unravels the fine-grained sub-trajectories in a 1.3-million-cell transcriptomic mouse atlas without losing
16 the global connectivity at such a high cell count. We further apply VIA to discovering elusive lineages
17 and less populous cell fates missed by other methods across a variety of data types, including single-cell
18 proteomic, epigenomic, multi-omics datasets, and a new in-house single-cell morphological dataset.

19 Background

20 Single-cell omics data captures snapshots of cells that catalog cell types and molecular states with high
21 precision. These high-content readouts can be harnessed to model evolving cellular heterogeneity and
22 track dynamical changes of cell fates in tissue, tumour, and cell population. However, current
23 computational methods face four critical challenges. First, it remains difficult to accurately reconstruct
24 high-resolution cell trajectories and detect the pertinent cell fates and lineages without relying on prior
25 knowledge of input parameter settings. This is a foundational but unmet attribute of trajectory inference
26 (TI) that could make lineage prediction less biased towards input parameters, and thus minimize the
27 confounding factors that impact the underlying hypothesis testing. However, even the few algorithms
28 which automate cell fate detection (e.g., SlingShot¹, Palantir² and Monocle3) exhibit low sensitivity to
29 cell fates and are highly susceptible to changes in input parameters. Second, current trajectory inference
30 (TI) methods predominantly work well on tree-like trajectories (e.g. Slingshot), but lack the
31 generalisability to infer disconnected, cyclic or hybrid topologies without imposing restrictions on
32 transitions and causality⁴. This attribute is crucial in enabling unbiased discovery of complex trajectories
33 which are commonly not well known a priori, especially given the increasing diversity of single-cell omic
34 datasets. Third, the growing scale of single-cell data, notably cell atlases of whole organisms^{6,7}, embryos^{8,9}
35 and human organs¹⁰, exceeds the existing TI capacity, not just in runtime and memory, but in preserving
36 both the fine-grain resolution of the embedded trajectories and the global connectivity among them. Very
37 often, such global information is lost in current TI methods after extensive dimension reduction or
38 subsampling. Fourth, fueling the advance in single-cell technologies is the ongoing pursuit to understand
39 cellular heterogeneity from a broader perspective beyond transcriptomics. A notable example is the
40 emergence of single-cell imaging technologies that now allow information-rich profiling of

41 morphological and biophysical phenotypes of single-cells and thus offer novel mechanistic cues to
42 cellular functions that cannot be solely inferred by proteomic or sequencing data (e.g. in cancer⁵⁹,
43 ageing⁶⁰, drug responses⁶¹). However, the applicability of TI to a broader spectrum of single-cell data has
44 yet to be fully exploited.

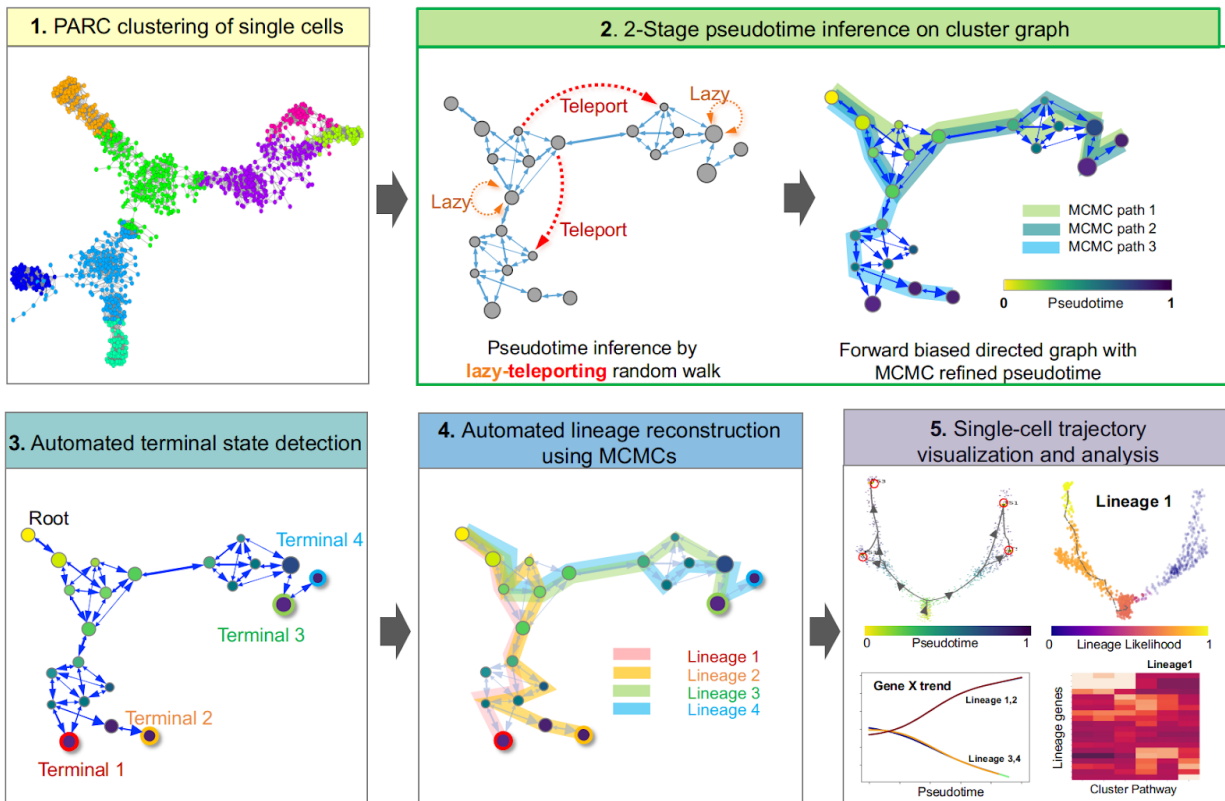
45

46 To overcome these recurring challenges, we present VIA, a graph-based TI algorithm that uses a new
47 strategy to compute pseudotime, and reconstruct cell lineages based on lazy-teleporting random walks
48 integrated with Markov chain Monte Carlo (MCMC) refinement (**Fig. 1**). VIA relaxes common
49 constraints on traversing the graph, and thus allows capture of cellular trajectories not only in
50 multi-furcations and trees, but also in disconnected and cyclic topologies. The lazy-teleporting MCMC
51 characteristics also make VIA robust to a wide range of pre-processing and input algorithmic parameters,
52 and allow VIA to consistently identify pertinent lineages that remain elusive or even lost in other
53 top-performing and popular TI algorithms, e.g. PAGA²⁸, Palintir, SlingShot, Monocle3 and CellRank¹³.
54 We validate the performance of VIA and thus its ability to offer better interpretation of the underlying
55 biology across a variety of transcriptomic, epigenomic and integrated multi-omic datasets (seven
56 biological datasets with a further two datasets presented in **Supplementary**). Notably, we show in
57 subsequent sections that VIA accurately detects minor dendritic sub-populations and their characteristic
58 gene expression trends in human hematopoiesis; automatically identifies pancreatic islets including rare
59 delta cells; and recovers endothelial and cardiomyocyte bifurcation in integrated data sets of single-cell
60 RNA-sequencing (scRNA-seq) and single-cell sequencing assay for transposase-accessible chromatin
61 (scATAC-seq).

62

63 Another defining attribute of VIA is its resilience in handling the wide disparity in single-cell data size,
64 structure and dimensionality across modalities. Specifically, VIA is highly scalable with respect to
65 number of cells (10^2 to $>10^6$ cells) and features, without requiring extensive dimensionality reduction or
66 subsampling which compromise global information. We showcase this scalability in analyzing the
67 fine-grained developmental sub-trajectories in the 1.3-million-cell mouse organogenesis atlas in terms of
68 fast runtime and preservation of global cell-type connectivity, which is otherwise lost in existing TI
69 methods. We also show that VIA is robust against the dimensionality drop (down to 10's - 100's
70 dimensions) in mass cytometry (proteomics) and imaging cytometry (morphological) data. For instance,
71 VIA consistently reconstructs the pseudotime that recapitulates murine embryonic stem cells (ESCs)
72 differentiation toward mesoderm cells in CyTOF data, where the lazy-teleporting MCMCs contribute to
73 the high accuracy of inference. Lastly, we hypothesize that VIA can also be applied to imaging cytometry
74 for gaining a mechanistic biophysical understanding of cellular progress. To this end, we profiled the
75 biophysical and morphological phenotypes of single-cell live breast cancer cells with our recently
76 developed high-throughput imaging flow cytometer, called FACED³³. Validated with the in-situ
77 fluorescence image capture, we found that VIA reliably reconstructs the continuous cell-cycle
78 progressions from G1-S-G2/M phase, and reveals subtle changes in cell mass accumulation.

Fig.1: VIA algorithm workflow



79

80 **Figure 1. General workflow of VIA algorithm. Step 1:** Single-cell level graph is clustered such that each node
 81 represents a cluster of single cells (computed by our clustering algorithm PARC¹¹). The resulting cluster graph forms
 82 the basis for subsequent random walks. **Step 2:** 2-stage pseudotime computation: (i) The pseudotime (relative to a
 83 user defined start cell) is first computed by the expected hitting time for a lazy-teleporting random walk along an
 84 undirected graph. At each step, the walk (with small probability) can remain (orange arrows) or teleport (red arrows)
 85 to any other state. (ii) Edges are then forward biased based on the expected hitting time (See forward biased edges
 86 illustrated as the imbalance of double-arrowhead size). The pseudotime is further refined on the directed graph by
 87 running Markov chain Monte Carlo (MCMC) simulations (See 3 highlighted paths starting at root). **Step 3:** Consensus
 88 vote on terminal states based on vertex connectivity properties of the directed graph. **Step 4:** lineage likelihoods
 89 computed as the visitation frequency under lazy-teleporting MCMC simulations. **Step 5:** visualization that combines
 90 network topology and single-cell level pseudotime/lineage probability properties onto an embedding using GAMs, as
 91 well as unsupervised downstream analysis (e.g. gene expression trend along pseudotime for each lineage).

92 Results

93 Algorithm

94 VIA first represents the single-cell data as a cluster graph (i.e. each node is a cluster of single cells),
 95 computed by our recently developed data-driven community-detection algorithm, PARC, which allows
 96 scalable clustering whilst preserving global properties of the topology needed for accurate TI¹¹ (**Step 1 in**
 97 **Fig. 1**). The cell fates and their lineage pathways are then computed by a two-stage probabilistic method,
 98 which is the key algorithmic contribution of this work (**Step 2 in Fig. 1**, see **Methods** for detailed
 99 explanation). In the first stage of Step 2, VIA models the cellular process as a modified random walk that
 100 allows degrees of *laziness* (remaining at a node/state) and *teleportation* (jumping to any other node/state)
 101 with pre-defined probabilities. The pseudotime, and thus the graph directionality, can be computed based

102 on the theoretical hitting times of nodes (See the theory and derivation in **Methods and Supplementary**
103 **Note 2**). The lazy-teleporting behavior prevents the expected hitting time from converging to a local
104 distribution in the graph as otherwise occurs in regular random walks, especially when the sample size
105 grows¹². More specifically, the laziness and teleportation factors regulate the weights given to each
106 eigenvector-value pair in the expected hitting time formulation such that the stationary distribution (given
107 by the local-node degree-properties in regular walks) does not overwhelm the global information
108 provided by other ‘eigen-pairs’. Moreover, the computation does not require subsetting the first k
109 eigenvectors (bypassing the need for the user to select a suitable threshold or subset of eigenvectors) since
110 the dimensionality is not on the order of number of cells, but is equal to the number of clusters. Hence all
111 eigenvalue-eigenvector pairs can be incorporated without causing a bottleneck in runtime. Consequently
112 in VIA, the modified walk on a cluster-graph not only enables scalable pseudotime computation for large
113 datasets in terms of runtime, but also preserves information about the global neighborhood relationships
114 within the graph. In the second stage of Step 2, VIA infers the directionality of the graph by biasing the
115 edge-weights with the initial pseudotime computations, and refines the pseudotime through
116 lazy-teleporting MCMC simulations on the forward biased graph.

117

118 Next (**Step 3 in Fig. 1**), the MCMC-refined graph-edges of the lazy-teleporting random walk enable
119 accurate predictions of terminal cell fates through a consensus vote of various vertex connectivity
120 properties derived from the directed graph. The cell fate predictions obtained using this approach are
121 more accurate and robust to changes in input data and parameters compared to other TI methods (**Fig.2**
122 simulated complex topologies **and Fig. S1** summary of lineage detection accuracy for all benchmarked
123 real datasets). Trajectories towards identified terminal states are then resolved using lazy-teleporting
124 MCMC simulations (**Step 4 in Fig. 1**). Together, these four steps facilitate holistic topological
125 visualization of TI on the single-cell level (e.g. using UMAP or PHATE^{14,15}) and other data-driven
126 downstream analyses such as recovering gene expression trends (**Methods**). (**Step 5 in Fig.1**).

127 VIA accurately captures complex topologies obscured in other TI methods

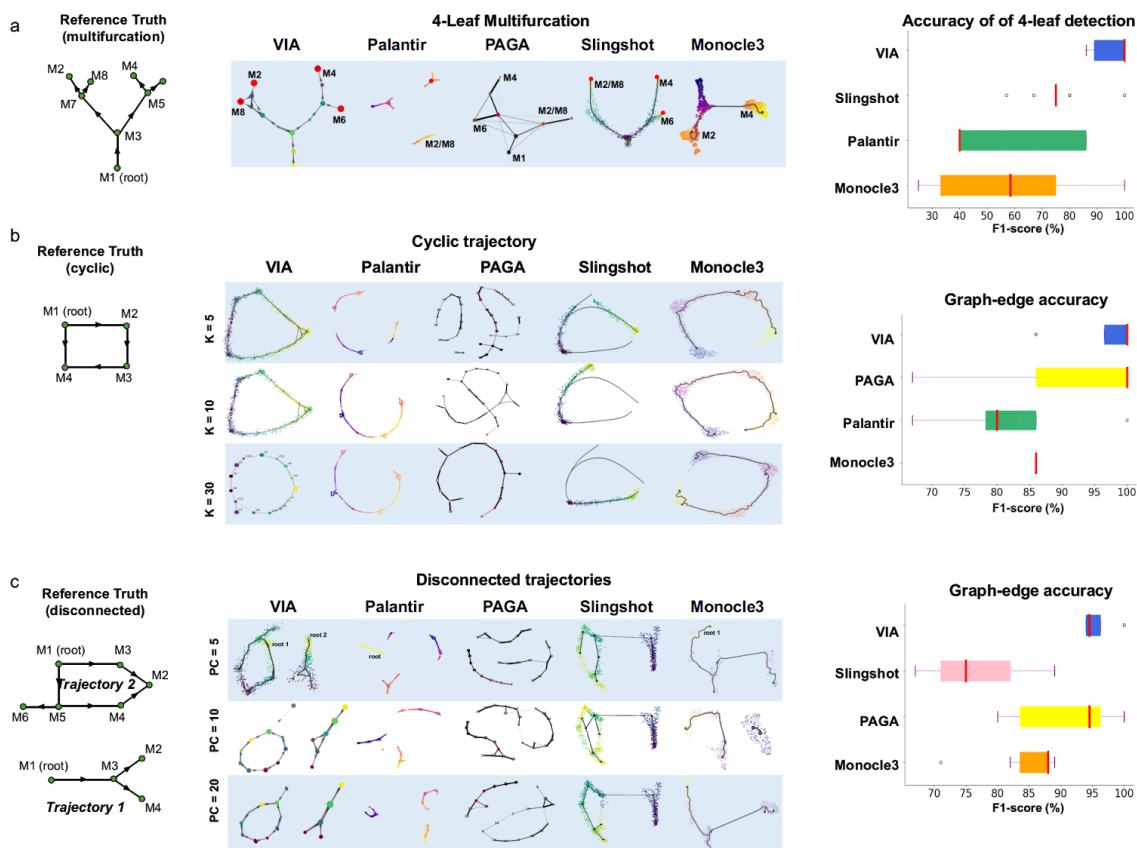
128 We first generate and analyze simulated datasets (see **Methods**) to demonstrate that VIA’s probabilistic
129 approach to graph-traversal allows it to infer cell fates when the underlying data spans combinations of
130 multifurcating trees and cyclic/disconnected topologies - topologies and lineages often obscured in
131 existing TI methods. In VIA, the relaxation of edge constraints in computing lineage pathways and
132 pseudotime enables accurate detection of cell fates and complex trajectories by avoiding prematurely
133 imposing constraints on node-to-node mobility. Other methods resort to constraints such as reducing the
134 graph to a tree, imposing unidirectionality by thresholding edges based on pseudotime directionality,
135 removing outgoing edges from terminal states^{13,2} and computing shortest paths for pseudotime^{2,1}.

136

137 In a 4-leaf multifurcation topology (**Fig. 2a**), VIA accurately captures the two cascading bifurcations
138 which lead to 4 leaf nodes. In particular, VIA detects the elusive ‘M2’ terminal state whereas other
139 methods (Palantir, PAGA, Slingshot and Monocle3) merge it with the ‘M8’ lineage. Monocle3 typically
140 only captures a single bifurcation and thus merges the pairs of leaves that otherwise arise from the second
141 layer of bifurcation (**Fig. 2a**). Even for the fairly simple cyclic topology (**Fig. 2b**), other methods tend to
142 fragment the structure to varying degrees depending on the parameter choice whereas VIA consistently

143 preserves the global cyclic structure. This is not to say VIA is invariant to parameter choice, but rather
 144 that VIA predictably modulates the graph resolution across a wide range of K without disrupting the
 145 underlying global topology (see the increase in the number of nodes in K=30 versus K=5 in **Fig. 2b**). This
 146 characteristic is important for robustly analyzing multiple levels of resolution in complex graph
 147 topologies, as also shown in our later investigation of the 1.3-million-cell mouse atlas. We quantify
 148 graph-edge accuracy in the cyclic and disconnected datasets by identifying false/true positive/negative
 149 edges relative to the reference truth in order to compute an F1-score. The performance comparison for the
 150 disconnected hybrid topologies (**Fig. 2c**) shows that VIA disentangles the cyclic and bifurcating lineages
 151 and captures the key leaf-states in the bifurcation as well as the ‘tail’ extending from the cyclic topology.
 152 Palantir overly fragments the two trajectories, whereas Monocle3 and Slingshot merge them.
 153

Fig. 2: VIA performance comparison for complex hybrid topologies



154

155 **Figure 2 Performance on complex hybrid topologies (a) Toy Multifurcating:** 1000 ‘cells’ multifurcating to four
 156 terminal states. One of the terminal states (M2) is very close to another terminal state (M8), and thus merged by other
 157 methods (Slingshot, Palantir, PAGA and Monocle3). The F1-scores show prediction accuracy of the 4 terminal states
 158 when the number of Principal Components varies (5-200 input PCs). PAGA does not automatically detect
 159 lineages/cell fates and is thus excluded from the F1-score analysis **(b) Toy Cyclic:** VIA recovers a cyclic network for
 160 a range of K (in KNN). Slingshot does not use a K(NN) parameter and identifies 3 different lineages (top to bottom).
 161 PAGA, Monocle3 and Palantir show linear or fragmented structures, however PAGA’s performance for this dataset
 162 improves for higher KNN as the underlying graph representation becomes more connected. (Right) Graph-edge
 163 accuracy compared to the reference truth for a varying number of K(NN), where true positive edges are those that
 164 connect milestones in the reference graph **(c) Disconnected:** This dataset has two disconnected trajectories (T1 and
 165 T2). T2 is cyclic with an extra branch (M5 to M6) and T1 has a bifurcation at M3. (Right) T1 performance comparison
 166 of graph accuracy across different numbers of input PCs . Palantir is heavily fragmented and hence excluded from
 167 graph-edge accuracy computations. Slingshot, Monocle3 and sometimes PAGA place an edge (false positive)
 168 between T1 and T2 connecting the two trajectories, and the bifurcation is typically merged.

169 VIA reveals rare lineages in epigenomic and transcriptomic landscapes of 170 human hematopoiesis.

171 To assess the performance of VIA on inferring real cellular trajectory, we first considered a range of
172 scRNA-seq datasets, including hematopoiesis^{2,27}, endocrine genesis, B-cell differentiation²⁶ and
173 embryonic stem (ES) cell differentiation in embryoid bodies¹⁵. We present the analyses of CD34+ human
174 hematopoiesis and endocrine differentiation here, whereas the generalizable performance of VIA on other
175 scRNA-seq datasets are presented in **Supplementary Fig. S1, S2 and S7**. We highlight human
176 hematopoiesis as it has been extensively studied not only with scRNA-seq, but also other single-cell
177 omics modalities, notably scATAC-seq. Hence, it allows us to reliably assess lineage identification
178 performance and downstream analyses using VIA.

179

180 First, we show that VIA consistently reveals from the scRNA-seq dataset the typical hierarchical
181 bifurcations during hematopoiesis that result in key committed lineages of hematopoietic stem cells
182 (HSCs) to monocytic, lymphoid, erythroid, classical and plasmacytoid dendritic cell (cDCs and pDCs)
183 lineages and megakaryocytes (**Fig. 3a**). The automated detection of these terminal states in VIA, as
184 quantified by F1-scores on the annotated cells, remains robust to varying the number of neighbors in the
185 KNN graph, and the number of principal components (PCs) (**Fig. 3c**). Specifically, VIA's sustained
186 sensitivity to rarer cell types (e.g. DCs and megakaryocytes) can be attributed to a better underlying graph
187 structure where nodes are well delineated by PARC (as rare cell types are well separated by graph pruning
188 in the clustering stage) and edges are not prematurely removed due to restrictions on causality.

189

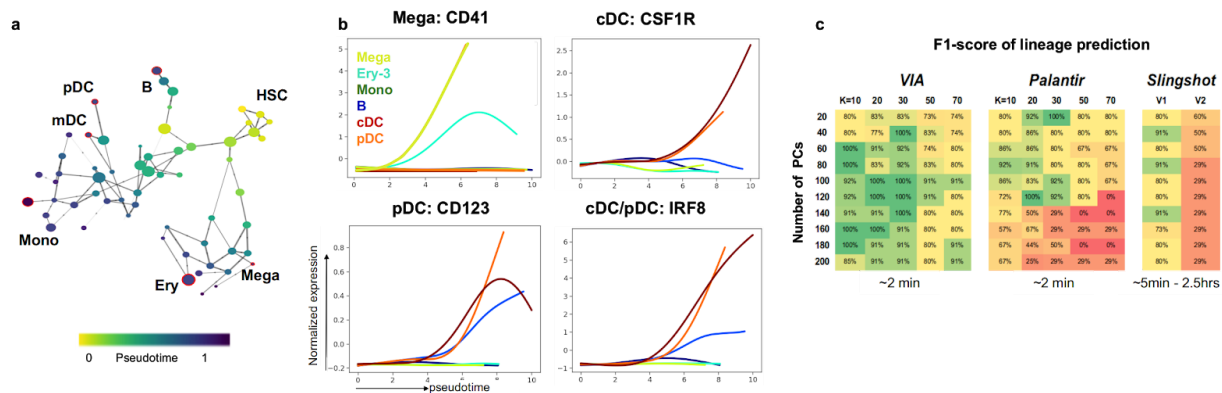
190 In contrast, the sensitivity of Palantir and Slingshot in detecting rarer lineages drops significantly outside
191 a favourable "sweet spot" of parameters. Slingshot can only recover the major cell populations
192 (monocytes, erythroid and B cells) and confuses the DC populations with the monocytes and the
193 megakaryocytes with the erythroid cells. Palantir can only identify the DCs and megakaryocytes for a
194 handful of parameter options, whereas VIA achieves this goal across a much wider range of parameters
195 (**Fig. 3c**). Since PAGA does not offer automated cell fate prediction or lineage paths, it is not
196 benchmarked on this dataset. To verify that VIA reliably delineates the megakaryocyte, cDC and pDC
197 lineages, we used VIA to automatically plot the lineage specific trends for selected marker genes. We
198 showed that while both DC lineages exhibit elevated *IRF8*, the *CSF1R* is specific to the cDC, and the
199 *CD123* remains elevated for pDCs whereas it is first up-regulated, then down-regulated in cDCs (**Fig.3b**
200 **and Fig. S3-S4**).

201

202 We find that VIA's interpretation of the human scATAC-seq profiles (**Fig. 3d**) mirrors the continuous
203 landscape of scRNA-seq human hematopoietic data (**Fig. 3a**). We use two common preprocessing
204 pipelines^{31,27} (see **Methods**), intended to alleviate challenges posed by the sparsity of scATAC-seq data, to
205 show that VIA consistently predicts the expected hierarchy of lineages furcating from hematopoietic
206 progenitors to their descendants. The graph topology of VIA (colored by pseudotime) captures the
207 progression of multipotent progenitors (MPPs) towards the lymphoid-primed MPPs (LMPP) and the
208 common myeloid progenitors (CMPs) which in turn give rise to the CLP and MEP lineages respectively.
209 The known joint contribution of LMPPs and CMPs towards the GMP lineage is also captured by the VIA

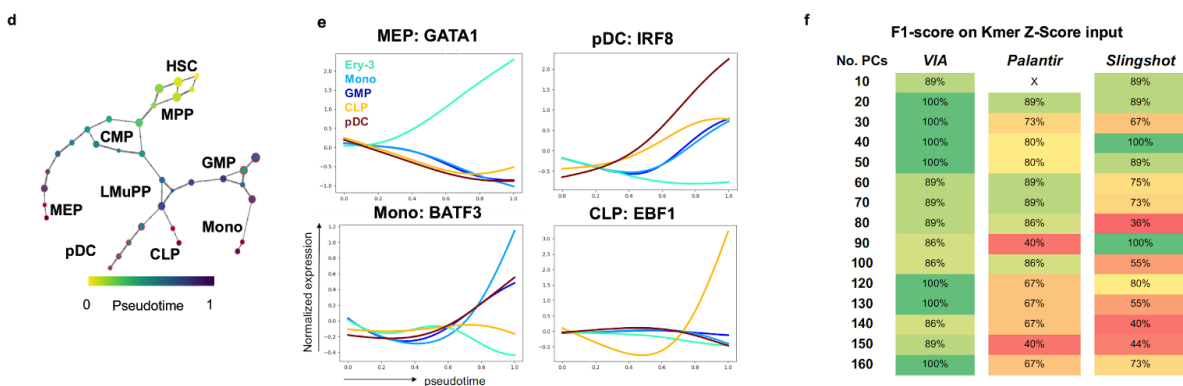
210 graph. We verified the lineages identified by VIA by analyzing the changes in the accessibility of TF
 211 motifs associated with known regulators of the lineage commitments, e.g. *GATA1* (erythroid), *CEBPD*
 212 (myeloid) and IRF8 (DCs) (**Fig 3e, Supplementary Fig. S5c**). Again, we note that the detection of these
 213 lineages is less straightforward in other methods, which generally face a sharp drop in accuracy of
 214 detecting relevant cell fates as the input number of PCs exceeds ~50PCs (e.g. Palantir often misses the
 215 CLP and monocyte lineages, see **Supplementary Fig. S6** for Palantir's outputs across parameters and
 216 **Fig. 3f** for the corresponding prediction accuracy). We emphasize that VIA's robustness in handling both
 217 of these scRNA-seq and scATAC-seq datasets demonstrates its unique ability to achieve stable prediction
 218 and thus faithful query of the underlying biology without biasing specific sets of input parameters which
 219 nontrivially vary across datasets - as also evident from our series of "stress tests" on VIA's performance
 220 (**Supplementary Fig. S1**).
 221

Fig.3: Detection of elusive cell types and their gene trends in scRNA-seq hematopoiesis



222

scATAC-seq: Human Hematopoiesis



223 **Figure 3 VIA analysis of human hematopoiesis based on scRNA-seq and scATAC-seq¹³ data** (a) VIA graph
 224 colored by inferred pseudotime. Identified terminal state nodes are outlined in red and labeled according to their
 225 representative annotated cell type (b) pseudo-temporal trends of marker genes for key minor populations (see
 226 Supplementary Fig. S3-S5 for gene trends of all lineages) (c) F1-scores for terminal state detection of mDC, pDC,
 227 Mega, Ery, Mono and B cell lineages (d) Graph topology of scATAC-seq hematopoietic data using Buenrostro¹³
 228 pre-processing protocol, nodes colored by inferred pseudotime (e) pseudo-temporal trends of transcription-factor
 229 motifs (f) F1-scores for terminal state detection of MEP, CLP, pDC and Mono lineages for fixed KNN=20 and different
 230 number of PCs. Pre-processed using *k-mer* Z Scores protocol which is a more challenging input as evidenced by the
 231 performance drop for other methods beyond 50PCs. VIA's F1-scores are more robust to choice of number of PCs

232 VIA detects small endocrine Delta lineages and Beta subtypes

233 We also use a scRNA-seq dataset of E15.5 murine pancreatic cells to again examine whether VIA can
234 automatically detect multiple lineages, in particular less populous ones. This data spans all developmental
235 stages from initial endocrine progenitor-precursor (EP) state (low level of *Ngn3*, or *Ngn3^{low}*), to
236 intermediate EP (high level of *Ngn3*, or *Ngn3^{high}*) and *Fev⁺* states, to terminal states of hormone-producing
237 alpha, beta, epsilon and delta cells⁵ (**Fig. 4a**).

238

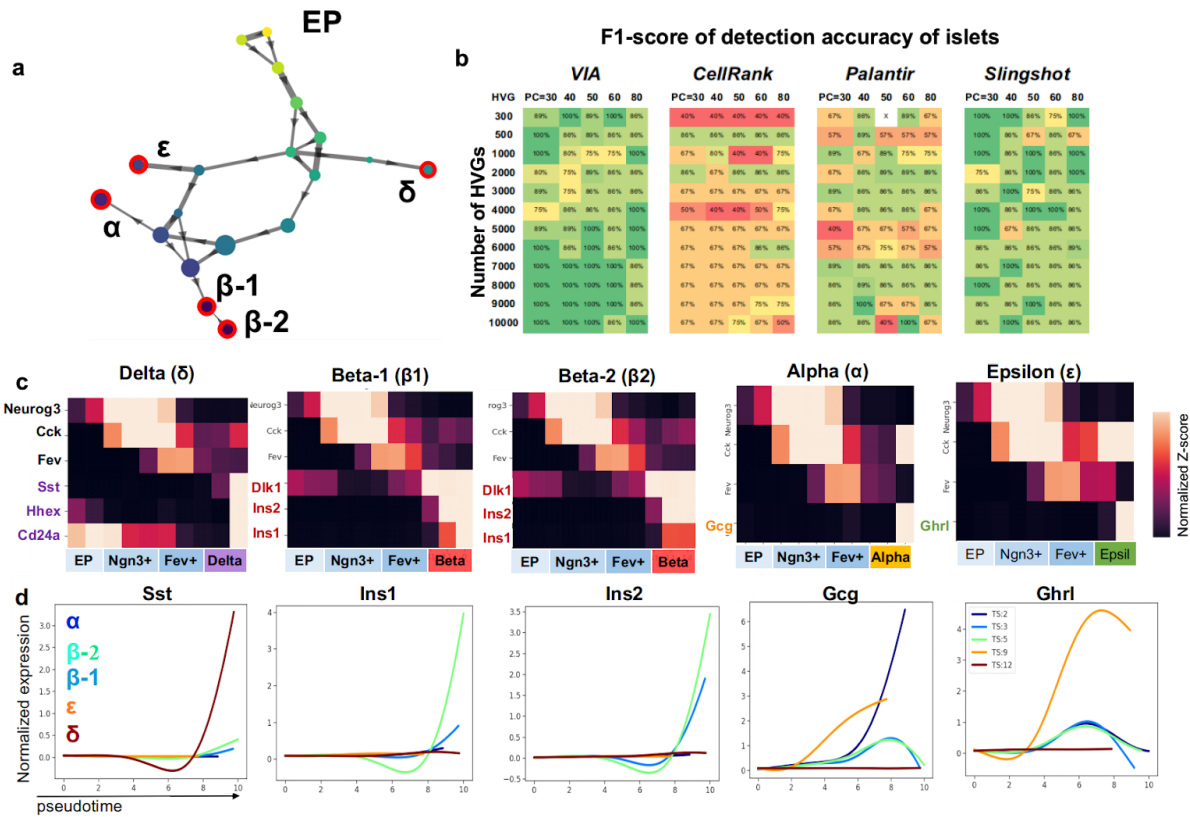
239 A key challenge in analyzing this dataset is the automated detection of the small delta-cell population (a
240 mere 3% of the total population), which otherwise requires manual assignment in CellRank and Palantir
241 (see **Supplementary Fig. S9-S10** for sample outputs at different parameters). In contrast, the
242 well-delineated nodes of the VIA cluster-graph (as a result of sensitive terminal state prediction enabled
243 by the lazy-teleporting MCMC property of VIA) lends itself to automatically detecting this small
244 population of delta cells, together with all other key lineages (alpha, beta and epsilon lineages) (**Fig.**
245 **4a-b**). As evidenced by the corresponding gene-expression trend analysis, VIA detects all of the
246 hormone-producing cells including delta cells which show exclusively elevated *Hhex*, *Sst* and *Cd24a*
247 (**Fig. 4c-d**). To show that this is not a co-incidence of parameter choice, we verify that these populations
248 can be identified for a wide range of chosen highly variable genes (HVGs) and number of PCs (**Fig. 4b**).

249

250 Interestingly, we find that VIA often automatically detects two Beta-cell subpopulations (Beta-1 and
251 Beta-2) (**Fig.4b-e**) that express the common Beta-cell markers, such as *Dlk1*, *Pdx1*, but differ in their
252 expressions of *Ins1* and *Ins2* (**Fig. 4c-d and Fig.S8d**). The pseudotime order within this Beta-cell
253 heterogeneity^{29,30}, undetectable by other TI methods on this dataset, can further be reconciled in the VIA
254 graph where the immature Beta-2 population precedes the mature Beta-1 population. We find that the
255 immature Beta-2 population strongly expresses *Ins2*, and weakly expresses *Ins1*, followed by the mature
256 Beta-1 population which expresses both types of *Ins*³⁰ (**Fig. 4c-d and Fig.S8d** for VIA graphs colored by
257 *Ins1* and *Ins2* further show the difference in *Ins* expression by the two Beta populations).

258

Fig. 4: Detection of endocrine Beta cell sub-types and rare Delta cell population



259

260 **Figure 4. VIA detects small populations in endocrine progenitor cells differentiation.** (a) VIA graph topology
 261 Pancreatic Islets: Colored by VIA pseudotime with detected terminal states shown in red and annotated based on
 262 known cell type as Alpha, Beta-1, Beta-2, Delta and Epsilon lineages where Beta-2 is *Ins1^{low}Ins2⁺* Beta subtype
 263 (**Supplementary Fig. S8** for graph node-level gene expression intensity of *Ins1* and *Ins2*). (b) Prediction accuracy of
 264 the 4 major endocrine cell types when varying the number of HVGs selected in pre-processing, and the number of
 265 PCs. (c) VIA inferred cluster-level pathway shows gene regulation along endocrine progenitor (EP) to *Fev⁺* cells
 266 followed by expression of islet specific genes. (d) shows gene-expression trends along pseudotime for each
 267 pancreatic islet.
 268

269 **VIA recovers *Isl1⁺* cardiac progenitor bifurcation in multi-omics data**

270 We next demonstrate the applicability of VIA in single-cell multi-omics analysis by investigating murine
 271 *Isl1⁺* cardiac progenitor cells (CPC) which are known to bifurcate towards endothelial and
 272 cardiomyocyte fates (**Fig. 5b-e**). VIA consistently uncovers the bifurcating lineages using both single-cell
 273 transcriptomic (scRNA-seq) and chromatin accessibility (scATAC-seq) information²⁰, as well as their data
 274 integration (see **Methods** for data integration using Seurat). Other methods such as Palantir and
 275 Slingshot, that are also applicable to non-transcriptomic data, fail to uncover the two main lineages.
 276

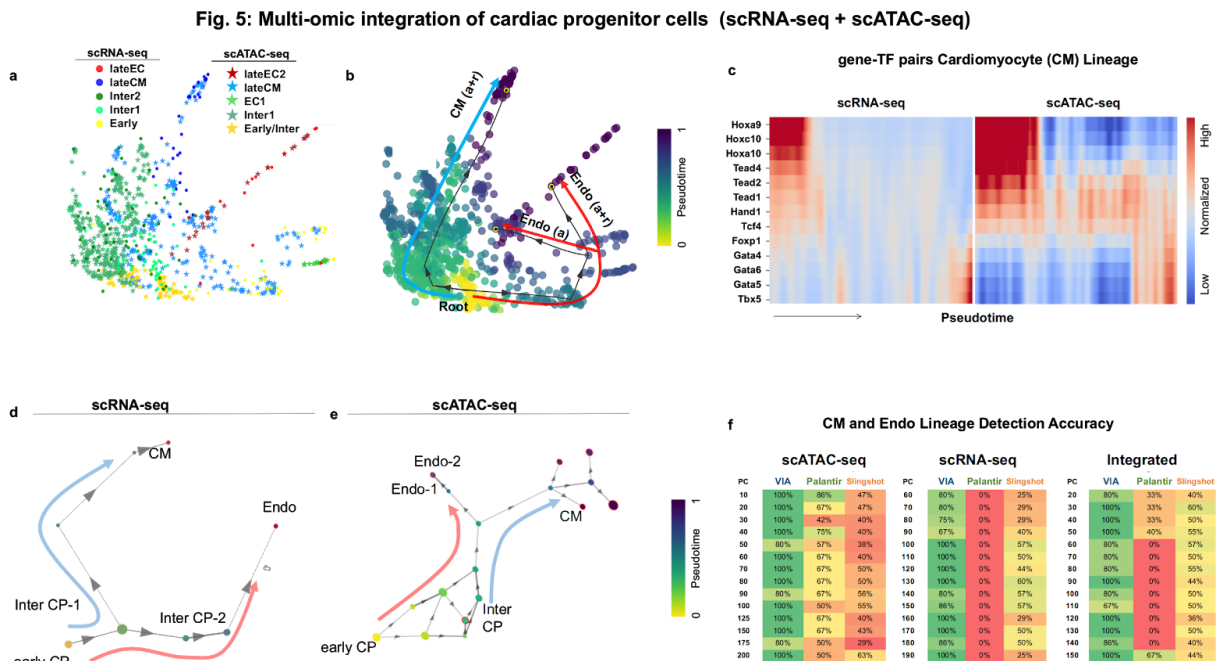
277 Palantir and Slingshot typically only detect the cardiomyocyte lineage (this is exacerbated when the
 278 number of input principal components (PCs) increases), and instead falsely detect several intermediate
 279 and early stages as final cell fates (see **Supplementary Fig. S12-S13** for outputs by Slingshot and
 280 Palantir, **and Fig. 5f** for the corresponding prediction accuracy). PAGA does not offer automated cell fate

281 prediction or lineage paths and is therefore not benchmarked for this dataset. The disparity in trajectory
 282 inference is most evident in the scRNAseq and integrated data where Slingshot and Palantir do not
 283 resolve either of the two cell fates (**Supplementary Fig. S12-S13** for sample outputs corresponding to the
 284 prediction accuracy shown in **Fig.5f**). We hypothesized that lowering the K (number of nearest neighbors)
 285 in Palantir and VIA would be more appropriate given the extremely low cell count (~200 cells) of the
 286 scRNA-seq dataset. Whilst this approach did not alter the outcome for Palantir, we found that VIA is able
 287 to capture the transition from early to intermediate CPCs and finally lineage committed cells.

288

289 More importantly, VIA automatically generates a pseudotemporal ordering of relevant cells (without
 290 requiring manual selection of relevant cells as done in Jia et al.²⁰) along each lineage and their marker-TF
 291 pairs (**Fig. 5f and Supplementary Fig. S11g** for differential gene expression analysis). Hence, VIA can
 292 be used to faithfully interpret relationships between transcription factor dynamics and gene expression in
 293 an unsupervised manner. The highlighted gene and TF pairs in the cardiac lineage show a strong
 294 correlation between expression and accessibility of *Gata* and Homeobox *Hox* genes which are known to
 295 be related to the regulation of cardiomyocyte proliferation^{23,24,25}. VIA's reliable performance against
 296 user-reconfiguration (number of PCs, individual or integrated omic data) suggests its utility in
 297 transferable interpretation between scRNA-seq and scATAC-seq data.

298



299

300 **Figure 5. Multi-omic integrated analysis of scRNA-seq and scATAC-seq cardiac progenitors** (a) scRNA-seq
 301 and scATAC-seq data of Isl1+ Cardiac Progenitors (CPs) integrated using Seurat3 before PHATE. Colored by
 302 annotated cell-type and experimental modality (b) Colored by VIA pseudotime with VIA-inferred trajectory towards
 303 Endothelial and Myocyte lineages projected on top. (c) gene-TF pair expression along VIA inferred pseudotime for
 304 each CM lineage (see **Supplementary Fig.S11** for Top 5 most differentially expressed genes for each VIA node
 305 along each lineage as well as node-level TF motif accessibility) (d) VIA graph for scRNA-seq data only and (e)
 306 scATAC-seq data only. (f) Accuracy of detecting the CM and Endo lineages in the individual and integrated data. This
 307 is challenging for Palantir and Slingshot which either detect several early and intermediate stages or no terminal
 308 states at all (see visual outputs for these methods in **Supplementary Fig.S12-S13**)

309 VIA preserves global connectivity when scaling to millions of cells

310 VIA is designed to be highly scalable and offers automated lineage prediction without extensive
311 dimension reduction or subsampling even at large cell counts. To showcase this, we use VIA to explore
312 the 1.3-million scRNA-seq mouse organogenesis cell atlas (MOCA)⁸. While this dataset is inaccessible to
313 most TI methods from a runtime and memory perspective, VIA can efficiently resolve the underlying
314 developmental heterogeneity, including 9 major trajectories (**Fig. 6a,b**) with a runtime of ~40 minutes,
315 compared to the next fastest method PAGA which has a runtime of 3 hours, and Palantir which takes over
316 4 hours. Other methods like Slingshot and CellRank were deemed infeasible due to extremely long
317 runtimes on much smaller datasets. (**Supplementary Table S3** for a summary of runtimes). Going
318 beyond the computational efficiency, VIA also preserves wider neighborhood information and reveals a
319 globally connected topology of MOCA which is otherwise lost in the Monocle3 analysis which first
320 reduces the input data dimensionality using UMAP.

321

322 The overall cluster graph of VIA consists of three main branches that concur with the known
323 developmental process at early organogenesis.¹⁶ (**Fig. 6a**). It starts from the root stem which has a high
324 concentration of E9.5 early epithelial cells made of multiple sub-trajectories (e.g. epidermis, and
325 foregut/hindgut epithelial cells derived from the ectoderm and endoderm). The stem is connected to two
326 distinct lineages: 1) mesenchymal cells originated from the mesoderm which arises from interactions
327 between the ectoderm and endoderm¹⁷ and 2) neural tube/crest cells derived from neurulation when the
328 ectoderm folds inwards¹.

329

330 The sparsity of early cells (only ~8% are E9.5) and the absence of earlier ancestral cells make it
331 particularly challenging to capture the simultaneous development of trajectories. However, VIA is able to
332 capture the overall pseudotime structure depicting early organogenesis (**Fig. 6b**). For instance, at the
333 junction of the epithelial-to-mesenchymal branch, we find early mesenchymal cells from E9.5-E10.5.
334 Cells from later mesenchymal developmental stages (e.g. myocytes from E12.5- E13.5) reside at the
335 leaves of the branch. Similarly, at the junction of epithelial-to-neural tube, we find dorsal tube neural cells
336 and notochord plate cells which are predominantly from E9.5-E10.5 and more developed neural cells at
337 branch tips (e.g. excitatory and inhibitory neurons appearing at E12.5-E13.5). In contrast, the pseudotime
338 gradient of PAGA's nodes offer little salient information at this scale, with 90% of cells predicted to be in
339 the first 10% of the pseudotime color scale (**see Supplementary Fig. S14c**).

340

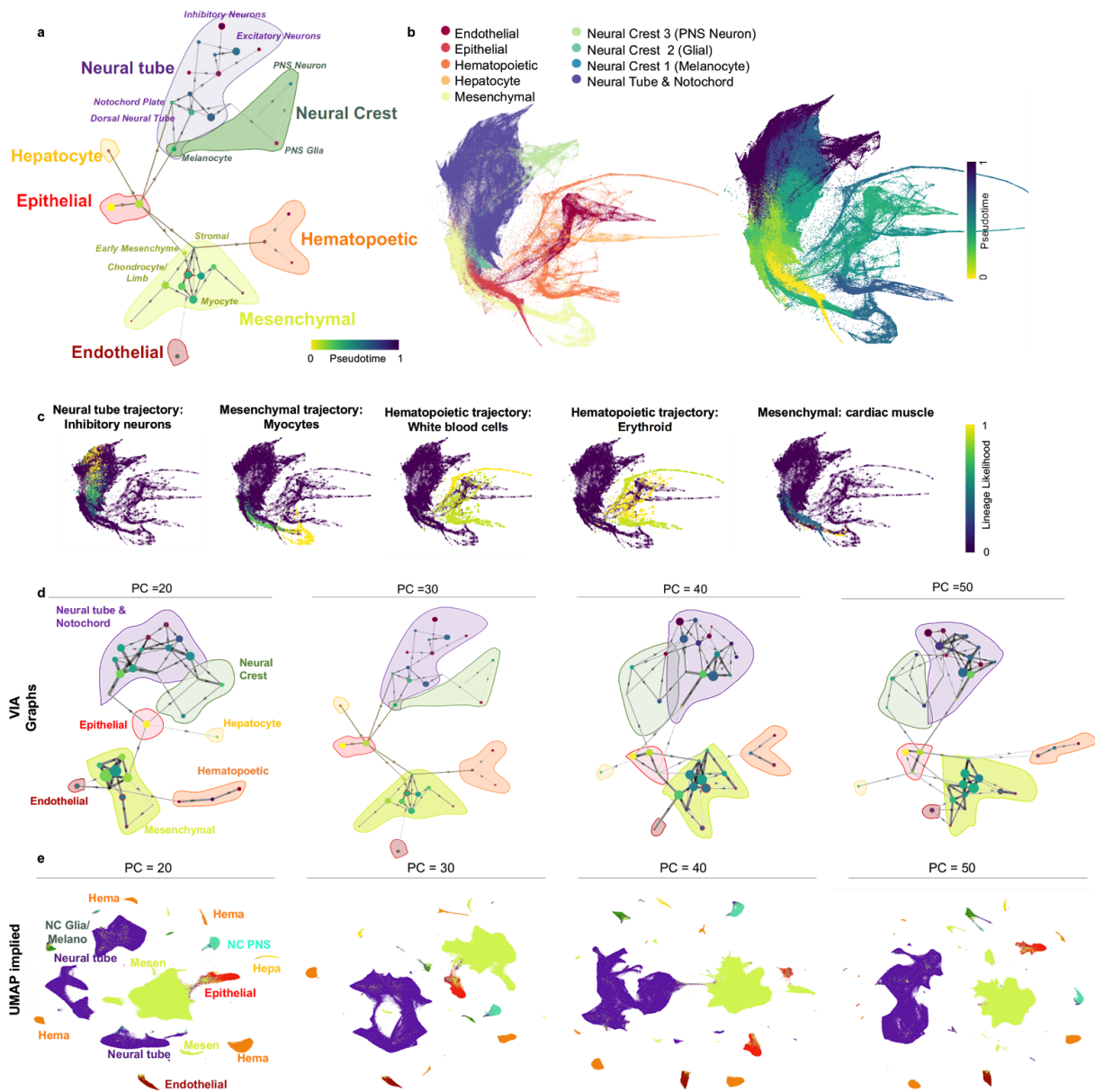
341 VIA also consistently places the other smaller dispersed groups of trajectories (e.g. endothelial,
342 hematopoietic) in biologically relevant neighborhoods (see **Supplementary Notes 3** for a detailed
343 explanation of VIA's structural connections supported by known transitions in organogenesis literature).
344 While VIA's connected topology offers a coarse-grained holistic view, it does not compromise the ability
345 to delineate individual lineage pathways, such as the erythroid and white blood cell lineages within the
346 hematopoietic super group (consistent with annotations made by Cao et al.,⁸) as shown in **Fig. 6c**.

347

348 As such, TI using VIA uniquely preserves both the global and local structures of the data. Whilst
349 manifold-learning methods are often used to extensively reduce dimensionality to mitigate the
350 computational burden of large single-cell datasets, they tend to incur loss of global information and be

351 sensitive to input parameters. VIA is sufficiently scalable to bypass such a step, and therefore retains a
352 higher degree of neighborhood information when mapping large datasets. This is in contrast to
353 Monocle3's⁸ UMAP-reduced inputs that reveal different disconnected super-groups and fluctuating
354 connectivity depending on input parameters. As shown in **Fig. 6d,e** (and **Fig. S14** for varying KNN),
355 methods such as Monocle3 and Slingshot which require on a low dimensional representation (e.g.
356 UMAP) for TI are susceptible to unpredictable changes in the composition of super cell groups, their
357 relative positions and inter-connectivity. For instance, in UMAP, the neural tube group is sometimes
358 shown as a single super group, and other times fragmented across the embedding without context of
359 neighboring groups. Similarly the hematopoietic supergroup is shown as a single, two or even three
360 separate groups dispersed across the embedding landscape (**Fig. 6e**). In contrast, VIA uncovers
361 biologically consistent structures across the same range of parameters. In VIA, the cells belonging to
362 these fine-grained supergroups remain connected and neighborhood relationships are preserved, for
363 instance the neural crest cells (containing Peripheral Nervous System neurons and glial cells) remain
364 adjacent to the neural tube (**Fig. 6f**).

Fig. 6: Large-scale (1.3 million cells) trajectory inference of mouse organogenesis



365

366 **Figure 6 VIA accurately infers global connectivity and sub-trajectories in the 1.3-million scRNA-seq mouse**
 367 **organogenesis cell atlas. (a)** MOCA graph trajectory (nodes colored by pseudotime) and shaded-colored regions
 368 corresponding to major cell groups. Stem branch consists of epithelial cells derived from ectoderm and endoderm,
 369 leading to two main branches: 1) the mesenchymal and 2) the neural tube and neural crest. Other major groups are
 370 placed in the biologically relevant neighborhoods, such as the adjacencies between hepatocyte and epithelial
 371 trajectories; the neural crest and the neural tube; as well as the links between early mesenchyme with both the
 372 hematopoietic cells and the endothelial cells (see Supplementary Note 3) **(b)** Colored by VIA pseudotime. **(c)** Lineage
 373 pathways and probabilities of neuronal, myocyte and WBC lineages (see Fig.S6 for other lineages). **(d)** VIA graph
 374 preserves key relationships across choice of number of principal components whereas **(e)** UMAP embedding is first
 375 step in the TI method Monocle3 and highly susceptible to choice of number of PCs (or K in KNN see Fig.S12-15)

376 VIA's lazy-teleporting MCMCs delineate mesoderm differentiation in mass 377 cytometry data

378 Broad applicability of TI beyond transcriptomic analysis is increasingly critical, but existing methods
379 have limitations contending with the disparity in the data structure (e.g. sparsity and dimensionality)
380 across a variety of single-cell data types and oftentimes are designed with a view to only handling
381 transcriptomic data. To this end, we investigated whether VIA can cope with the significant drop in data
382 dimensionality (10-100), as often presented in flow/mass cytometry data, and still delineate continuous
383 biological processes.

384

385 We applied VIA on a time-series mass cytometry data (28 antibodies, 90K cells) capturing murine
386 embryonic stem cells (ESCs) differentiation toward mesoderm cells³². The mESCs are captured at 12
387 intervals within the first 11 days and hence provide sufficiently granular temporal annotation to allow a
388 correlation assessment of the inferred pseudotimes. We quantified that the pseudotimes computed by VIA
389 shows a Pearson correlation of ~88% with the actual annotated days. We further verified that VIA's
390 performance is critically improved by the lazy-teleporting MCMCs (**Fig. 7d**), without which the
391 correlation drops closer to PAGA's. Palantir suffers from low connectivity of cells between the Day 0-1
392 and the subsequent early stages, and thus results in loss of pseudotime gradient and low correlation to the
393 true annotations.

394

395 More importantly, unlike previous analysis³² of the same data which required chronological labels to
396 visualize the chronological developmental hierarchy, we ran VIA without such supervised adjustments
397 and accurately captured the sequential development. Not only can it achieve faster runtime (running in 2
398 minutes on the full antibody-feature set versus Slingshot which required 6 hours even on a subset of first
399 5 PCs **see Table S3** for more runtime comparisons), VIA detects 3 terminal states corresponding to cells
400 in the final developmental stages of Day 10-11 which are indicated by upregulation of *Pdgfra*, *Cd44* and
401 *Gata4* mesodermal markers (**Fig. 7f**). In contrast, other methods struggle to identify the correct terminal
402 states (e.g. Palantir and Slingshot **Fig. 7e**) and do not depict salient structures (e.g. PAGA) (**Fig. 7e**).

Fig. 7: CyTOF ESC to Mesoderm

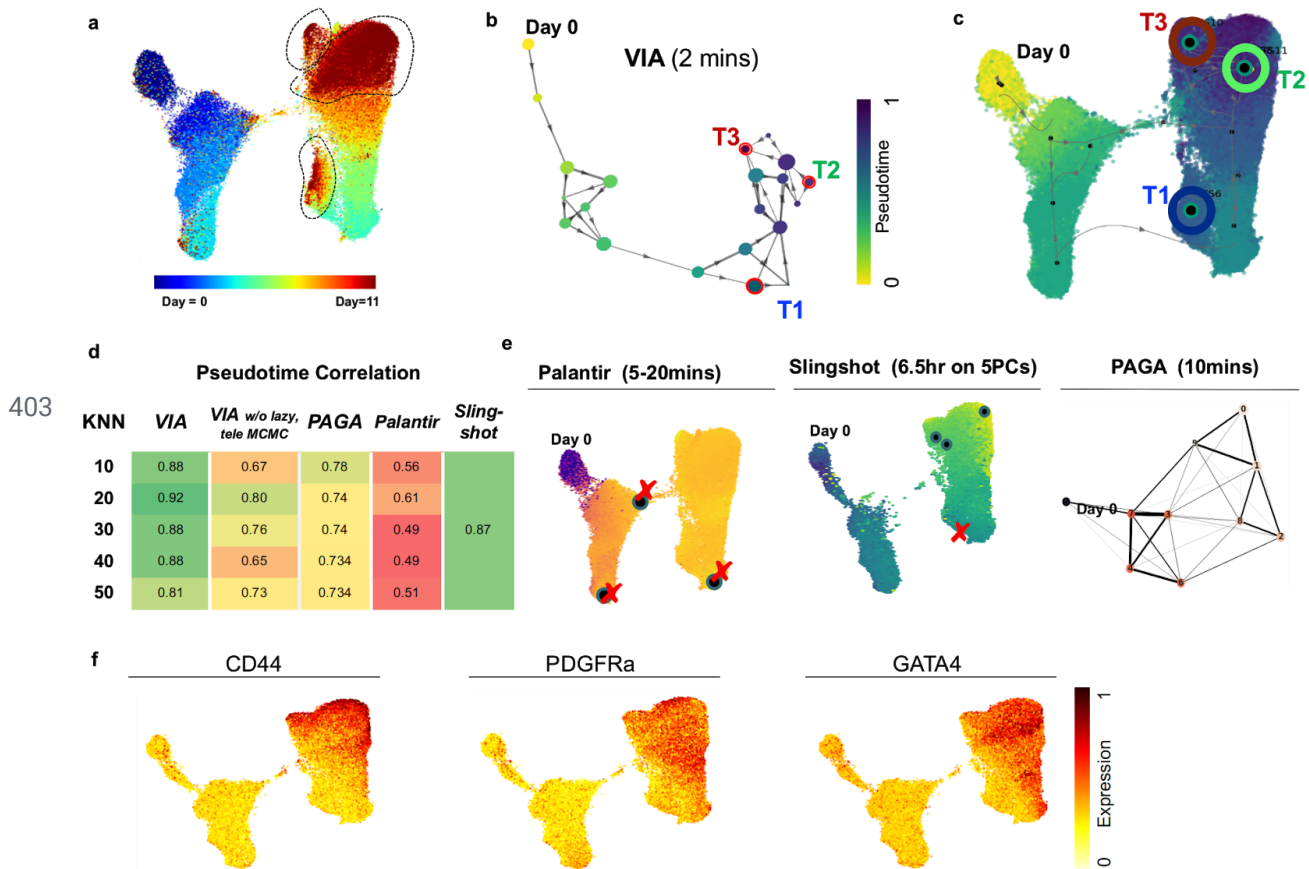


Figure 7 VIA analysis of mESC differentiation toward mesoderm cells from mass cytometry: (a) UMAP plot colored by annotated days 0-11. Three regions of Day 10-11 marked in dotted black lines. (b) VIA cluster-graph colored by pseudotime (c) Terminal states and VIA output projected onto UMAP. Terminal states are located in the areas containing Day 10-11 cells. (d) Comparison of correlation of pseudotime and annotated Days across T1 methods for varying number of K number of nearest neighbors. PAGA and Palantir's pseudotime computation is misguided by the weak link connecting Day 0 cells to other early cells. The effect is that Day 0 cells appear exaggeratedly far, while the remaining early and late cells temporally squeezed. VIA's 2-step pseudotime computation produces a pseudotime scale closer to the annotated dates. (e) Example outputs of Palantir, PAGA and Slingshot with the terminal states (circles) predicted by Slingshot and Palantir. Red 'X' denotes incorrect (false positive) or missing (false negative) terminal state. (f) Gene expression of key mesodermal markers

414 VIA captures morphological trends of live cells in cell cycle progression

415

416 Apart from the omics technologies, optical microscopy is a powerful parallel advance in single-cell
417 analysis for generating the “fingerprint” profiles of cell morphology. Such spatial information is typically
418 obscured in sequencing data, but can effectively underpin the cell states and functions without costly and
419 time-consuming sequencing protocols. However, trajectory predictions based on morphological profiles
420 of single cells have only been scarcely studied until recently, but advancements in high-throughput
421 imaging cytometry are now making large-scale image data generation and related studies feasible. We
422 thus sought to test if VIA can predict biologically relevant progress based on single-cell morphological
423 snapshots captured by our recently developed high-throughput imaging flow cytometer, called FACED³³. -
424 a technology that is at least 100 times faster than state-of-the-art imaging flow cytometry (**Fig. 8a**).

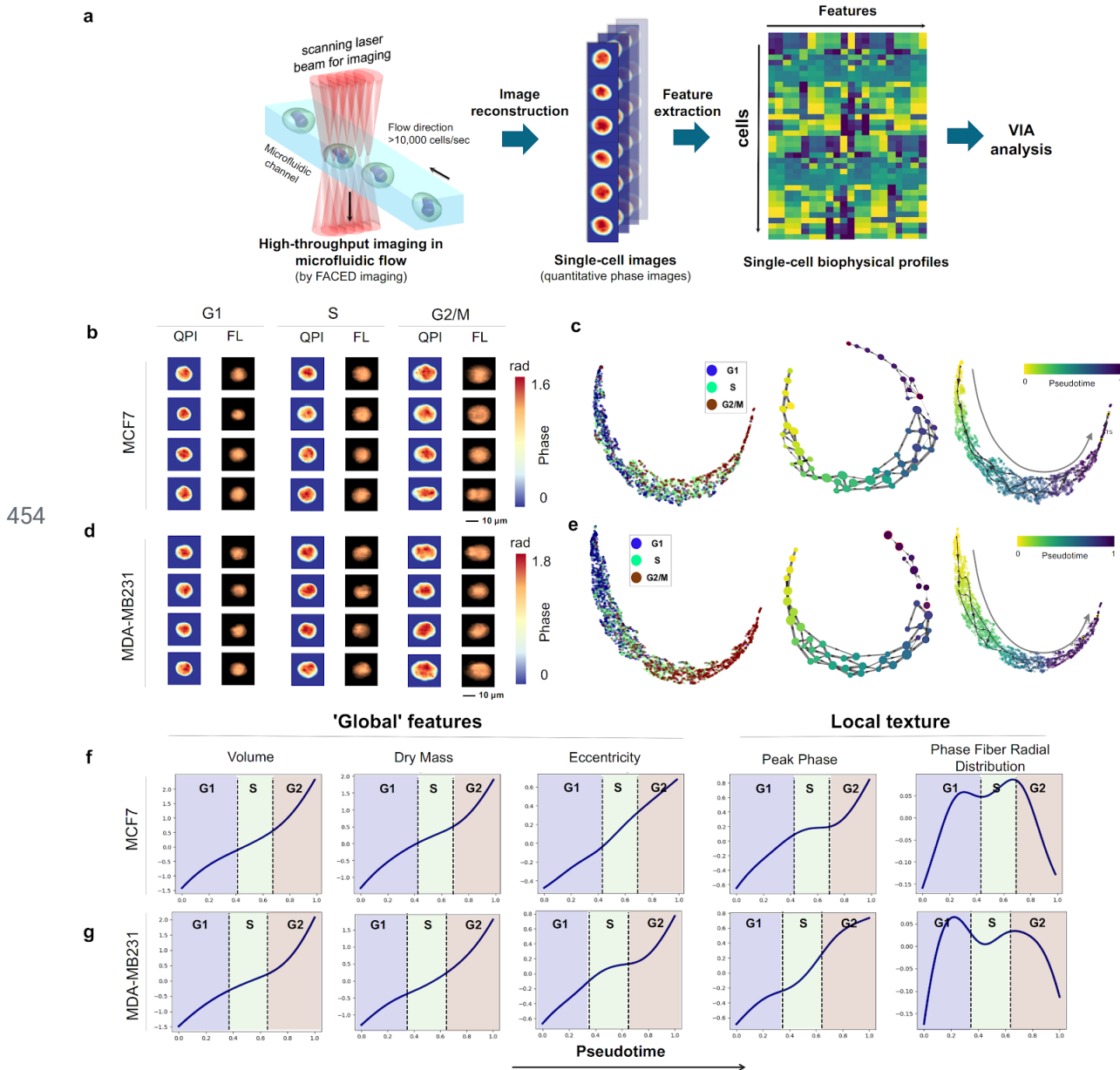
425

426 Our FACED imaging platform captured multiple image contrasts of single cells, including fluorescence
427 (FL), and quantitative phase images (QPI), which measure high-resolution biophysical properties of cells,
428 which are otherwise inaccessible in other methods⁶². Using the QPIs captured by FACED, we first
429 generated spatially-resolved single-cell biophysical profiles of two live breast cancer cell types
430 (MDA-MB231 and MCF7) undergoing cell cycle progressions (38 features including cell shape, size, dry
431 mass density, optical density and their subcellular textures (**see Supplementary Table S4 and Table S5**
432 for definitions of features)). The QPI together with the FL images of individual cells were also used to
433 train a convolutional neural network (CNN)-based regression model for predicting the DNA content. We
434 first validated that there is a high correlation (Pearson’s correlation coefficient $r = 0.72$) between the
435 actual DNA content determined by the FL images and DNA content predicted by the QPI
436 (**Supplementary Fig. S16a**). In addition, the predicted percentages of cells in each cell cycle phases (i.e.
437 G1, S and G2/M) by the biophysical profile are highly consistent with the ground truth defined by the
438 DNA dye (**Supplementary Fig. S16b**). Based on the biophysical profiles as validated by the above tests,
439 VIA reliably reconstructed the continuous cell-cycle progressions from G1-S-G2/M phase of both types
440 of live breast cancer cells (**Methods**)(**Fig. 8b-g**).

441

442 Intriguingly, according to the pseudotime ordered by VIA, not only does it reveal the known cell growth
443 in size and mass³⁴, and general conservation of cell mass density³⁵ (as derived from the FACED images
444 (**Methods**)) throughout the G1/S/G2 phases, but also a slow-down trend during the G1/S transition in
445 both cell types, consistent with the lower protein-accumulation rate during S phase³⁶ (**Fig. 8f-g**). The
446 variation in biophysical textures (e.g. peak phase, and phase fiber radial distribution) along the VIA
447 pseudotime likely relates to known architectural changes of chromosomes and cytoskeletons during the
448 cell cycles (**Fig. 8f-g**). We find that Palantir is very sensitive to the choice of early cells even when
449 choosing from the pool of annotated G1 cells, showing a bifurcating topology unless the early cell is
450 carefully designated based on the diffusion map location of G1 cells (**see Fig. S15** for Palantir and PAGA
451 outputs). The slowdown during the S-phase is also not detected by Palantir’s gene trends. These results
452 further substantiate the growing body of work^{37,38,39,40} on imaging biophysical cytometry for gaining a
453 mechanistic understanding of biological systems, especially when combined with omics analysis⁴¹.

Fig. 8: FACED Imaging Cytometry Cell Cycle



455 **Figure 8 VIA predicts cell cycle progression based on single-cell biophysical morphology** (a) FACED
 456 high-throughput imaging flow cytometry of MDA-MB231 and MCF7 cells, followed by image reconstruction and
 457 biophysical feature extraction. See **Methods** detailed experimental workflow. (b) Randomly sampled quantitative
 458 phase images (QPI) and fluorescence images (FL) of MCF7 cells and (d) MDA-MB231 cells. (c) Single-cell UMAP
 459 embedding colored by the known cell-cycle phase (left), given by DNA-labelled fluorescence images. VIA inferred
 460 cluster-graph topology, nodes colored by pseudotime (mid). UMAP colored by VIA pseudotime for MCF7 (e) VIA
 461 analysis repeated for MDA-MB231 cells. (f) Unsupervised image-feature-trends of global and local biophysical
 462 textures against VIA pseudotime for MCF7 and (g) MDA-MB231 cells (see **Supplementary Table S4 for feature**
 463 **definitions**). Cell cycle pseudotime boundaries are defined here as the intersection of the pseudotime probability
 464 density functions of each cell cycle stage (annotated based on fluorescence intensity).

465 Conclusion

466 With the growing scale and complexity of single-cell datasets, there is an unmet need for accurate cell
467 fate prediction and lineage detection in the complex topologies of interest in biology (not limited to trees).
468 This challenge, broadly faced by the current TI methods, is further compounded by susceptibility to
469 algorithmic parameter changes, limited scalability to large data size; and insufficient generalizability to
470 multi-omic data beyond transcriptomic data. We introduced VIA that alleviates these challenges by fast
471 and scalable construction of cluster-graph of cells, followed by pseudotime, and reconstructing cell
472 lineages based on lazy-teleporting random walks and MCMC simulations. This unique strategy critically
473 relaxes common constraints on graph traversal and causality that impede accurate prediction of elusive
474 lineages and less populous cell fates. We validated the efficacy of these measures in terms of detecting
475 various challenging topologies on simulated data, as well as accurate and robust prediction of cell fates on
476 a variety biological processes (spanning epigenomic, transcriptomic, integrated omic, as well as imaging
477 and mass cytometric data) to show that VIA detects pertinent biological lineages that remain undetected
478 by other methods.

479

480 Notably, VIA distinguished between dendritic subtypes in an scRNA-seq hematopoiesis dataset;
481 identified the rare delta cell islet in pancreatic development, a population requiring manual assignment in
482 other TI methods; and revealed the bifurcation towards cardiomyocyte and endothelial lineage
483 commitment in a multi-omic scATAC-seq and scRNA-seq dataset which proved challenging for other
484 methods. In order to demonstrate that these biological findings are robust to user parameter tuning, we
485 conducted a series of ‘stress tests’ on both simulated and biological data which show that VIA behaves
486 more predictably (allowing controllable degrees of analytical granularity) and accurately than other
487 methods. In other methods, user parameter choice can incur fragmentation or spurious linkages in the
488 modeled topology, and consequently only yield biologically sensible lineages for a narrow sweet spot of
489 parameters (See the summary in **Supplementary Fig. S1** and sample outputs by other methods in
490 **Supplementary Fig. S6, S9, S10, S12 and S13**).

491

492 We also demonstrated on the 1.3 million MOCA dataset that VIA is highly scalable with a runtime of ~40
493 minutes (compared to 3-4 hours on the next fastest method). Importantly, VIA not only recovers the
494 fine-grained sub-trajectories, but also maintains global connectivity between related cell types and thus
495 captures key relationships among lineages in early embryogenesis. It also computes a more salient
496 pseudotime measure supported by lazy-teleporting MCMCs, compared to other methods whose
497 pseudotime scale was distorted at such high cell counts. We also showed that methods which require
498 UMAP (or t-SNE) before parsing MOCA are highly susceptible to user defined input parameters that can
499 significantly and unpredictably fragment the global topology.

500

501 We also assessed whether VIA can be generalized to other single-cell datasets, especially those with
502 significant dimensionality disparity compared to sequencing data. We first applied VIA to the mESC
503 CyTOF dataset and showed that the lazy-teleporting MCMCs strategy in VIA enables it to outperform
504 other methods in correctly correlating the pseudotime of the mesoderm development to the annotated
505 dates. We finally explored the utility of VIA in analyzing emerging image-based single-cell biophysical

506 profile data. We showed that VIA not only successfully identified the progression of G1/S/G2 stages, but
507 also revealed the subtle changes in biophysical-related cellular properties, which are otherwise obscured
508 in other methods. VIA could thus motivate new strategies in single-cell analysis that link cellular
509 biophysical phenotypes and biochemical/biomolecular information - discovering how molecular
510 signatures translate into the emergent cellular biophysical properties, which has already shown effective
511 in studies of cancer, ageing, and drug responses. Overall, VIA offers an advancement to TI methods to
512 robustly study a diverse range of single-cell data. Together with its scalable computation and efficient
513 runtime, VIA could be useful for multifaceted exploratory analysis to uncover novel biological processes,
514 potentially those deviated from the healthy trajectories

515 **Methods**

516 **VIA Algorithm**

517 VIA applies a scalable probabilistic method to infer cell state dynamics and differentiation hierarchies by
518 organizing cells into trajectories along a pseudotime axis in a nearest-neighbor graph which is the basis
519 for subsequent random walks. Single cells are represented by graph nodes that are connected based on
520 their feature similarity, e.g. gene expression, transcription factor accessibility motif, protein expression or
521 morphological features of cell images. A typical routine in VIA mainly consists of four steps:

522

523 **1. Accelerated and scalable cluster-graph construction.** VIA first represents the single-cell data in a
524 k-nearest-neighbor (KNN) graph where each node is a cluster of single cells. The clusters are
525 computed by our recently developed clustering algorithm, PARC¹¹. In brief, PARC is built on
526 hierarchical navigable small world (HNSW⁵⁸) accelerated KNN graph construction and a fast
527 community-detection algorithm (Leiden method⁴²), which is further refined by data-driven pruning.
528 The combination of these steps enables PARC to outperform other clustering algorithms in
529 computational run-time, scalability in data size and dimension (without relying on subsampling of
530 large-scale, high-dimensional single-cell data (>1 million cells)), and sensitivity of rare-cell detection.
531 We employ the cluster-level topology, instead of a single-cell-level graph, for TI as it provides a
532 coarser but clearer view of the key linkages and pathways of the underlying cell dynamics without
533 imposing constraints on the graph edges. Together with the strength of PARC in clustering scalability
534 and sensitivity, this step critically allows VIA to faithfully reveal complex topologies namely cyclic,
535 disconnected and multifurcating trajectories (**Fig. 2**).

536

537 **2. Probabilistic pseudotime computation.** The trajectories are then modeled in VIA as (i)
538 lazy-teleporting random walk paths along which the pseudotime is computed and further refined by
539 (ii) MCMC simulations. The root is a single cell chosen by the user. These two sub-steps are detailed
540 as follows:

541 (i) *Lazy-teleporting random walk.* We first compute the pseudotime as the expected hitting time
542 of a *lazy-teleporting* random walk on an undirected cluster-graph generated in Step 1. The
543 lazy-teleporting nature of this random walk ensures that as the sample size grows, the expected
544 hitting time of each node does not converge to the stationary probability given by local node
545 properties, but instead continues to incorporate the wider global neighborhood information¹².

546 Here we highlight the derivation of the closed form expression of the hitting time of this modified
 547 random walk with a detailed derivation in **Supplementary Note 2**.
 548

549 The cluster graph constructed in VIA is defined as a weighted connected graph $\mathbf{G}(V, E, W)$ with
 550 a vertex set V of n vertices (or nodes), i.e. $V = \{v_1, \dots, v_n\}$ and an edge set E , i.e. a set of
 551 ordered pairs of distinct nodes. W is an $n \times n$ weight matrix that describes a set of edge weights
 552 between node i and j , $w_{ij} \geq 0$ are assigned to the edges (v_i, v_j) . For an undirected graph,
 553 $w_{ij} = w_{ji}$, the $n \times n$ probability transition matrix, P , of a standard random walk on G is given by

$$554 \quad P = D^{-1}W \quad (1)$$

555 where D is the $n \times n$ degree matrix, which is a diagonal matrix of the weighted sum of the degree
 556 of each node, i.e. the matrix elements are expressed as

$$557 \quad d_{ij} = \begin{cases} \sum_k w_{ik} & , i = j \\ 0 & , i \neq j \end{cases} \quad (2)$$

558 where k are the neighbouring nodes connected to node i . Hence, d_{ii} (which can be reduced as d_i)
 559 is the degree of node i . We next consider a *lazy* random walk, defined as Z , with probability
 560 $(1 - x)$ of being lazy (where $0 < x < 1$), i.e. staying at the same node, then

$$561 \quad Z = xP + (1 - x)I \quad (3)$$

562 where I is the identity matrix. When teleportation occurs with a probability $(1 - \alpha)$, the modified
 563 lazy-teleporting random walk Z' can be written as follows, where J is an $n \times n$ matrix of ones.

$$564 \quad Z' = \alpha Z + (1 - \alpha) \frac{1}{n} J \quad (4)$$

565 Here we adapt the concept of personalized PageRank vector, originally used for recording (or
 566 *ranking*) personal preferences of a web-surfer toward particular website pages⁴³, to *rank* the
 567 importance of other nodes (clusters of cells) to a given node, depending on the similarities among
 568 nodes (related to P in the graph), and the lazy-teleporting random walk characteristics in the
 569 graph (set by probabilities of teleporting and being lazy). Based on this concept, one could model
 570 the likelihood to transit from one node (cluster of cells) to another, and thus construct the
 571 pseudotime based on the hitting time, which is a parameter describing the expected number of
 572 steps it takes for a random walk that starts at node i and visit node j for the first time. Consider
 573 the teleporting probability of $(1 - \alpha)$ and a seed vector s specifying the initial probability

574 distribution across the n nodes (such that $\sum_m s_m = 1$, where s_m is the probability of starting at
 575 node m) the personalized PageRank vector $pr_\alpha(s)$ (which is defined as a column vector) is the
 576 unique solution to⁵⁶

$$577 \quad pr_\alpha(s)^T = \alpha pr_\alpha(s)^T Z + (1 - \alpha)s^T. \quad (5)$$

578 Substituting Z (Eq. (3)) into Eq. (5), we can express the personalized PageRank vector $pr_\alpha(s)$ in
 579 terms of the inverse of the β -normalized Laplacian, $R_{\beta, NL}$ of the modified random walk
 580 (**Supplementary Note 2**), i.e.

$$581 \quad pr_\alpha(s)^T = \beta s^T D^{-0.5} R_{\beta, NL} D^{0.5}, \quad (6)$$

587 where $\beta = \frac{2(1-\alpha)}{(2-\alpha)}$, and $R_{\beta,NL} = \sum_{m=1} \frac{\Phi_m \Phi_m^T}{[\beta + 2\alpha(1-\beta)\eta_m]}$. Φ_m and η_m are the m^{th} eigenvector and
 588 eigenvalue of the normalized Laplacian. In the expression of $R_{\beta,NL}$, the β and α regulate the
 589 weight of contribution in each eigenvalue-eigenvector pair of the summation such that the first
 590 eigenvalue-eigenvector pair (corresponding to the stationary distribution and given by the
 591 local-node degree-properties) remains included in the overall expression, but does not overwhelm
 592 the global information provided by subsequent ‘eigen-pairs’. Moreover, computation of $R_{\beta,NL}$ is
 593 not limited to a subset of the first k eigenvectors (bypassing the need for the user to select a
 594 suitable threshold or subset of eigenvectors) since the dimensionality is not on the order of
 595 number of cells, but equal to the number of clusters and hence all eigenvalue-eigenvector pairs
 596 can be incorporated without causing a bottleneck in runtime.
 The expected hitting time from node q to node r is given by⁴⁴,

$$h_{\alpha}(q, r) = \frac{[pr_{\alpha}(e_r)^T](r)}{d_r} - \frac{[pr_{\alpha}(e_r)^T](q)}{d_q} \quad (7)$$

where e_i is an indicator vector with 1 in the i^{th} entry and 0 elsewhere (i.e. $s_m = 1$ if $m = i$ and $s_m = 0$ if $m \neq i$). We can substitute Eq. (6) into Eq. (7), making use of the fact that $\frac{1}{d_r} = [D^{-1}e_r]^T(r)$, and $D^{-0.5}R_{\beta,NL}D^{-0.5}$ is symmetric, to obtain a closed form expression of the hitting time in terms of $R_{\beta,NL}$

$$h_{\alpha}(q, r) = \beta(e_r - e_q)^T D^{-0.5}R_{\beta,NL}D^{-0.5}e_r \quad (8)$$

(ii) *MCMC simulation*: The hitting time metric computed in Step-1 is used to infer graph-directionality. Instead of pruning edges in the ‘reverse’ direction, edge-weights are biased based on the time difference between nodes using the logistic function with growth factor $b=1$.

$$f(t) = \frac{1}{1 + e^{-b(t_1 - t_0)}}$$

We then recompute the pseudotimes on the forward biased graph: Since there is no closed form solution of hitting times on a *directed* graph, we perform MCMC simulations (parallely processed to enable fast simulations of 1000s of teleporting, lazy random walks starting at the root node of the cluster graph) and use the first quartile of the simulated pseudotime values for a respective node as the refined pseudotime for that node relative to the root. This refinement step ensures that the pseudotime is robust to the spurious links (or conversely, links that are too weakly weighted) that can distort calculations based purely on the closed form solution of hitting times (**Supplementary Fig. 7d**). By using this 2-step pseudotime computation, VIA mitigates the issues of convergence issues and spurious edge-weights, both of which are common in random-walk pseudotime computation on large and complex datasets¹².

3. **Automated terminal-state detection**. The algorithm uses the refined directed and weighted graph (edges are re-weighted using the refined pseudotimes) to predict which nodes represent the terminal states based on a consensus vote of pseudotime and multiple vertex connectivity properties, including out-degree (i.e. the number of edges directed out of a node), closeness $C(q)$, and betweenness $B(q)$.

$$C(q) = \frac{1}{\sum_{q \neq r} l(q,r)}$$

597

$$B(q) = \sum_{r \neq q \neq t} \frac{\sigma_{rt}(q)}{\sigma_{rt}}$$

598 $l(q, r)$ is the distance between node q and node r (i.e. the sum of edges in a shortest path connecting
599 them). σ_{rt} is the total number of shortest paths from node r to node t . $\sigma_{rt}(q)$ is the number of these
600 paths passing through node q . The consensus vote is performed on nodes that score above (or below
601 for out-degree) the median in terms of connectivity properties. We show on multiple simulated and
602 real biological datasets that VIA more accurately predicts the terminal states, across a range of input
603 data dimensions and key algorithm parameters, than other methods attempting the same
604 **(Supplementary Fig. S1)**.
605

606 **4. Automated trajectory reconstruction.** VIA then identifies the most likely path of each lineage by
607 computing the likelihood of a node traversing towards a particular terminal state (e.g. differentiation).
608 These lineage likelihoods are computed as the visitation frequency under lazy-teleporting MCMC
609 simulations from the root to a particular terminal state, i.e. the probability of *node i* reaching
610 *terminal-state j* as the number of times *cell i* is visited along a successful path (i.e. *terminal-state j* is
611 reached) divided by the number of times *cell i* is visited along all of the simulations. In contrast to
612 other trajectory reconstruction methods which compute the shortest paths between root and terminal
613 node^{1,2}, the lazy-teleporting MCMC simulations in VIA offer a probabilistic view of pathways under
614 relaxed conditions that are not only restricted to the random-walk along a tree-like graph, but can also
615 be generalizable to other types of topologies, such as cyclic or connected/disconnected paths. In the
616 same vein, we avoid confining the graph to an absorbing Markov chain^{13,3} (AMC) as this places
617 prematurely strict / potentially inaccurate constraints on node-to-node mobility and can impede
618 sensitivity to cell fates (as demonstrated by VIA's superior cell fate detection across numerous
619 datasets **(Supplementary Fig. S1)**).

620 Downstream visualization and analysis

621 VIA generates a visualization that combines the network topology and single-cell level
622 pseudotime/lineage probability properties onto an embedding based on UMAP or PHATE. Generalized
623 additive models (GAMs) are used to draw edges found in the high-dimensional graph onto the lower
624 dimensional visualization **(Fig. 1)**. An unsupervised downstream analysis of cell features (e.g. marker
625 gene expression, protein expression or image phenotype) along pseudotime for each lineage is performed
626 **(Fig. 1)**. Specifically, VIA plots the expression of features across pseudotime for each lineage by using
627 the lineage likelihood properties to weight the GAMs. A cluster-level lineage pathway is automatically
628 produced by VIA to visualize feature heat maps at the cluster-level along a lineage-path to see the
629 regulation of genes. VIA provides the option of gene imputation before plotting the lineage specific gene
630 trends. The imputation is fast as it relies on the single-cell KNN (scKNN) graph computed in Step 1.
631 Using an affinity-based imputation method⁴⁵, this step computes a “diffused” transition matrix on the
632 scKNN graph used to impute and denoise the original gene expressions.

633 **Benchmarked Methods**

634 The methods were mainly chosen based on their superior performance in a recent large-scale
635 benchmarking study⁴, including a select few recent methods claiming to supersede those in the study.
636 Specifically, recent and popular methods exhibiting reasonable scalability, and automated cell fate
637 prediction in multi-lineage trajectories were favoured as candidates for benchmarking (See
638 **Supplementary Table S1** for the key characteristics of methods). Performance stress-tests in terms of
639 lineage detection of each biological dataset, and pseudotime correlation for time-series data were
640 conducted over a range of key input parameters (e.g. numbers of k-nearest neighbors, highly variable
641 genes (HVGs), principal components (PCs)) and pre-processing protocols (see **Supplementary Fig. 1**).
642 All comparisons were run on a computer with an Intel(R) Xeon (R) W-2123 central processing unit
643 (3.60GHz, 8 cores) and 126 GB RAM.

644

645 Quantifying terminal state prediction accuracy for parameter tests was done using the F1-score, defined
646 as the harmonic mean of recall and precision and calculated as:

$$647 \quad F_1 = \frac{tp}{tp + 0.5(fp + fn)}$$

648 Where tp is a true-positive: the identification of a terminal cluster that is in fact a final differentiated cell
649 fate; fp is a false positive identification of a cluster as terminal when in fact it represents an intermediate
650 state; and fn is a false negative where a known cell fate fails to be identified

651

652 **PAGA**²⁸. It uses a cluster-graph representation to capture the underlying topology. PAGA computes a
653 unified pseudotime by averaging the single-cell level diffusion pseudotime computed by DPT, but
654 requires manual specification of terminal cell fates and clusters that contribute to lineages of interest in
655 order to compare gene expression trends across lineages.

656

657 **Palantir**². It uses diffusion-map⁴⁶ components to represent the underlying trajectory. Pseudotimes are
658 computed as the shortest path along a KNN-graph constructed in a low-dimensional diffusion component
659 space, with edges weighted such that the distance between nodes corresponds to the diffusion
660 pseudotime⁴⁷ (DPT). Terminal states are identified as extrema of the diffusion maps that are also outliers
661 of the stationary distribution. The lineage-likelihood probabilities are computed using Absorbing Markov
662 Chains (constructed by removing outgoing edges of terminal states, and thresholding reverse edges).

663

664 **Slingshot**¹. It is designed to process low-dimensional embeddings of the single-cell data. By default
665 Slingshot runs clustering based on Gaussian mixture modeling and recommends using the first few PCs as
666 input. Slingshot connects the clusters using a minimum spanning tree and then fits principle curves for
667 each detected branch. It uses the orthogonal projection against each principal curve to fit a separate
668 pseudotime for each lineage, and hence the gene expressions cannot be compared across lineages. Also,
669 the runtimes are prohibitively long for large datasets or high input dimensions.

670

671 **CellRank**¹³. This method combines the information of RNA velocity (computed using scVelo⁴⁸) and
672 gene-expression to infer trajectories. Given it is mainly suited for the scRNA-seq data, with the

673 RNA-velocity computation limiting the overall runtime for larger dataset, we limit our comparison to the
674 pancreatic dataset which the authors of CellRank used to highlight its performance.

675

676 **Monocle3**³⁶. The workflow consists of three steps: the first is to project the data to two or three
677 dimensions using UMAP (this is a strict requirement), followed by Louvain clustering on a K-Nearest
678 Neighbor graph constructed in the low-dimensional UMAP space. A cluster-graph is then created and
679 partitioned to deduce disconnected trajectories. Subsequently, it learns a principal graph in the
680 low-dimensional space along which it calculates pseudotimes as the geodesic distance from root to cell.

681 **Simulated Data**

682 We employed the DynToy⁴ (<https://github.com/dynverse/dyntoy>) package, which generates synthetic
683 single-cell gene expression data (~1000 cells x 1000 ‘genes’), to simulate different complex trajectory
684 models. Using these datasets, we tested that VIA consistently and more accurately captures both tree and
685 non-tree like structures (multifurcating, cyclic, and disconnected) compared to other methods (**Fig.2**). All
686 methods are subject to the same data pre-processing steps, PCA dimension reduction and root-cell to
687 initialize the path. Graph edge accuracy is computed based on an F1-score of connectivity in the TI
688 generated versus reference graphs. For example, an edge is considered a true positive if it connects two
689 states that are made of the same cell type or of two cell types that are connected in the reference truth. A
690 false negative is the lack of an edge to connect to cell types that are connected in the reference.

691

692 **Multifurcating structure.** This dataset consists of 1000 ‘cells’ multifurcating into 4 terminal states. VIA
693 robustly captures all four terminal cell fates across a range of input PCs and the pseudotimes are well
694 inferred relative to the root node (**Fig. 2a**). Note that two terminal states (M2 and M8), which are very
695 close to each other, are easily merged by the other methods (Slingshot, Palantir, Monocle3, and PAGA).

696 **Cyclic structure.** We ran VIA and other methods for different values of K nearest neighbors. VIA
697 unambiguously shows a cyclic network for a range of K (in KNN). Slingshot does not use a KNN
698 parameter and shows 3 fragmented different lineages (top to bottom). PAGA fails to capture the
699 connected cyclic structure at K = 10 and 5, while Palantir visually shows a linear (K = 10, 30) or
700 disconnected structure (K = 5). Monocle recovers a linear trajectory, failing to detect the loop closure.
701 Van den Berge et al⁵⁷ also find that Monocle3 consistently fragments or fits branching structures onto
702 cyclic simulated datasets.

703 **Disconnected structure.** This dataset comprises two disconnected trajectories (T1 and T2). T1 is cyclic
704 with an extra branch (M5 to M6), T2 has a bifurcation at M3 (**Fig. 2c**). VIA captures the two
705 disconnected structures as well as the M6 branch in the cyclic structure, and the bifurcation in the smaller
706 structure. PAGA captures the underlying structure at PC = 20 but becomes fragmented for other numbers
707 of PCs. Palantir also yields multiple fragments and is not able to capture the overall structure, while
708 Slingshot (using the default clustering based on Gaussian mixture modeling) connects T1 and T2, and
709 only captures one of the bifurcations in T1.

710 **Biological Data**

711 The pre-processing steps described below for each dataset are not included in the reported runtimes as
712 these steps are typically very fast, (typically less than 1-10% of the total runtime depending on the
713 method. E.g. only a few minutes for pre-processing 100,000s of cells) and only need to be performed
714 once as they remain the same for all subsequent analyses. It should also be noted that visualization (e.g.
715 UMAP, t-SNE) are not included in the runtimes. VIA provides a subsampling option at the visualization
716 stage to accelerate this process for large datasets without impacting the previous computational steps.
717 However, to ensure fair comparisons between TI methods (e.g. other methods do not have an option to
718 compute the embedding on a subsampled input and transfer the results between the full trajectory and the
719 sampled visualization, or rely on a slow version of tSNE), we simply provide each TI method with a
720 pre-computed visualization embedding on which the computed results are projected.

721

722 **ScRNA-seq of mouse pre-B cells.** This dataset²⁶ models the pre-BI cell (Hardy fraction C') process
723 during which cells progress to the pre-BII stage and B cell progenitors undergo growth arrest and
724 differentiation. Measurements were obtained at 0, 2, 6, 12, 18 and 24 hours (h) for a total of 313 cells x
725 9,075 genes. We follow a standard Scanpy preprocessing recipe⁴⁹ that filters cells with low counts, and
726 genes that occur in less than 3 cells. The filtered cells are normalized by library size and log transformed.
727 The top 5000 highly variable genes (HVG) are retained. Cells are renormalized by library count and
728 scaled to unit variance and zero mean. VIA identifies the terminal state at 18-24 h and accurately
729 recapitulates the gene expression trends²⁶ along inferred pseudotime of *Igll1*, *Slc7a5*, *Foxo1*, *Myc*, *Ldha*
730 and *Lig4*. (**Supplementary Fig. S2a**). We show the results generalize across a range of PCs for two
731 values of K of the graph with higher accuracy in locating the later cell fates than Slingshot and Palantir.
732 (**Supplementary Fig. S2b**).

733

734 **ScRNA-seq of human CD34+ bone marrow cells.** This is a scRNA-seq dataset of 5800 cells
735 representing human hematopoiesis². We used the filtered, normalized and log-transformed count matrix
736 provided by Setty et al², with PCA performed on all the remaining genes. The cells were annotated using
737 SingleR⁵⁰ which automatically labeled cells based on the hematopoietic reference dataset Novershtern
738 Hematopoietic Cell Data - GSE24759⁵¹. The annotations are in agreement with the labels inferred by
739 Setty et al. for the 7 clusters, including the root HSCs cluster that differentiates into 6 different lineages:
740 monocytes, erythrocytes, and B cells, as well as the less populous megakaryocytes, cDCs and pDCs. VIA
741 consistently identifies these lineages across a wider range of input parameters and data dimensions (e.g.
742 the number of K and PCs provided as input to the algorithms see **Fig. 2p**, and **Supplementary Fig. S3c**).
743 Notably, the upregulated gene expression trends of the small populations can be recovered in VIA, i.e.
744 pDC and cDC show elevated CD123 and CSF1R levels relative to other lineages, and the upregulated
745 CD41 expression in megakaryocytes (**Supplementary Fig. S3-S4**).

746

747 **ScRNA-seq of human embryoid body.** This is a mid-sized scRNA-seq dataset of 16,825 human cells in
748 embryoid bodies (EBs)¹⁵. We followed the same pre-processing steps as Moon et al. to filter out dead
749 cells and those with too high or low library count. Cells are normalized by library count followed by
750 square root transform. Finally the transformed counts are scaled to unit variance and zero mean. The

751 filtered data contained 16825 cells \times 17580 genes. PCA is performed on the processed data before
752 running each TI method. VIA identifies 6 cell fates, which, based on the upregulation of marker genes as
753 cells proceed towards respective lineages, are in accord with the annotations given by Moon et al., (See
754 the gene heatmap and changes in gene expression along respective lineage trajectories in **Supplementary**
755 **Fig. S5**). Note that Palantir and Slingshot do not capture the cardiac cell fate, and Slingshot also misses
756 the neural crest (see the F1-scores summary for terminal state detection **Supplementary Fig. S5**).

757

758 **ScRNA-seq of mouse organogenesis cell atlas.** This is a large and complex scRNA-seq dataset of mouse
759 organogenesis cell atlas (MOCA) consisting of 1.3 million cells⁶. The dataset contains cells from 61
760 embryos spanning 5 developmental stages from early organogenesis (E9.5-E10.5) to organogenesis
761 (E13.5). Of the 2 million cells profiled, 1.3 million are ‘high-quality’ cells that are analysed by VIA. The
762 runtime is approximately 40 minutes which is in stark contrast to the next fastest tool Palantir which takes
763 4 hours (excluding visualization). The authors of MOCA manually annotated 38 cell-types based on the
764 differentially expressed genes of the clusters. In general, each cell type exclusively falls under one of 10
765 major and disjoint trajectories inferred by applying Monocle3 to the UMAP of MOCA. The authors
766 attributed the disconnected nature of the 10 trajectories to the paucity of earlier stage common
767 predecessor cells. We followed the same steps as Cao et al.⁶ to retain high-quality cells (i.e. remove cells
768 with less than 400 mRNA, and remove doublet cells and cells from doubled derived sub-clusters). PCA
769 was applied to the top 2000 HVGs with the top 30 PCs selected for analysis. VIA analyzed the data in the
770 high-dimensional PC space. We bypass the step in Monocle3⁶ which applies UMAP on the PCs prior to
771 TI as this incurs an additional bias from choice of manifold-learning parameters and a further loss in
772 neighborhood information. As a result, VIA produces a more connected structure with linkages between
773 some of the major cell types that become segregated in UMAP (and hence Monocle3), and favors a
774 biologically relevant interpretation (**Fig. 2, Supplementary Fig. S11**). A detailed explanation of these
775 connections (graph-edges) extending between certain major groups using references to literature on
776 organogenesis is presented in **Supplementary Note 3**.

777

778 **ScRNA-seq of murine endocrine development**⁵. This is an scRNA-seq dataset of E15.5 murine
779 pancreatic cells spanning all developmental stages from an initial endocrine progenitor-precursor (EP)
780 state (low level of *Ngn3*, or *Ngn3^{low}*), to the intermediate EP (high level of *Ngn3*, or *Ngn3^{high}*) and *Fev⁺*
781 states, to the terminal states of hormone-producing alpha, beta, epsilon and delta cells⁵. Following steps
782 by Lange et al¹³, we preprocessed the data using scVelo to filter genes, normalize each cell by total counts
783 over all genes, keep the top most variable genes, and take the log-transform. PCA was applied to the
784 processed gene matrix. We assessed the performance of VIA and other TI methods (CellRank, Palantir,
785 Slingshot) across a range of number of retained HVGs and input PCs (**Fig. 2m, Supplementary Fig. S6**).

786

787 **ScATAC-seq of human bone marrow cells.** This scATAC-seq data profiles 3072 cells isolated from
788 human bone marrow using fluorescence activated cell sorting (FACS), yielding 9 populations²⁷: HSC,
789 MPP, CMP, CLP, LMPP, GMP, MEP, mono and plasmacytoid DCs (**Fig. 3a and Supplementary Fig.**
790 **S7**). We examined TI results for two different preprocessing pipelines to gauge how robust VIA is on the
791 scATAC-seq analysis which is known to be challenging for its extreme intrinsic sparsity. We used the
792 pre-processed data consisting of PCA applied to the z-scores of the transcription factor (TF) motifs used
793 by Buenrostro et al²⁷. Their approach corrects for batch effects in select populations and weighting of PCs

794 based on reference populations and hence involves manual curation. We also employed a more general
795 approach used by Chen et al.³¹ which employs ChromVAR to compute k-mer accessibility z-scores across
796 cells. VIA infers the correct trajectories and the terminal cell fates for both of these inputs, again across a
797 wide range of input parameters (**Fig. 3d and Supplementary Fig. S7**).

798

799 **scRNA-seq and scATAC-seq of *Isl1*+ cardiac progenitor cells.** This time-series dataset captures
800 murine *Isl1*+ cardiac progenitor cells (CPCs) from E7.5 to E9.5 characterized by scRNA-seq (197 cells)
801 and scATAC-seq (695 cells)²⁰. The *Isl1*+ CPCs are known to undergo multipotent differentiation to
802 cardiomyocytes or endothelial cells. For the scRNA-seq data, the quality filtered genes and the size-factor
803 normalized expression values are provided by Jia et al.²⁰ as a “Single Cell Expression Set” object in R.
804 Similarly, the cells in the scATAC-seq experiment were provided in a “SingleCellExperiment” object with
805 low quality cells excluded from further analysis. The accessibility of peaks was transformed to a binary
806 representation as input for TF-IDF (term frequency-inverse document frequency) weighting prior to
807 singular value decomposition (SVD). The highlighted TF motifs in the heatmap (**Fig. 2j**) correspond to
808 those highlighted by Jia et al. We tested the performance when varying the number of SVDs used. We
809 also considered the outcome when merging the scATAC-seq and scRNA-seq data using Seurat3⁵².
810 Despite the relatively low cell count of both datasets, and the relatively under-represented scRNA-seq cell
811 count, the two datasets overlapped reasonably well and allowed us to infer the expected lineages in an
812 unsupervised manner (**Fig. 2d and Supplementary Fig. S8**). In contrast, Jia et al., performed a supervised
813 TI by manually selecting cells relevant to the different lineages (for the scATAC-seq cells) and choosing
814 the two diffusion components that best characterize the developmental trajectories in low dimension²⁰.

815

816 **Mass cytometry data of mouse embryonic stem cells (mESC).** This is a mass cytometry (or CyTOF)
817 dataset, consisting of 90,000 cells and 28 antibodies (corresponding to ~7000 cells each from Day 0-11
818 measurements), that represents differentiation of mESC to mesoderm cells³². An arcsinh transform with a
819 scaling factor of 5 was applied on all features - a standard procedure for CyTOF datasets, followed by
820 normalization to unit variance and zero mean. All 28 antibodies are used by the TI methods (with the
821 exception of Slingshot which requires PCA followed by subsetting of the first 5 PCs in order to
822 computationally handle the high cell count) (**Supplementary Fig. S9**). To improve Palantir performance
823 we used 5000 waypoints (instead of default 1200) but this takes almost 20 minutes to complete
824 (excluding time taken for embedding the visualization). VIA runs in ~3 minutes and produces results
825 consistent with the known ordering and identifies regions of Day 10-11 cells.

826

827 **Single-cell biophysical phenotypes derived from imaging flow cytometry.** This is the in-house dataset
828 of single-cell biophysical phenotypes of two different human breast cancer types (MDA-MB231 and
829 MCF7). Following our recent image-based biophysical phenotyping strategy^{53,54}, we defined the
830 spatially-resolved biophysical features of a cell in a hierarchical manner based on both bright-field and
831 quantitative phase images captured by the FACED imaging flow cytometer (i.e., from the bulk features to
832 the subcellular textures). At the bulk level, we extracted the cell size, dry mass density, and cell shape. At
833 the subcellular texture level, we parameterized the global and local textural characteristics of optical
834 density and mass density at both the coarse and fine scales (e.g., local variation of mass density, its
835 higher-order statistics, phase entropy radial distribution etc.). This hierarchical phenotyping approach^{53,54}
836 allowed us to establish a single-cell biophysical profile of 38 features, which were normalized based on

837 the z-score (See Supplementary Table S4 and Table S5). All these features, without any PCA, are used
838 as input to VIA. In order to weigh the features, we use a mutual information classifier to rank the features,
839 based on the integrated fluorescence intensity of the fluorescence FACED images of the cells (which
840 serve as the ground truth of the cell-cycle stages). Following normalization, the top 3 features (which
841 relate to cell size) are weighted (using a factor between 3-10).

842 **Imaging flow cytometry experiment**

843 **FACED imaging flow cytometer setup**

844 A multimodal FACED imaging flow cytometry (IFC) platform was used to obtain the quantitative phase
845 and fluorescence images of single cells in microfluidic flow at an imaging throughput of ~70,000
846 cells/sec. The light source consisted of an Nd:YVO picosecond laser (center wavelength = 1064 nm,
847 Time-Bandwidth) and a periodically-poled lithium niobate (PPLN) crystal (Covesion) for second
848 harmonic generation of a green pulsed beam (center wavelength = 532 nm) with a repetition rate of 20
849 MHz. The beam was then directed to the FACED module, which mainly consists of a pair of
850 almost-parallel plane mirrors. This module generated a linear array of 50 beamlets (foci) which were
851 projected by an objective lens (40X, 0.6NA, MRH08430, Nikon) on the flowing cells in the microfluidic
852 channel for imaging. Each beamlet was designed to have a time delay of 1 ns with the neighboring
853 beamlet in order to minimize the fluorescence crosstalk due to the fluorescence decay. Detailed
854 configuration of the FACED module can be referred to Wu et al.³³. The epi-fluorescence image signal
855 was collected by the same objective lens and directed through a band-pass dichroic beamsplitter (center:
856 575nm, bandwidth: 15nm). The filtered orange fluorescence signal was collected by the photomultiplier
857 tube (PMT) (rise time: 0.57 ns, Hamamatsu). On the other hand, the transmitted light through the cell was
858 collected by another objective lens (40X, 0.8NA, MRD07420, Nikon). The light was then split equally by
859 the 50:50 beamsplitter into two paths, each of which encodes different phase-gradient image contrasts of
860 the same cell (a concept similar to Scherlien photography⁵⁵). The two beams are combined,
861 time-interleaved, and directed to the photodetector (PD) (bandwidth: >10 GHz, Alphas) for detection.
862 The signals obtained from both PMT and PD were then passed to a real-time high-bandwidth digitizer (20
863 GHz, 80 GS/s, Lecroy) for data recording.

864

865 **Cell culture and preparation**

866 MDA-MB231 (ATCC) and MCF7 (ATCC), which are two different breast cancer cell lines, were used for
867 the cell cycle study. The culture medium for MDA-MB231 was ATCC modified RPMI 1640 (Gibco)
868 supplemented with 10% fetal bovine serum (FBS) (Gibco) and 1% antibiotic-antimycotic (Anti-Anti)
869 (Gibco), while that for MCF7 was DMEM supplemented with 10% FBS (Gibco) and 1% Anti-Anti
870 (Gibco). The cells were cultured inside an incubator under 5% CO₂ and 37°C, and subcultured twice a
871 week. 1e6 cells were pipetted out from each cell line and stained with Vybrant DyeCycle orange stain
872 (Invitrogen).

873 Data Availability

874 Data used in Figures 1-3 as well as Supplementary Figures S1-S15) is available on:

- 875 1. Pancreatic data: Gene Expression Omnibus (GEO) under accession code GSE132188.
 - 876 2. Cardiac progenitor data is available from the ENA repository under the accession code
877 PRJEB23303 or from [<https://github.com/loosolab/cardiac-progenitors>].
 - 878 3. B-cell: STATegraData GitHub repository. [<https://github.com/STATegraData/STATegraData>]
 - 879 4. Mass cytometry mesoderm: Cytobank
880 [<https://community.cytobank.org/cytobank/experiments/71953>].
 - 881 5. Raw and processed data for scRNA-seq Human Hematopoiesis are available through the Human
882 Cell Atlas data portal at
883 <https://data.humancellatlas.org/explore/projects/091cf39b-01bc-42e5-9437-f419a66c8a45>.
 - 884 6. Embryoid Body: Mendeley Data repository at <https://doi.org/10.17632/v6n743h5ng.1>.
 - 885 7. Mouse Organogenesis : NCBI Gene Expression Omnibus under accession number GSE119945
 - 886 8. FACED cell cycle: <https://github.com/ShobiStassen/VIA> and on FigShare
887 <https://doi.org/10.6084/m9.figshare.13601405.v1>
 - 888 9. scATAC-seq Hematopoiesis: GEO: GSE96772. Processed scATAC-seq data, which include PC
889 values and TF scores per cell can be found in Data S1. of
890 <https://doi.org/10.1016/j.cell.2018.03.074>
 - 891 10. Toy Data: <https://github.com/ShobiStassen/VIA>
- 892

893 Code Availability

894 VIA is available as a pip installable python library “pyVIA” with tutorials and sample data available on
895 <https://github.com/ShobiStassen/VIA> and <https://pypi.org/project/pyVIA/>
896

897 References

- 898 1. Street, K. et al. Slingshot: cell lineage and pseudotime inference for single-cell transcriptomics.
899 BMC Genomics 19, 477 (2018).
- 900 2. Setty M, Kisieliovas V, Levine J, Gayoso A, Mazutis L, Pe'er D. Characterization of cell fate
901 probabilities in single-cell data with Palantir [published correction appears in Nat Biotechnol.
902 2019 Oct;37(10):1237]. Nat Biotechnol. 2019;37(4):451-460. doi:10.1038/s41587-019-0068-4
- 903 3. Qiu, X., Mao, Q., Tang, Y. *et al.* Reversed graph embedding resolves complex single-cell
904 trajectories. *Nat Methods* 14, 979–982 (2017). <https://doi.org/10.1038/nmeth.4402>
- 905 4. Saelens, W., Cannoodt, R., Todorov, H. et al. A comparison of single-cell trajectory inference
906 methods. *Nat Biotechnol* 37, 547–554 (2019). <https://doi.org/10.1038/s41587-019-0071-9>
- 907 5. Bastidas-Ponce, A. et al. Comprehensive single cell mRNA profiling reveals a detailed roadmap
908 for pancreatic endocrinegenesis. *Development* 146, (2019).
- 909 6. Cao, J. et al. Comprehensive single- cell transcriptional profiling of a multicellular organism.
910 *Science* 357,661–667 (2017).

- 911 7. Packer, J. S. et al. A lineage- resolved molecular atlas of *C. elegans* embryogenesis at single- cell
912 resolution. *Science* 365, eaax1971 (2019).
- 913 8. Cao, J., Spielmann, M., Qiu, X. et al. The single-cell transcriptional landscape of mammalian
914 organogenesis. *Nature* 566, 496–502 (2019).
- 915 9. Briggs, J. A. et al. The dynamics of gene expression in vertebrate embryogenesis at single- cell
916 resolution. *Science* 360, eaar5780 (2018).
- 917 10. Litviňuková, M., Talavera-López, C., Maatz, H. et al. Cells of the adult human heart. *Nature*
918 (2020).
- 919 11. Stassen SV, Siu DMD, Lee KCM, Ho JWK, So HKH, Tsia KK. PARC: ultrafast and accurate
920 clustering of phenotypic data of millions of single cells. *Bioinformatics*. 2020 May
921 1;36(9):2778-2786. doi: 10.1093/bioinformatics/btaa042.
- 922 12. Ulrike von Luxburg, Agnes Rad, Matthias Hein. Hitting and Commute Times in Large Random
923 Neighborhood Graphs. *Journal of Machine Learning Research* 15, 1751-1798 (2014)
- 924 13. Marius Lange, Volker Bergen, Michal Klein, Manu Setty, Bernhard Reuter, Mostafa Bakhti,
925 Heiko Lickert, Meshal Ansari, Janine Schniering, Herbert B. Schiller, Dana Pe'er, Fabian J.
926 Theis. CellRank for directed single-cell fate mapping. *bioRxiv* 2020.10.19.345983; doi:
927 <https://doi.org/10.1101/2020.10.19.345983>
- 928 14. McInnes, L., Healy, J., Saul, N. & Großberger, L. UMAP: uniform manifold approximation and
929 projection. *J. Open Source Software*. 3, 861 (2018).
- 930 15. Moon, K.R., van Dijk, D., Wang, Z. et al. Visualizing structure and transitions in
931 high-dimensional biological data. *Nat Biotechnol* 37, 1482–1492 (2019).
932 <https://doi.org/10.1038/s41587-019-0336-3>
- 933 16. Tam PP, Behringer RR. Mouse gastrulation: the formation of a mammalian body plan. *Mech Dev*.
934 1997;68(1-2):3-25. doi:10.1016/s0925-4773(97)00123-8
- 935 17. Chin AM, Hill DR, Aurora M, Spence JR. Morphogenesis and maturation of the embryonic and
936 postnatal intestine. *Semin Cell Dev Biol*. 2017 Jun;66:81-93. doi: 10.1016/j.semcd.2017.01.011.
937 Epub 2017 Feb 1.
- 938 18. Gilbert SF. *Developmental Biology*. 6th edition. Sunderland (MA): Sinauer Associates; 2000. The
939 Neural Crest. Available from: <https://www.ncbi.nlm.nih.gov/books/NBK10065/>
- 940 19. The human body at cellular resolution: the NIH Human Biomolecular Atlas Program, *Nature*
941 (2019) <https://doi.org/10.1038/s41586-019-1629-x>
- 942 20. Jia G, Preussner J, Chen X, Guenther S, Yuan X, Yekelchik M, Kuenne C, Looso M, Zhou Y,
943 Teichmann S, Braun T. Single cell RNA-seq and ATAC-seq analysis of cardiac progenitor cell
944 transition states and lineage settlement. *Nat Commun*. 2018 Nov 19;9(1):4877.
- 945 21. Tanya E. Foley, Bradley Hess, Joanne G. A. Savory, Randy Ringuette, David Lohnes. Role of Cdx
946 factors in early mesodermal fate decisions. *Development* 2019 146: dev170498 doi:
947 10.1242/dev.170498 Published 1 April 2019
- 948 22. Yao Y, Yao J, Boström KI. SOX Transcription Factors in Endothelial Differentiation and
949 Endothelial-Mesenchymal Transitions. *Front Cardiovasc Med*. 2019;6:30. Published 2019 Mar
950 28. doi:10.3389/fcvm.2019.00030
- 951 23. Potta SP, Liang H, Winkler J, Doss MX, Chen S, Wagh V, Pfannkuche K, Hescheler J, Sachinidis
952 A. Isolation and functional characterization of alpha-smooth muscle actin expressing

- 953 cardiomyocytes from embryonic stem cells. *Cell Physiol Biochem*. 2010;25(6):595-604. doi:
954 10.1159/000315078. Epub 2010 May 18. PMID: 20511704.
- 955 24. Warkman AS, Whitman SA, Miller MK, Garriock RJ, Schwach CM, Gregorio CC, Krieg PA.
956 Developmental expression and cardiac transcriptional regulation of Myh7b, a third myosin heavy
957 chain in the vertebrate heart. *Cytoskeleton (Hoboken)*. 2012 May;69(5):324-35. doi:
958 10.1002/cm.21029. Epub 2012 Apr 30. Erratum in: *Cytoskeleton (Hoboken)*. 2012
959 Dec;69(12):1086. PMID: 22422726; PMCID: PMC4734749.
- 960 25. Mahmoud AI, Kocabas F, Muralidhar SA, et al. Meis1 regulates postnatal cardiomyocyte cell
961 cycle arrest. *Nature*. 2013;497(7448):249-253. doi:10.1038/nature12054
- 962 26. Gomez-Cabrero, D., Tarazona, S., Ferreirós-Vidal, I. et al. STATegra, a comprehensive
963 multi-omics dataset of B-cell differentiation in mouse. *Sci Data* 6, 256 (2019).
964 <https://doi.org/10.1038/s41597-019-0202-7>
- 965 27. Jason D. Buenrostro, M. Ryan Corces, Caleb A. Lareau, Beijing Wu, Alicia N. Schep, Martin J.
966 Aryee, Ravindra Majeti, Howard Y. Chang, William J. Greenleaf, Integrated Single-Cell Analysis
967 Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation, *Cell*, 173,
968 1535-1548.e16, (2018) <https://doi.org/10.1016/j.cell.2018.03.074>.
- 969 28. Wolf, F. A. et al. PAGA: graph abstraction reconciles clustering with trajectory inference through
970 a topology preserving map of single cells. *Genome Biol*. 20, 59 (2019).
- 971 29. Gutierrez GD, Gromada J, Sussel L. Heterogeneity of the Pancreatic Beta Cell. *Front Genet*.
972 2017;8:22. Published 2017 Mar 6. doi:10.3389/fgene.2017.00022
- 973 30. Krentz NAJ, Lee MYY, Xu EE, Sproul SLJ, Maslova A, Sasaki S, Lynn FC. Single-Cell
974 Transcriptome Profiling of Mouse and hESC-Derived Pancreatic Progenitors. *Stem Cell Reports*.
975 2018 Dec 11;11(6):1551-1564. doi: 10.1016/j.stemcr.2018.11.008. PMID: 30540962; PMCID:
976 PMC6294286.
- 977 31. Chen, H., Lareau, C., Andreani, T. *et al.* Assessment of computational methods for the analysis of
978 single-cell ATAC-seq data. *Genome Biol* 20, 241 (2019).
979 <https://doi.org/10.1186/s13059-019-1854-5>
- 980 32. Ko, M.E., Williams, C.M., Fread, K.I. *et al.* FLOW-MAP: a graph-based, force-directed layout
981 algorithm for trajectory mapping in single-cell time course datasets. *Nat Protoc* 15, 398–420
982 (2020). <https://doi.org/10.1038/s41596-019-0246-3>
- 983 33. Wu J. L., Xu Y. Q., Xu J. J., Wei X. X., Chan A. C. S., Tang A. H. L., Lau A. K. S., Chung B. M.
984 F., Cheung Shum H., Lam E. Y., Wong K. K. Y., Tsia K. K., “Ultrafast laser-scanning time-stretch
985 imaging at visible wavelengths,” *Light Sci. Appl.* 6(1), e16196 (2016).10.1038/lsa.2016.196
- 986 34. Popescu G, Park Y, Lue N, Best-Popescu C, Deflores L, Dasari RR, Feld MS, Badizadegan K.
987 Optical imaging of cell mass and growth dynamics. *Am J Physiol Cell Physiol*. 2008
988 Aug;295(2):C538-44. doi: 10.1152/ajpcell.00121.2008. Epub 2008 Jun 18.
- 989 35. Kyoohyun Kim, Jochen Guck The Relative Densities of Cytoplasm and Nuclear Compartments
990 Are Robust against Strong Perturbation *Biophysical Journal*. Volume 119, Issue 10, 17 November
991 2020, Pages 1946-1957
- 992 36. Kafri R, Levy J, Ginzberg MB, Oh S, Lahav G, Kirschner MW. Dynamics extracted from fixed
993 cells reveal feedback linking cell growth to cell cycle. *Nature*. 2013 Feb 28;494(7438):480-3. doi:
994 10.1038/nature11897. PMID: 23446419; PMCID: PMC3730528.

- 995 37. Park SR, Namkoong S, Friesen L, Cho CS, Zhang ZZ, Chen YC, Yoon E, Kim CH, Kwak H,
996 Kang HM, Lee JH. Single-Cell Transcriptome Analysis of Colon Cancer Cell Response to
997 5-Fluorouracil-Induced DNA Damage. *Cell Rep*. 2020 Aug 25;32(8):108077. doi:
998 10.1016/j.celrep.2020.108077.
- 999 38. Zangle TA, Teitell MA. Live-cell mass profiling: an emerging approach in quantitative
1000 biophysics. *Nat Methods*. 2014 Dec;11(12):1221-8. doi: 10.1038/nmeth.3175. PMID: 25423019;
1001 PMCID: PMC4319180.
- 1002 39. Tse HT, Gossett DR, Moon YS, Masaeli M, Sohsman M, Ying Y, Mislick K, Adams RP, Rao J,
1003 Di Carlo D. Quantitative diagnosis of malignant pleural effusions by single-cell
1004 mechanophenotyping. *Sci Transl Med*. 2013 Nov 20;5(212):212ra163. doi:
1005 10.1126/scitranslmed.3006559. PMID: 24259051.
- 1006 40. Otto, O., Rosendahl, P., Mietke, A. *et al.* Real-time deformability cytometry: on-the-fly cell
1007 mechanical phenotyping. *Nat Methods* 12, 199–202 (2015). <https://doi.org/10.1038/nmeth.3281>
- 1008 41. Kimmerling, R.J., Prakadan, S.M., Gupta, A.J. *et al.* Linking single-cell measurements of mass,
1009 growth rate, and gene expression. *Genome Biol* 19, 207 (2018).
1010 <https://doi.org/10.1186/s13059-018-1576-0>
- 1011 42. Traag, V.A., Waltman, L. & van Eck, N.J. From Louvain to Leiden: guaranteeing well-connected
1012 communities. *Sci Rep* 9, 5233 (2019). <https://doi.org/10.1038/s41598-019-41695-z>
- 1013 43. Langville, Amy N., and Carl D. Meyer. *Google's PageRank and Beyond: The Science of Search*
1014 *Engine Rankings*. Princeton University Press, 2006.
- 1015 44. Chung F., Zhao W. (2010) PageRank and Random Walks on Graphs. In: Katona G.O.H.,
1016 Schrijver A., Szónyi T., Sági G. (eds) *Fete of Combinatorics and Computer Science*. Bolyai
1017 Society Mathematical Studies, vol 20. Springer, Berlin, Heidelberg.
- 1018 45. van Dijk D, Sharma R, Nainys J, et al. Recovering Gene Interactions from Single-Cell Data
1019 Using Data Diffusion. *Cell*. 2018;174(3):716-729.e27. doi:10.1016/j.cell.2018.05.061
- 1020 46. Coifman, R. R. et al. Geometric diffusions as a tool for harmonic analysis and structure definition
1021 of data: diffusion maps. *Proc. Natl Acad. Sci. USA* 102,7426–7431 (2005).
- 1022 47. Haghverdi L, Büttner M, Wolf FA, Buettner F, Theis FJ. Diffusion pseudotime robustly
1023 reconstructs lineage branching. *Nat Methods*. 2016;13(10):845-848. doi:10.1038/nmeth.3971
- 1024 48. Bergen, V., Lange, M., Peidli, S. et al. Generalizing RNA velocity to transient cell states through
1025 dynamical modeling. *Nat Biotechnol* 38, 1408–1414 (2020).
1026 <https://doi.org/10.1038/s41587-020-0591-3>
- 1027 49. Zheng GX, Terry JM, et al. Massively parallel digital transcriptional profiling of single cells. *Nat*
1028 *Commun*. 2017 Jan 16;8:14049. doi: 10.1038/ncomms14049.
- 1029 50. Aran D et al., (2019). “Reference-based analysis of lung single-cell sequencing reveals a
1030 transitional profibrotic macrophage.” *Nat. Immunol.*, 20, 163-172
- 1031 51. Novershtern N. et al., Densely interconnected transcriptional circuits control cell states in human
1032 hematopoiesis. *Cell*. 2011 Jan 21;144(2):296-309.
- 1033 52. Stuart T, Butler A, Hoffman P, Hafemeister C, Papalexi E, Mauck WM 3rd, Hao Y, Stoeckius M,
1034 Smibert P, Satija R. Comprehensive Integration of Single-Cell Data. *Cell*. 2019
1035 Jun13;177(7):1888-1902.e21. doi: 10.1016/j.cell.2019.05.031. Epub 2019 Jun 6.

- 1036 53. Siu, KCM Lee, MCK Lo, SV Stassen, M Wang, IZQ Zhang, HKH So. Deep-learning-assisted
1037 biophysical imaging cytometry at massive throughput delineates cell population heterogeneity.
1038 *Lab on a Chip* 20 (20), 3696-3708
- 1039 54. KCM Lee, M Wang, KSE Cheah, GCF Chan, HKH So, KKY Wong, KK Tsia.. Quantitative
1040 phase imaging flow cytometry for ultra-large-scale single-cell biophysical phenotyping.
1041 *Cytometry Part A* 95 (5), 510-520
- 1042 55. Wenwei Yan Jianglai Wu Kenneth K. Y. Wong Kevin K. Tsia, A high-throughput all-optical
1043 laser-scanning imaging flow cytometer with biomolecular specificity and subcellular resolution,
1044 *J. Biophotonics* (2017) <https://onlinelibrary.wiley.com/doi/abs/10.1002/jbio.201700178>
- 1045 56. F. Chung and S.-T. Yau, Discrete Green's Functions. *Journal of Combinatorial Theory, Series A*,
1046 91(1-2) (2000), pp. 191–214
- 1047 57. Van den Berge, K., Roux de Bézieux, H., Street, K. et al. Trajectory-based differential expression
1048 analysis for single-cell sequencing data. *Nat Communications*
- 1049 58. Yury A. Malkov, D. Yashunin. Efficient and Robust Approximate Nearest Neighbor Search Using
1050 Hierarchical Navigable Small World Graph, *Computer Science, Medicine, Mathematics, IEEE*
1051 *Transactions on Pattern Analysis and Machine Intelligence*, 2020
- 1052 59. *eLife* 2018;7:e29213 DOI: 10.7554/eLife.29213
- 1053 60. Phillip, J., Wu, PH., Gilkes, D. et al. Biophysical and biomolecular determination of cellular age
1054 in humans. *Nat Biomed Eng* 1, 0093 (2017). <https://doi.org/10.1038/s41551-017-0093>
- 1055 61. *eLife* 2017;6:e24060 DOI: 10.7554/eLife.24060
- 1056 62. Park, Y., Depeursinge, C. & Popescu, G. Quantitative phase imaging in biomedicine. *Nature*
1057 *Photon* 12, 578–589 (2018).
1058