

1 **Quantifying Cell-State Densities in Single-Cell Phenotypic Landscapes using Mellon**

2

3 Dominik Otto^{1,2,3}, Cailin Jordan^{1,2,3,4*}, Brennan Dury^{1,2,3*}, Christine Dien^{1,2,3}, Manu Setty^{1,2,3+}

4 ¹ Basic Sciences Division, Fred Hutchinson Cancer Center, Seattle WA

5 ² Computational Biology Program, Public Health Sciences Division, Seattle WA

6 ³ Translational Data Science IRC, Fred Hutchinson Cancer Center, Seattle WA

7 ⁴ Molecular and Cellular Biology Program, University of Washington, Seattle WA

8

9 * Equal contribution

10 + Corresponding author: Manu Setty (msetty@fredhutch.org)

11

12

13 **Abstract**

14 Cell-state density characterizes the distribution of cells along phenotypic landscapes and is crucial for
15 unraveling the mechanisms that drive cellular differentiation, regeneration, and disease. Here, we present
16 Mellon, a novel computational algorithm for high-resolution estimation of cell-state densities from single-
17 cell data. We demonstrate Mellon's efficacy by dissecting the density landscape of various differentiating
18 systems, revealing a consistent pattern of high-density regions corresponding to major cell types
19 intertwined with low-density, rare transitory states. Utilizing hematopoietic stem cell fate specification to
20 B-cells as a case study, we present evidence implicating enhancer priming and the activation of master
21 regulators in the emergence of these transitory states. Mellon offers the flexibility to perform temporal
22 interpolation of time-series data, providing a detailed view of cell-state dynamics during the inherently
23 continuous developmental processes. Scalable and adaptable, Mellon facilitates density estimation
24 across various single-cell data modalities, scaling linearly with the number of cells. Our work underscores
25 the importance of cell-state density in understanding the differentiation processes, and the potential of
26 Mellon to provide new insights into the regulatory mechanisms guiding cellular fate decisions.

27

28 Introduction

29 Cell differentiation is a dynamic process that underpins the development and function of all multicellular
30 organisms. Understanding how cells are distributed along differentiation trajectories is critical for
31 deciphering the mechanisms that drive cellular differentiation, pinpointing the key regulators and
32 characterizing the dysregulation of these processes in disease. Cell-state density is a representation of
33 this distribution of cells and is impacted by biological process spanning proliferation to apoptosis (**Fig.**
34 **1B**, **Supplementary Fig. 1A-C**). For instance, proliferation can increase the number of cells in a state,
35 resulting in high cell-state density (**Fig. 1B**). Cells converge to checkpoints that ensure the fidelity of the
36 differentiation process, also leading to high cell-state density (**Fig. 1B**). In contrast, transcriptional
37 acceleration, as seen in rare transitory cells, lead to lower cell-state density (**Fig. 1B**). Finally, apoptosis
38 decreases the number of cells in a state, also resulting in low cell-state density (**Fig. 1B**). As a result of
39 these influences, cell-state density of differentiation landscapes is likely not uniform (**Fig. 1A**) but exhibit
40 rich heterogeneity of high- and low-density regions (**Fig. 1B**).

41 Single-cell studies have underscored the importance of the heterogeneous nature of cell-state density in
42 single-cell phenotypic landscapes¹⁻³. Rapid and coordinated transcriptional acceleration leading to low-
43 density transitory states connecting higher-density regions have been demonstrated to be a hallmark of
44 developmental progression in diverse biological contexts from plants to humans⁴⁻⁷. Rare transitory cells
45 have also emerged as critical entities in the processes of differentiation¹, reprogramming⁸, and the
46 emergence of metastasis⁹. Despite the central importance of cell-state density, current approaches for
47 density estimation in single-cell data often produce noisy results and struggle to provide biologically
48 meaningful interpretation (**Supplementary Fig. 2**).

49 Here, we introduce Mellon, a novel computational algorithm to estimate cell-state density from single-cell
50 data (**Fig. 1C-G**). The core principle of Mellon is based on the intrinsic relationship between neighbor
51 distances and density, whereby distribution of nearest neighbor distances is linked with cell-state density
52 using a Poisson distribution (**Fig. 1C**). Mellon then connects densities between highly similar cell-states
53 using Gaussian processes to accurately and robustly compute cell-state densities that characterize
54 single-cell phenotypic landscapes (**Fig. 1D-E**). Unlike existing approaches that interpret single-cell
55 datasets solely as a collection of discrete cell states, Mellon infers a *continuous density function* across
56 the high-dimensional cell-state space (**Fig. 1F**), capturing the essential characteristics of the cell
57 population in its entirety. The density function can also be used to determine cell-state densities at single-
58 cell resolution (**Fig. 1G**). Mellon is designed to efficiently scale to increasingly prevalent atlas-scale
59 single-cell datasets and can be employed to infer cell-state density from diverse single-cell data
60 modalities.

61 We applied Mellon to dissect the density landscape of human hematopoiesis, revealing numerous high-
62 density regions corresponding to major cell types, intertwined with low-density, rare transitory cells. We
63 discovered a strong correlation between low-density regions and cell-fate specification, suggesting that
64 that lineage specification in hematopoiesis is driven by accelerated transcriptional changes. Exploration
65 of the open chromatin landscape during lineage specification hinted at the role of enhancer priming in
66 facilitating these transcriptional changes. Furthermore, extending Mellon's framework to time-series
67 datasets enabled us to compute time-continuous cell-state densities and interpolate cell-state densities
68 between observed timepoints, providing a high-resolution view of the cell-state dynamics during erythroid
69 differentiation in mouse gastrulation. Mellon, a scalable and user-friendly open-source software package,
70 complete with documentation and tutorials, is available at github.com/settylab/Mellon.

72 **Results**

73 **The Mellon modeling approach**

74 Mellon aims to compute cell-state densities within the intricate, high-dimensional single-cell phenotypic
75 landscapes. Two major challenges need to be resolved to estimate cell-state densities: First, the high-
76 dimensionality of single-cell data is an inherent computational obstacle, which Mellon overcomes by
77 leveraging the relationship between density and neighbor distances (**Methods**). The second challenge
78 lies in ensuring the precise and reliable density estimation in low-density states, which often represent
79 rare, transitory cells that play critical roles in a range of biological systems^{1,8-10}. To address this, Mellon
80 employs a strategy of estimating a continuous density function over the entire single-cell landscape. This
81 approach enhances both the accuracy and robustness of density estimation (**Methods**). Moreover, the
82 density function encapsulates a smooth and continuous portrayal of the high-dimensional phenotypic
83 landscape, enabling density estimation not only for individual measured cells—thus achieving single-cell
84 resolution—but also for unobserved cell-states, offering a comprehensive depiction of the entire cell
85 population (**Fig 1, Supplementary Fig. 1D**).

86 Mellon's utilization of neighbor distances and inference of continuous density function is underpinned by
87 two well-established principles of single-cell analysis. First, Mellon assumes that distances between cells
88 in the chosen representation of the phenotypic landscape are biologically meaningful and thus represent
89 a valid measure of cell-to-cell similarity. We refer to such a space as *cell-state space*, where each point
90 signifies a distinct cell state. To construct such a representation, we employ diffusion maps¹¹, a non-linear
91 dimensionality reduction technique that has been demonstrated to reliably and robustly represent the
92 single-cell phenotypic landscape^{12,13}. Moreover, distances within diffusion space are considered more
93 biologically informative than relying on gene expression-based distances (such as PCA)¹²⁻¹⁴ due to its
94 consideration of potential cell-state transition trajectories.

95 The second assumption Mellon relies on is that density changes from cell-to-cell are smooth and
96 continuous in nature i.e., Mellon assumes that cells with high degrees of similarity possess similar
97 densities. The inherent molecular heterogeneity of cells, primarily due to the subtle differences in gene
98 expression, supports these smooth density transitions. Further, single-cell studies have revealed that
99 cells experience gradual, rather than abrupt, changes in gene expression, providing empirical support for
100 this assumption^{1,15-17}.

101 Mellon first computes distance to the nearest neighbor for each cell in the cell-state space. We then
102 capitalize on the stochastic relationship between density and neighbor distances, where cells in higher
103 density states tend to exhibit shorter distances to their nearest neighbors, whereas cells in lower density
104 states tend to have longer distances (**Supplementary Fig 3A**). Formally, Mellon relates the nearest
105 neighbor distance to local density through the nearest neighbor distribution by employing a Poisson point
106 process (**Fig. 1C-D, Methods**). Nearest neighbor distribution describes the probability of another cell-
107 state existing within some distance of a given cell-state. Intuitively, regions with higher density of cell-
108 states correspond to tighter nearest-neighbor distributions, while lower densities result in broader
109 distributions (**Fig. 1D**).

110 Mellon then connects densities between highly related cells to estimate a continuous density function.
111 The true density function can be arbitrarily complex depending on the biological system. Mellon therefore
112 employs Gaussian Process (GP) in a Bayesian model to approximate this function without assuming a
113 specific functional form (**Fig. 1D**). GPs are a mathematical framework to model the patterns and
114 relationships among data points and, are highly effective for scenarios where the true functional form is
115 intricate or unknown and where observations are limited^{18,19}. GPs are thus ideally suited for density

116 estimation from noisy single-cell data. They achieve their robustness by incorporating the smoothness
117 assumption through a covariance kernel, facilitating sharing of information between adjacent
118 observations. In Mellon, the covariance kernel of the GP encodes cell-state similarity and determines the
119 influence of nearby cells on density estimates at a specific state (**Methods**). This covariance kernel is
120 effectively computed for all pairs of cells and thus ensures the appropriate weightage of nearby cells in
121 both high- and low-density states (**Supplementary Fig 3B-F**). Finally, Mellon adopts a scalable Bayesian
122 inference approach, tailored for atlas-scale single-cell datasets. The scalability is in large part achieved
123 by employing a sparse Gaussian Process that approximates the full covariance structure using a set of
124 landmark points (**Methods**).

125 The density function derived by Mellon is a continuous representation of the single-cell phenotypic
126 landscape (**Fig. 1E, Supplementary Fig. 1D**). This function enables density estimation at single-cell
127 resolution (**Fig. 1F**). Visualizing cell-state densities with methods such as UMAPs (**Fig. 1F**) simplifies the
128 exploration of high- and low-density cell states in differentiation landscapes. Within the cell-state density
129 landscape, we discerned what we term *regions* – connected subsets within the cell-state space with
130 similar density characteristics. Such regions represent collections of closely related cell states that cells
131 appear to inhabit (*high-density* regions) or traverse (*low-density* regions) during their differentiation
132 journey (**Fig. 1F**).

133 To assess Mellon's accuracy, we generated simulated datasets composed of either discrete clusters or
134 continuous trajectories, using mixtures of Gaussians in ten to twenty dimensions (**Supplementary Fig.**
135 **4A, D, G**). Comparing the ground truth density from the Gaussian mixtures to Mellon-inferred density
136 demonstrated strong agreement, showcasing Mellon's capability to accurately estimate cell-state
137 densities in high-dimensional spaces (**Supplementary Fig. 4**).

138

139 Density landscape of hematopoiesis with Mellon

140 Hematopoiesis, the process through which the blood and immune cells differentiate from hematopoietic
141 stem cells (HSCs), provides an ideal paradigm to understand and model differentiation²⁰. We therefore
142 utilized a previously generated single-cell multiome dataset of T-cell depleted bone marrow²¹
143 representing human hematopoietic differentiation (**Fig. 2A**) to evaluate the performance of Mellon and
144 interpret the inferred cell-state densities.

145 We used diffusion maps to derive a representation of hematopoietic cell-states and applied Mellon to
146 infer density in this high-dimensional cell-state space (**Fig. 2B**). The resulting density landscape exhibited
147 considerable heterogeneity, with numerous high-density regions, corresponding to major cell types,
148 interconnected by low-density regions indicative of rare transitory cells (**Fig. 2B**). Monocytes, for instance,
149 exhibited the highest cell-state density (**Fig. 2C**), which is consistent with their status as the most
150 prevalent leukocyte in hematopoiesis and their emergence from bone marrow in a naïve state²².
151 Intriguingly, we observed noticeable fluctuation in density within several cell-type clusters, suggesting an
152 inherent heterogeneity, a nuance often masked when cells are grouped together by clustering methods
153 (**Fig. 2C**).

154 For a more comprehensive understanding of the hematopoietic density landscape, we utilized our
155 trajectory detection algorithm Palantir¹⁴ to determine a pseudo-temporal ordering of cells representing
156 developmental progression and cell-fate propensities that quantify the probability of each cell to
157 differentiate to a terminal cell-type (**Supplementary Fig. 5A-B**). We compared cell-state density along
158 pseudotime for each lineage and observed that the increase in fate propensity towards the lineage is
159 strongly correlated with the first low-density region in each lineage (**Fig. 2D-E, Supplementary Fig. 5C-**

160 **D)**. Low-density regions therefore appear to be a hallmark of cell-fate specification in hematopoiesis.
161 These low-density regions from HSCs to fate-committed cells encompasses <0.4% of the data and under
162 <0.01% of bone marrow cells, demonstrating the ability of Mellon to identify low-frequency rare transitory
163 cells (**Fig. 2E**).

164 The occurrence of low-density regions in density landscapes can be attributed to accelerated gene
165 expression changes, divergence, or apoptosis (**Supplementary Fig. 1A-C**). Apoptosis during
166 hematopoietic cell-fate commitment has been shown to be minimal²³. Further, divergence or spread of
167 cell states, while theoretically possible, would likely result in a broader distribution rather than the
168 observed tight trajectories. Therefore, our results strongly suggest that hematopoietic lineage
169 specification events occur through low-density regions induced by rapid and accelerated gene expression
170 changes.

171 We next devised a gene change analysis procedure to identify genes with high expression change in
172 low-density regions (**Methods**). Our procedure consists of two steps: (1) We first determine local
173 variability for each gene, which represents the change in expression of the gene in a cell-state. Local
174 variability for a gene is determined as follows: For each state, we computed the absolute difference in
175 gene expression to its neighbors. The differences are normalized by distance between states and the
176 maximum of these normalized differences is nominated as the local variability of the gene. (2) Genes are
177 then ranked by the weighted average of local variability across cells spanning a low-density region and
178 the flanking high-density regions. Inverse of density are used as weights to ensure genes with higher
179 expression change in low-density regions are ranked higher. Thus, gene change analysis quantifies the
180 influence of a gene in driving state transitions in low-density regions (**Methods**).

181 We applied the gene change analysis procedure to identify genes that drive hematopoietic fate
182 specifications by selecting cells spanning hematopoietic stem-cells to fate committed cells along each
183 lineage (**Supplementary Fig. 6A**). Upregulated genes in each lineage transition were enriched for
184 lineage identity genes whereas downregulated genes across lineages were associated with stem cell
185 programs (**Supplementary Fig. 6B-C**), indicating that changes that underlie cell-fate specification in
186 hematopoiesis occur in low-density regions. Notably, we observed genes with higher expression levels
187 specifically in low-density states, suggesting that despite their transitory nature, certain gene regulatory
188 programs are uniquely adapted to facilitate these transitions (**Supplementary Fig. 6**).

189 We next utilized Mellon densities and associated genes to investigate B-cell fate specification. Genes
190 with high change scores in this low-density region were enriched for modulators of B-cell lineage
191 specification with their roles traversing transcriptional regulation, intracellular signaling and cell migration.
192 Transcription factor EBF1 had the highest change score (**Supplementary Fig. 7A**), aligning with its role
193 as the master regulator of B-cell differentiation²⁴. In fact, the upregulation of EBF1 is exquisitely localized
194 to the low-density region between stem and B-lineage committed cells (**Fig. 2F**), with similar dynamics
195 observed in other critical B-cell commitment regulators such as PAX5 and IL7R (**Supplementary Fig.**
196 **7B-C**). From a signaling point of view, we observed an upregulation of IL-7 responsive Stat signaling
197 targets in the same low-density cells concurrent with IL7R upregulation (**Fig. 2G, Supplementary Fig.**
198 **7C**). These observations are consistent with previous studies that have illustrated the vital role of IL-7
199 driven activation of STAT5 in a rare precursor population for B-cell specification¹. Finally, genes such as
200 NEGR1, with documented roles in cell adhesion and migration²⁵, also score high (**Supplementary Fig.**
201 **7B**), demonstrating that the spatio-temporal continuum of B-cell differentiation within the bone marrow is
202 executed as rapid transcriptional changes through low-density cell-states.

203 These findings underscore the potential of Mellon to uncover rare, biologically significant cell populations.
204 They also demonstrate that rapid transcriptional changes that drive state transitions in low-density

205 regions are shaped by an intricate interplay of cell-autonomous and extrinsic factors, highlighting how
206 Mellon can help unravel this complexity.

207 Following fate specification, B-cell development is a highly orchestrated process where cells transition
208 through checkpoints as they gain functional and non-self-reactive B-cell receptors²⁶. We analyzed Mellon
209 densities along pseudotime and observed that B-cell differentiation is defined by alternating high- and
210 low-density regions (**Fig. 2H**). Using marker gene expression and gene change score analysis, we
211 inferred that every high-density peak represents a well-characterized checkpoint, and every checkpoint
212 corresponds to a high-density peak (**Fig. 2H, Supplementary Fig 7D**). This also suggests that
213 checkpoint releases manifest as low-density states. Since apoptosis has only been extensively observed
214 in the transition from Pre-Pro B-cells to Pre-B-cells^{1,26}, our results suggest that cells rapidly change their
215 state upon checkpoint release until they reach the next checkpoint, where they converge to create high-
216 density regions.

217 As a test of robustness of these results, we assessed Mellon's reproducibility by computing cell-state
218 densities for single-cell datasets of bone marrow cells from eight independent donors from the Human
219 Cell Atlas²⁷. Densities were highly consistent across the donors, demonstrating consistent observation of
220 high- and low-density regions across the hematopoietic landscape (**Supplementary Fig. 8A-B**).
221 Moreover, density patterns along the B-cell differentiation trajectories were also consistent between the
222 samples, reinforcing the reliability and reproducibility of Mellon density estimates (**Supplementary Fig.**
223 **8C**).

224 Versatility of Mellon cell-state densities

225 We investigated whether cell-state density is a fundamental property of the homeostatic system by
226 investigating whether cell-state density is restored upon regeneration. We utilized a single-cell dataset of
227 lung regeneration where lungs were profiled using scRNA-seq following injuries induced with bleomycin
228 (**Fig. 2I**)²⁸. We applied Mellon to compute cell-state densities before injury and upon recovery.
229 Remarkably, we observed that the density landscape reverts to the homeostatic state upon regeneration
230 from injury (**Fig. 2J, Fig., Supplementary Fig. 9**). This observation suggests that cell-state density, while
231 fundamental to tissue homeostasis, is also reflective of the tissue regenerate state. As the tissue recovers
232 from injury, the restoration of the original cell-state density landscape could serve as an indicator of
233 successful tissue regeneration.

234 We further explored Mellon's versatility by applying it to a variety of homeostatic biological systems such
235 as pancreatic development²⁹, endoderm differentiation³⁰ and spatial organization of intestinal tissues³¹.
236 The recurring observation of high- and low-density regions across these diverse systems suggests that
237 these patterns are a ubiquitous feature of single-cell phenotypic landscapes (**Supplementary Fig. 10**).
238 These density variations supply a wealth of biological insight beyond abstract quantities: High-density
239 regions across systems typically correspond to key developmental checkpoints or bottlenecks, while low-
240 density regions often represent rare transitory cells undergoing rapid transcriptional changes to bridge
241 the denser areas (**Supplementary Fig. 10**).

242 These findings emphasize the effectiveness of Mellon for accurately characterizing differentiation
243 landscapes and highlight the importance of scrutinizing both high- and low-density regions for a holistic
244 understanding of the differentiation processes. Mellon's fine-grained resolution also aids the identification
245 of rare transitory cells, a critical element of diverse biological phenomena.

246 Mellon produces robust cell-state densities

247 We next assessed the robustness of Mellon cell-state densities across different parameters. The number
248 of cells measured in a dataset is a crucial factor affecting the accuracy and reliability of density estimates.

249 We performed subsampling experiments and compared the results to those obtained using the full
250 dataset by leveraging the continuous nature of Mellon. Our subsampling experiments show that Mellon's
251 density estimates are highly robust to subsampling, even when reducing the number of cells by an order
252 of magnitude across different datasets (**Supplementary Fig. 11,12**). Density estimates are also robust
253 to variations in the number of diffusion components (**Supplementary Fig. 13**), dimensionality
254 (**Supplementary Fig. 14**), the number of landmarks (**Supplementary Fig. 15**), and the length-scale
255 heuristic employed for scalability (**Supplementary Fig. 16**). These findings underscore the reliability of
256 Mellon's density estimation approach, which can provide accurate and robust results even with limited
257 data.

258
259 Finally, we compared Mellon to existing approaches for cell-state density estimation. Densities have been
260 approximated as the inverse of distance to k th nearest neighbor^{2,14} due to computational complexity.
261 However, due to the inherent noise and sparsity of scRNA-seq data, these approaches often fail to
262 generate robust density estimates (**Supplementary Fig. 2A-B**). The characteristic high- and low-density
263 regions identified by Mellon could not be demarcated by densities estimated solely from nearest neighbor
264 distances (**Supplementary Fig. 2B-D**). Given this noise, 2D embeddings, such as UMAPs, have been
265 widely utilized for density computation. While such embeddings are effective for visualization, the low-
266 dimensionality restricts their capacity to encapsulate all biologically significant variability. The UMAP
267 density estimates for the T-cell depleted bone marrow data are dominated by the most dominant cell-
268 type, i.e., monocytes (**Supplementary Fig. 2C**) with no discernable variability in the other lineages
269 (**Supplementary Fig. 2D**). Thus cell-state density estimation using Mellon substantially outperforms
270 existing approaches in accuracy, and biological interpretability.

271
272 **Enhancer priming as a catalyst for rapid transcriptional changes in low density cell-**
273 **states.**

274 We next turned our attention to the mechanisms that regulate the rapid transcriptional changes that
275 generate rare transitory cells during lineage specification. Previous studies have identified extensive
276 priming of lineage-specifying genes in hematopoietic stem cells, where gene loci are maintained in an
277 open chromatin state through pre-established enhancers, even in the absence of gene expression³²⁻³⁴.
278 Moreover, enhancer priming has been implicated to play a role in rapid transcriptional responses to stimuli
279 in hematopoietic cells³⁵.

280
281 Building on these findings, we hypothesize that the rapid upregulation of lineage specifying genes as
282 cells transition from HSCs (a high-density region) to fate committed cells (another high-density region) is
283 in part driven by enhancer priming. In this scenario, the loci of lineage-specifying genes are maintained
284 in an accessible state in HSCs through pre-established enhancers. A combination of cell-autonomous
285 and extrinsic factors trigger the upregulation of a small set of master regulators, which in turn rapidly
286 upregulate the expression of lineage-specifying genes in a coordinated manner. Thus, the combined
287 activity of pre-established enhancers in HSCs and lineage-specific enhancers established by master
288 regulators could produce the rapid transcriptional changes that underpin rare transitory cells in low-
289 density states (**Supplementary Fig. 17**).

290
291 We used the transition from HSC to B-cells as the case-study (**Fig. 3C**) to test our hypothesis. We
292 leveraged the multiomic nature of our T cell depleted bone marrow dataset, with measurements of both
293 expression (RNA) and chromatin accessibility (ATAC) in the same single cells (**Supplementary Fig. 18**).
294 The first step is to delineate the primed and lineage-specific peaks associated with a gene. The noise

295 and sparsity of scATAC data means that determination of individual peak accessibility at single-cell level
296 is extremely unreliable³⁶. Therefore, we devised a procedure to disentangle primed and lineage-specific
297 peaks associated with a gene using different levels of abstractions (**Supplementary Fig. 19A, Methods**):
298 First, we used our SEACells algorithm²¹ to aggregate highly-related cells into metacells and identified the
299 set of peaks with accessibility that significantly correlate with gene expression (**Fig. 3A**). We then
300 identified the subset of these peaks with greater accessibility in B-cells compared to other lineages by
301 comparing accessibility between cell-types at the metacell resolution. This approach ensures the
302 exclusion of ubiquitous and low-signal peaks while retaining peaks that are important for B-cell fate
303 specification. Finally, we classified each peak as primed if it was accessible in HSCs, and as lineage-
304 specific if its accessibility was B-cell restricted (**Fig. 3B**). We verified that the accessibility of lineage-
305 specific peaks was near exclusive to B-cells and that of primed peaks were higher in HSCs and B cells
306 (**Supplementary Fig. 19B-C**).

307
308 We identified the set of genes with high change scores in B-cell specification using our gene change
309 analysis procedure (**Supplementary Fig. 20A**). We then used the subset of these genes with
310 upregulation in B-cell lineage and those with at least 5 peaks correlated with expression to test our
311 hypothesis. >80% of these genes were associated with at least one primed peak and one lineage-specific
312 peak (**Supplementary Fig. 20A**), implicating enhancer priming as vital to their upregulation. In contrast,
313 none of the genes associated with the erythroid fate specification demonstrated B-cell primed peaks. To
314 characterize the dynamics of these peaks during lineage specification, we computed two accessibility
315 scores for each gene at single-cell resolution: (i) primed score, defined as the aggregated accessibility of
316 all primed peaks correlated with the gene and (ii) lineage-specific score, defined as the aggregated
317 accessibility of all lineage-specific peaks correlated with the gene (**Methods**). We first used these scores
318 to examine the dynamics of EBF1, the gene with the highest change score in the low-density region of
319 B-cell fate specification (**Fig. 3D**) and the master regulator of B-cell differentiation²⁴. We observed that
320 primed peaks were open in stem cells as expected and increased in accessibility as B-cell fate was
321 specified (**Fig. 3E**, orange line). This was followed by the establishment and stabilization of lineage-
322 specific peaks (**Fig. 3E**, blue line) and finally lineage-specific upregulation of EBF1, highlighting the role
323 played by enhancer priming in the upregulation of EBF1. We next examined the genes upregulated in B-
324 cell lineage specification with primed and lineage-specific peaks, and observed a similar pattern to EBF1,
325 along with a coordinated upregulation that follows EBF1 expression (**Fig. 3E-F, Supplementary Fig.**
326 **20C**). Finally, we used in-silico ChIP³⁷ to identify that almost every gene in our gene set is a predicted
327 target of either EBF1 or PAX5 (**Fig. 3H**), consistent with our hypothesis and the proposed role of EBF1
328 as a trigger for a “big-bang” of B-cell development³⁸.

329
330 Our results support a mechanism where enhancer priming and subsequent activation of master
331 regulators lead to a rapid and coordinated upregulation of genes, resulting in the emergence of rare
332 transitory cells that confer lineage identity. These results highlight the importance of taking cell-state
333 density into consideration for understanding gene regulatory networks that drive cell-fate specification.
334 Our approach to determine primed and lineage-specific accessibility scores for each gene utilizes the
335 history of peak establishment, a feature unaccounted for by most current techniques, which tend to
336 aggregate all peaks in proximity of a gene to derive a single gene score^{32,36} (**Supplementary Fig. 20D-**
337 **E**). Finally, the expression and accessibility trends were determined using Gaussian process with the
338 function estimator implemented in Mellon, highlighting another utility of the Mellon framework
339 (**Supplementary Fig. 21, Methods**).

340

341 Identification of master regulators with Mellon

342 While master regulators have been identified for several hematopoietic lineages, the mechanisms
343 controlling lineage-specific upregulation of these master regulators remain largely elusive.
344 To investigate whether the regulation of EBF1 could be clarified through cell-state density, we adapted
345 our approach to compute gene-change scores to rank the EBF1 correlated peaks by their accessibility
346 change score in the low-density region of B-cell specification (**Methods**). Interestingly, the top peak in
347 this analysis was almost exclusively accessible in the low-density region (**Fig. 3I**). We employed in silico
348 ChIP to identify transcription factors with a strong predicted signal to bind this peak and observed that
349 top 10 enriched motifs were exclusively comprised of IRF and SOX motifs. Interestingly, the increase in
350 accessibility in the peaks is concurrent with upregulation of the transcription factor SOX4 (**Fig. 3J**), a
351 known regulator of EBF1 during B-cell development³⁹. These results clarify the temporal order of
352 transcriptional events where upregulation of SOX4 leads to lineage-specific expression of EBF1 to
353 establish B-cell fate and also suggest the specific set of regulatory elements that drive this mechanism.

354
355 The strong association of EBF1 expression with low-density transition (**Supplementary Fig. 20A**) and
356 the high number of expression-correlated peaks (**Supplementary Fig. 19E**), coupled with the
357 coordinated upregulation of its targets (**Fig. 3H**), suggests a paradigm for identifying master regulators
358 via Mellon densities. Additionally, identifying peaks whose accessibility changes are strongly associated
359 with cell-state density can offer insights into the regulation of the master regulators themselves.

360 361 Exploring Time-Series Single-Cell Datasets with Mellon to Understand Mouse 362 Gastrulation

363 Time-series single-cell datasets are invaluable for understanding the intricate dynamic processes driving
364 development, as they provide snapshots of the changes in cell-type and cell-state compositions during a
365 fast-changing process. Although various computational methods exist to model trajectories using time-
366 series data^{8,40-43}, they typically represent these changes as discrete steps between measured timepoints
367 and thus are limiting when studying inherently continuous processes like embryonic development. To
368 better represent these processes, we investigated if we could utilize Mellon's continuous density functions
369 to construct a *time-continuous* cell-state density function. This function will span not just the observed
370 timepoints, but can also interpolate densities at unobserved times, enabling a truly continuous view of
371 the shifting cell-state density landscape during development.

372
373 We used the mouse gastrulation atlas⁴⁴, a single-cell dataset of 116,312 cells spanning gastrulation and
374 early organogenesis (E6.5-E8.5) (**Fig. 4A**) for exploration of time-continuous densities. We first applied
375 Mellon to each timepoint individually and observed considerable variability in cell-state densities over
376 time (**Fig. 4B, Supplementary Fig. 22**). Interestingly, we observed that the emergence of new cell types
377 or lineages was often marked by a low-density transition (**Fig. 4B**), highlighting the "fits and starts" nature
378 of developmental progression⁵.

379
380 Mellon's capacity to generate a continuous density function enables it to estimate densities for cell states
381 that were not part of the training data (**Fig. 1**). To demonstrate this, we used the density function
382 associated with each specific time point to calculate densities across cells of all timepoints (**Fig. 4C**). In
383 essence, we estimated the likelihood of each cell state being observed at a different time point. This
384 unique feature of Mellon allows for comparison of cell-state densities across various developmental
385 stages by calculating the correlation between the pair of time-point densities within the same cell state
386 (**Supplementary Fig. 23**). Interestingly, embryonic stage E7.75 was least similar to neighboring

387 timepoints, indicating the completion of gastrulation and onset of organogenesis at E7.75
388 (**Supplementary Fig. 23A-D**).

389
390 We next constructed a time-continuous cell-state density of mouse gastrulation by incorporating
391 measurement time as a covariate. We devised a procedure to ensure that the covariance of measurement
392 times between cells reflects the empirically observed correlation between timepoint densities (**Methods**,
393 **Supplementary Fig. 23E-F**) to construct a density function that is continuous in both *time* and *cell-state*.
394 Therefore, we can estimate cell-state densities at any desired timepoint situated between the measured
395 instances (**Fig. 4D, Supplementary Video 1**). Thus, by leveraging the temporally related data, we
396 enhanced the cell-state distribution of individual time points - a cell state present in preceding and
397 following time points is likely to exist in the current time point, even if it hasn't been directly observed.

398
399 To validate this approach, we performed leave-one-out cross-validation by comparing density computed
400 exclusively from a timepoint with the interpolated density computed by omitting the same timepoint. The
401 two densities are highly correlated even for timepoint E7.75 (**Supplementary Fig. 24**), which is least
402 similar to its neighbors, providing a clear validation of our approach.

403
404 Importantly, our time-continuous approach also enables the quantification of rates of density change. By
405 taking the derivative of the time-continuous density along the time axis, we can assess the rates of
406 enrichment or depletion for every cell state at any time (**Fig. 4D, Supplementary Video 1**). Our analysis
407 reveals that the initial phase of gastrulation is predominantly characterized by growth, with a nearly
408 constant abundance of epiblast and primitive streak cells—a finding in line with prior studies noting high
409 proliferation⁴⁵ (**Supplementary Fig. 25A**). Following this phase, a sharp transition occurs at E7.5, where
410 a rapid decline of epiblast and primitive streak cells signals the completion of the gastrulation process
411 (**Supplementary Fig. 25B-C**). Finally, another transition at E8.375 marks the emergence of ectodermal
412 and endodermal structures, accompanied by a concomitant decline in their respective progenitors
413 (**Supplementary Fig. 25D**). These findings underscore the power and potential of employing time-
414 continuous cell-state density modeling to provide a high-resolution depiction of the developmental
415 process in its entirety.

416
417 The application of time-continuous cell-state densities can also offer insight into the dynamics of cell
418 abundance along specific developmental lineages. As a case study, we chose to investigate
419 erythropoiesis during gastrulation, given its well-understood process. Using the full gastrulation atlas and
420 Palantir¹⁴ we identified cells predisposed to differentiate into erythroid lineage and derived a pseudo-
421 temporal ordering of these cells (**Supplementary Fig 26**). Leveraging our time-continuous cell-state
422 density function, we approximated densities along pseudo-time, which revealed a continuous progression
423 of cells toward the erythroid state (**Fig. 4F**). Interestingly, there is a strong alignment between pseudotime
424 and real time indicating a linear dependency in the erythroid lineage. Note that the persistent high density
425 of early epiblast cells likely represents cells differentiating into cell types other than the erythroid lineage.

426
427 Further, the dynamics of cell-type proportion along real time can be investigated by computing the
428 marginal of the density representation that contrasts real-time versus pseudo-time (**Fig. 4G**). This
429 visualization allowed us to precisely pinpoint the timespan during which hematoendothelial progenitors,
430 the earliest precursors of erythroid cells, emerge from the nascent mesoderm (**Fig. 4G**). Notably, the
431 proportion of hematoendothelial cells remains relatively stable across time, indicating their transient
432 presence without expansion in the cell population. In stark contrast, blood progenitor cells (Type 2)
433 undergo a substantial increase in their proportion following their emergence, suggesting a period of

434 accelerated cell division. Therefore, our time interpolation offers valuable insights into the progression of
435 cell type abundances and allows for high resolution predictions of the emergence of specific cell types.

436

437 Our results showcase Mellon's capability to provide a comprehensive, time-continuous perspective on
438 cell-state densities during development and reprogramming.

439

440 Mellon infers densities from single-cell chromatin data.

441 Single-cell chromatin profiling techniques such as single-cell ATAC-seq¹⁶, CUT&Tag^{46,47}, and sortChIC-
442 seq⁴⁸ are revolutionizing the study of interplay between gene expression and chromatin landscape in
443 disease and differentiation. We developed Mellon to be adaptable to different single-cell modalities,
444 making it a valuable addition to the computational toolkit for these emerging techniques. Given their
445 robust representation of cell-states, we use diffusion maps for deriving a cell-state space for density
446 inference through Mellon. Diffusion maps rely only on distances between similar cells and thus can be
447 constructed for most single-cell data modalities following appropriate pre-processing²¹.

448 To evaluate Mellon's adaptability to scATAC-seq data, we computed diffusion maps from the ATAC
449 modality of the T-cell depleted bone marrow dataset²¹ and applied Mellon for cell-state density inference.
450 Similar to gene expression, Mellon reveals substantial chromatin-state density variability
451 (**Supplementary Fig. 27A**). High- and low-density states corresponded respectively to major cell-types
452 or checkpoints and rare transitory cells (**Supplementary Fig. 27A**). Applied to a mouse model of lung
453 adenocarcinoma⁴⁹, we observed extensive chromatin-state density variability amongst cells of the
454 primary tumors, with the transition to metastasis associated with a sharp decrease in density
455 (**Supplementary Fig. 27B**). We made similar observations with a larger-scale scRNA-seq dataset of the
456 same mouse model (**Supplementary Fig. 27C**)⁹, consistent with previous studies which have
457 demonstrated that metastases are seeded by small group of cells⁵⁰.

458 While diffusion maps provide desirable properties for state representation from single-cell data, Mellon's
459 effectiveness is not tied to their specific properties. Mellon is capable of estimating densities in any
460 representation with a meaningful distance metric. To demonstrate this, we used Mellon to infer cell-state
461 densities from a MIRA representation⁵¹ of a multimodal dataset of skin differentiation (**Supplementary**
462 **Fig. 27D**). Similar to observations with single modality datasets, we observed extensive variability in
463 densities with low-density regions corresponding to exit from the stem-cell state and specification of
464 different lineages (**Supplementary Fig. 27D**).

465 We also tested Mellon's ability to recover chromatin-state densities using single-cell histone modification
466 data. We applied Mellon to compute densities using a single-cell sortChIC dataset of histone
467 modifications in mouse hematopoiesis⁴⁸ using H3K4me1, a histone modification that marks enhancers
468 and H3K9me3, that marks heterochromatin (**Fig. 5A-D**). H3K4me1 densities demonstrated extensive
469 heterogeneity similar to single-cell RNA and ATAC (**Fig. 5E**). On the other hand, H3K9me3 densities are
470 relatively uniform and lower compared to H3K4me1 (**Fig. 5F**). While the lower density is likely reflective
471 of the noise in heterochromatin marks which tend to occur in broad megabase size domains, the relative
472 uniformity is likely reflective of the underlying biology: Active chromatin marks like H3K4me1 accurately
473 have been shown to distinguish cell types and states whereas heterochromatin mark H3K9me3 struggles
474 to achieve the same resolution⁴⁸. This follows the function of H3K9me3 to aid in general repression of
475 other cell-fates rather than to actively establish cell-type identity⁴⁸. To quantify the difference in
476 heterogeneity between the two marks, we subsampled cells from each mark and compared the rank of
477 the covariance matrices (**Fig. 5D**). The covariance rank is a measure of information content where greater
478 the rank, higher is the complexity of the system. The distribution of ranks is significantly higher for

479 H3K4me1 compared to H3K9me3 (p-value < 1e-30, Wilcoxon rank-sum test) demonstrating a greater
480 complexity across cell-types for the H3K4me1 histone modification.

481 Our results demonstrate the versatility of Mellon with diverse single-cell data modalities and data
482 representations. Mellon's ability to robustly and accurately identify cell-state densities from single-cell
483 chromatin data suggests a key utility in mechanistic investigations with emerging technologies that
484 concurrently measure active and repressive modifications^{52,53}.

485 Highly efficient and scalable: Mellon's power in atlas-scale single-cell analysis

486 There is a growing trend towards generation of atlas-scale datasets that profile millions of cells, as well
487 as integration of smaller datasets into large-scale data repositories^{54,55}. To enable density computation
488 in these massive datasets, Mellon incorporates several features that enable efficient scalability: First,
489 Mellon uses a sparse Gaussian process, leveraging landmark points to approximate the covariance
490 matrix, facilitating the efficient handling of high-dimensional data, and reducing the computational
491 overhead associated with large datasets. Second, Mellon requires a single computation of the covariance
492 matrix, removing the need for continuous updates in every iteration and thus improving computational
493 efficiency. Finally, Mellon is built on the JAX python library, which is well-known for its high-performance
494 computing capabilities⁵⁶. The utilization of JAX allows Mellon to optimize available hardware resources,
495 further enhancing its scalability and computational efficiency.

496 Mellon's architecture is designed to scale near linearly in time and memory requirements i.e., the runtime
497 grows proportionally with the number of cells when the number of landmarks is kept constant (**Fig. 6**). To
498 demonstrate the scalability, we used the T-cell depleted bone marrow (8.6k cells)²¹, CD34+ bone marrow
499 (6.8k cells)²¹, mouse gastrulation (116k)⁴⁴ and the iPS reprogramming dataset (250k cells)⁸ spanning
500 datasets of different sizes and characteristics. For example, using a single CPU core and a default of
501 5000 landmarks, Mellon required only 100 seconds to process 10k cells and 17 minutes to process 100k
502 cells of the iPS dataset, including the computation of diffusion maps (**Fig. 6A**). Additionally, Mellon
503 benefits from parallelization, enabling even faster processing times (**Fig.6B**). This highlights Mellon's
504 efficiency in handling datasets of different sizes.

505 To further evaluate Mellon's scalability, we performed benchmarking on simulated datasets, where we
506 utilized a single CPU core and excluded the time for diffusion map computation. Mellon demonstrated its
507 capability to handle large-scale datasets by accurately computing densities on a dataset of ~6 million
508 simulated cells in less than 12 hours (**Supplementary Fig. 28A-D**). In addition, we investigated Mellon's
509 scalability with reduced numbers of landmarks. When the number of landmarks was decreased to 1000,
510 Mellon maintained its accuracy while requiring less than 7 hours to accurately estimate densities for a
511 simulated dataset of around 10 million cells (**Supplementary Fig. 28E-H**). These results highlight
512 Mellon's remarkable scalability to tackle the burgeoning demands of increasingly large single-cell
513 datasets.

514 **Discussion**

515 Rapid transcriptional changes that lead to rare transitory cells and thus induce differences in cell-state
516 density have been well-documented as a fundamental property of developmental systems from plants to
517 mammals⁵. Single-cell studies have reinforced the critical nature of rare transitory cells in diverse
518 biological contexts such as development^{42,44}, differentiation¹, reprogramming⁸, plasticity of tumors¹⁰ and
519 metastasis⁹. However, existing approaches for estimating cell-state densities have fundamental
520 limitations: They either rely on noisy neighborhood-based estimates around individual cells or utilize 2D
521 dimensional embeddings that do not capture the full complexity of cell-states. Mellon addresses this gap
522 by providing a robust and accurate framework for estimating cell-state densities from high-dimensional
523 cell-state representations. Mellon can be applied to dissect the density landscapes not only in
524 differentiation and development but also during reprogramming, regeneration, and disease. We extended
525 the Mellon framework to estimate time-continuous cell-state density for temporal interpolation of time-
526 series data. The computational efficiency of Mellon allows for rapid density computations, enabling the
527 analysis of large-scale single-cell datasets containing hundreds of thousands of cells within minutes.
528 Furthermore, Mellon's flexibility supports density estimation for diverse single-cell data modalities,
529 making it a versatile tool for investigating cell-state densities across various biological systems.

530 Mellon's innovative approach involves formalizing the connection between density and nearest neighbor
531 distances using a Poisson process and establishing a link between cell-state similarity and density
532 through Gaussian processes. This unique combination overcomes computational challenges in high-
533 dimensional spaces and enhances the robustness and accuracy of density estimation. The scalability of
534 Mellon is achieved through the utilization of sparse Gaussian processes, heuristic for length-scale
535 optimization to avoid redundant computations of the covariance matrix, and implementation using
536 efficient JAX libraries.

537 Our work underscores the significance of cell-state density in understanding differentiation trajectories
538 and the potential of Mellon to provide new insights into the regulatory mechanisms guiding cell-fate
539 decisions. We have demonstrated the effectiveness of Mellon in estimating cell-state density using
540 hematopoietic differentiation. Mellon's ability to capture the heterogeneous density landscapes, where
541 high-density regions correspond to major cell-types and low-density regions represent rare transitory
542 cells, is particularly evident. By incorporating our trajectory detection algorithm, Palantir, we have been
543 able to observe a correlation between low-density regions and lineage specification. Further, our gene
544 change analysis procedure helps identify gene expression changes that drive low-density transitions and
545 can help elucidate the underlying molecular mechanisms. This was particularly insightful during our
546 investigation into B-cell fate specification, where the detection of low-density regions played a central role
547 in identifying the importance of enhancer priming and characterizing the regulation of the master regulator
548 EBF1. The pattern of alternating high- and low-density regions, observed during the process of B-cell
549 development, further highlighted the dynamic nature of differentiation. Importantly, the consistency of our
550 density estimates across independent donor samples highlights reproducibility and reliability of Mellon.
551 Therefore, these findings provide a strong foundation for further exploration into the intricacies of cellular
552 differentiation.

553 An important consideration for estimation of cell-state density is the inherent dimensionality of the cell-
554 state space. Mellon by default uses the dimensionality of the cell-state space i.e., number of diffusion
555 components for density estimation, but the intrinsic dimensionality is likely substantially lower. In other
556 words, not all diffusion components are relevant to describe any given region or point in the cell-state
557 space. With measures of intrinsic dimensionality⁵⁷, one can produce explicit units of density and make
558 statements about how many more cells per volume can be expected at a given state. High-dimensionality

559 of the state-space also presents a challenge for automatic determination of high- and low-density regions.
560 Therefore, we compared densities with lower-dimensional projections such as pseudo-time to identify
561 such regions. Incorporation of density as a feature for clustering algorithms or the use of local context
562 density could lead to direct computation of such regions in the high-dimensional state-space.

563 We anticipate that the time-continuous cell-state densities for temporal interpolation will be a powerful
564 addition to the computational toolkit for modeling cell-state dynamics using time-series single-cell
565 datasets. Mellon provides capabilities to interpolate cell-state density and since the density function is
566 differentiable, it also supports the computation of density change at all times between measured time
567 points. Thus Mellon densities can serve as inputs for development of computational algorithms leveraging
568 advances in optimal transport⁴³ for a high-resolution characterization of cell-fate choices using time-series
569 data.

570 We have demonstrated that cell-state density is a fundamental property of the differentiation landscape
571 by observing that homeostatic density is re-established upon lung regeneration (**Fig. 2I-J**). Thus, the
572 Mellon cell-state density function can itself serve as a phenotype of that differentiation landscape that is
573 altered upon perturbation. Single-cell datasets in unperturbed and perturbed conditions can be jointly
574 embedded into a common state space such as diffusion maps and density functions can be computed
575 separately for each condition in the common space. Comparison of densities from different conditions
576 can not only provide estimates of differential abundance at unprecedented resolution but can also be
577 utilized to develop summary statistics that describe and quantify the nature of the perturbation across the
578 entire differentiation landscape.

579 Fundamentally, the density function estimated by Mellon provides a comprehensive description of the
580 differentiation landscape, representing the probability distribution of cells within different states. Unlike
581 many existing approaches that rely solely on the measured cell-states and number of cells, which can
582 introduce technical biases and impact the interpretation, Mellon's density function reflects the inherent
583 complexity of the biological system. As the number of measured cells increases, the density function
584 converges in complexity, allowing for a more accurate representation of the relative abundances of all
585 possible cell states. Mellon can be extended to support online learning, enabling the incremental
586 refinement of the density function with new data. Monte Carlo sampling approaches can leverage
587 Mellon's cell-state density function to generate synthetic cell-state data, which can greatly enhance data-
588 intensive machine learning models. By incorporating the richness of the density function, these synthetic
589 data can augment training sets and improve the performance and robustness of downstream analyses.
590 Further, the differentiability of Mellon's density function opens up possibilities for utilization of partial
591 differential equations. This enables the modeling of the differentiation process as a dynamical system
592 and facilitates the inference of regulatory dynamics underlying cellular transitions. By integrating Mellon's
593 density function within differential equation frameworks, one can gain deeper insights into the regulatory
594 mechanisms governing cellular differentiation and uncover key factors driving the dynamic processes.

595

596 **Data Availability**

597 All datasets used in the manuscript have been previously published and the accession numbers are listed
598 in **Supplementary Table 1**. Mellon results and cell-type metadata information for the T-cell depleted
599 bone marrow and the mouse gastrulation data are available on Zenodo at
600 <https://doi.org/10.5281/zenodo.8118722>.

601

602 **Code Availability**

603 Mellon is available as a Python module at <https://github.com/settylab/Mellon>. Jupyter notebooks detailing
604 the usage of Mellon including cell-state density estimation, gene change computation, time-continuous
605 cell-state density estimation, and enhancer classification are available at
606 <https://mellon.readthedocs.io/en/latest/>. Pipelines for running SEACells, computing gene-peak
607 correlations, primed and lineage-specific accessibility scores are available at
608 https://github.com/settylab/atac_metacell_utilities.

609

610 **Author Contributions**

611 D. O. and M. S. conceived and designed the study, developed Mellon, developed additional analysis
612 methods and statistical tests. D. O. and B. D. developed the heuristics, performed robustness analyses
613 and implemented the framework. C. J. and M. S. performed analysis of enhancer dynamics. C. D.
614 supported enhancer dynamics analysis. D. O., C. J. and M. S. wrote the manuscript.

615

616 **Acknowledgements**

617 We thank members of the Setty lab for discussions and comments on the manuscript. This study was
618 supported by National Institute of General Medical Studies grant R35 GM147125 and Brotman Baty
619 Institute Pilot Award to MS; National Institutes of Health grant ORIP S10OD028685 to support high-
620 performance computing at the Fred Hutchinson Cancer Research Center.

621

622 **Competing Interests**

623 The authors declare no competing interests.

624

625

626 **Methods**

627 **Mellon Algorithm**

628 Mellon is a computational tool designed to infer cell-state densities from high-dimensional single-cell data.
629 The objective of Mellon is to characterize the complex density landscapes of single-cell data (**Fig. 1A-B**)
630 with density estimates that are robust even in low-density regions, while maintaining computational
631 efficiency.

632 Mellon's computational model is grounded on two core assumptions. Firstly, within the chosen
633 representation of cell states, smaller distances between cell states signify higher biological similarity. In
634 other words, we assume that biological dissimilarity can be effectively quantified by the Euclidean
635 distance within this representation. Secondly, we assume that cell-to-cell density changes are smooth
636 and continuous, meaning that cell states of high similarity are expected to have similar state densities.

637 The input to Mellon is a high-dimensional representation of the cell-states (e.g., Diffusion maps). The
638 Euclidean distance between these cell-states serves as a measure of biological dissimilarity. Mellon
639 outputs a *continuous density function* that allows evaluation of cell-state densities at single-cell resolution
640 (**Fig. 1E-F**). The densities are computed in the high-dimensional cell-state space and visualized using
641 low-dimensional embedding techniques such as UMAPs.

642 The Mellon framework contains the following major components:

- 643 • Mellon first calculates the distance to the nearest neighbor for each cell in the cell-state space,
644 following the first assumption.
- 645 • The distances are linked to density via the *Nearest-Neighbor Distribution* (**Fig. 1C**).
- 646 • Densities between highly related cell-states are connected by the *Gaussian Process* and the
647 associated kernel function (**Fig. 1C-D**).
- 648 • A *Bayesian Model* (**Fig. 1D**) is deployed, integrating the nearest-neighbor distribution, kernel
649 function, and Gaussian Process to compute the continuous cell-state density function (**Fig. 1E**).

650 We next describe each of these components in detail along with our approach to scale Mellon for large
651 datasets.

652

653 **Nearest-Neighbor Distribution**

654 The core principle of Mellon relies on the relationship between nearest-neighbor distances and density,
655 as depicted in **Supplementary Figure 3**. This connection can be formalized using a Poisson point
656 process to define a nearest-neighbor *distribution*, which describes the probability of another cell-state
657 existing within some distance of a reference cell-state. Intuitively, regions with a higher density of cell-
658 states correspond to tighter nearest-neighbor distributions, while low-density regions result in broader
659 distributions (**Fig. 1C**).

660 For distance r and density ρ , the probability density function of the Nearest-Neighbor distribution
661 $f_{\text{NN}}: \mathbb{R}^d \rightarrow \mathbb{R}^+$ is given by

$$662 \quad f_{\text{NN}}(r|\rho) = \exp(-\rho \cdot b(r, d)) \cdot \rho \frac{db(r, d)}{dr}$$

663 where $b(r, d)$ is the volume of a d -dimensional ball with radius r . For a cell-state $x \in \mathbb{R}^d$ with a nearest
664 neighbor distance $\text{dn}(x)$, this probability density function gives rise to the following maximum likelihood

665 estimate for density if no prior is employed, formalizing the inverse relationship between nearest neighbor
666 distances and density as

667
$$\hat{\rho}(\text{dn}(x)|d) = \frac{(d-1) \cdot \Gamma\left(\frac{d}{2} + 1\right)}{d \cdot \text{dn}(x)^d \cdot \pi^{\frac{d}{2}}}.$$

668 The derivation is detailed in **Supplementary Note 1**. The use of Poisson point process is facilitated by
669 the two key assumptions of Mellon: The use of Euclidean distance in the cell-state space is a critical
670 requirement for defining the probability density function. The second is the assumption of smoothness in
671 cell-to-cell density changes. This is crucial as it allows us to assume that the density at a given cell-state
672 corresponds to the average density within a sphere centered at that state, with the radius of the sphere
673 defined by the nearest neighbor distance.

674 Gaussian Process

675 Building upon the foundational connection between nearest neighbor distance and density, Mellon utilizes
676 Gaussian process (GP) priors to establish a relation between the densities of highly-similar cell-states,
677 facilitating a continuous density function estimation. Similarity between cell-states is encoded using the
678 covariance function of the Gaussian process. The random variable of the GP, denoted as $f(x)$, serves
679 as the approximation of the logarithm of the cell-state density. Two properties of GPs make them ideally
680 suitable for cell-state density estimation from single-cell data: (i) GPs can be used to describe arbitrarily
681 complex functional spaces where the true functional form is unknown and (ii) GPs provide robust
682 estimates even when small number of observations are available.

684 The GP is defined as follows:

685
$$f(x) \sim \text{GP}(m, \text{Matern52}(l))$$

686 where m and $\text{Matern52}(l)$ are the mean function and the Matern covariance function respectively. A more
687 detailed assessment is provided in **Supplementary Note 2**.

688 Mean function

689 The true log-cell-state density approaches negative infinity away from any observed cell state. However,
690 functions sampled by the Gaussian process approach the chosen mean. To approximate the true
691 behavior of density functions we choose a very small value for the mean m that implies a vanishingly
692 small probability for a distant cell state. This mean function is given by the constant:

693
$$m = P_{1\%}[\hat{\rho}(\text{dn}(x_i)|d)_{i \in \{1, \dots, n\}}] - 10$$

694 where $P_{1\%}[\cdot]$ is the 1st percentile of the given data, $\hat{\rho}$ is the heuristic maximum likelihood estimate for
695 density, and $\text{dn}(x_i)$ is Nearest-Neighbor Distance of cell-state x_i in \mathbb{R}^d . The choice of this mean is
696 discussed in **Supplementary Note 3**.

697 Covariance Function and length scale

698 Similarities between cell-states are encoded through the GP covariance function or kernel. Specifically,
699 the kernel function defines the covariate structure between cell-states which translates to the smoothness
700 of the density function. Some commonly used kernels are arbitrarily smooth and allow arbitrary
701 differentiability. Assuming such smoothness can, however, lead to unrealistic results¹⁹. We therefore
702 chose to use the Matern covariance function with $\nu = \frac{5}{2}$ as the kernel, which is exactly twice differentiable

703 and thus constrains the degree of smoothness of the density function. The Matern52 kernel for a pair of
704 cell states $x, y \in \mathbb{R}^d$ is defined as:

705
$$\text{Matern52}(l)(x, y) = \left(1 + \frac{\sqrt{5} \|x - y\|_2}{l} + \frac{5 \|x - y\|_2^2}{3l^2} \right) \exp\left(-\frac{\sqrt{5} \|x - y\|_2}{l} \right)$$

706 The covariate structure of the Gaussian process is governed by the length scale parameter, denoted as
707 l , which essentially determines the radius of influence around each cell state. Conceptually, the length
708 scale sets the reach of influence for each cell, defining the range within which other cells contribute to
709 the local density estimate (**Supplementary Fig. 3**). In areas of lower density, fewer but more
710 representative cells influence the density estimate, while in higher density areas, a larger number of cells
711 contribute. This scenario gives rise to an effective number of neighbors that is density-dependent, which
712 is a direct result of the distance-mediated impact on the local density estimate (**Supplementary Fig. 3**).

713 This method not only increases the reliability and robustness of density estimates, but it also enables the
714 creation of a continuously changing density function between cell states, offering a nuanced
715 representation of biological phenomena. Unlike the k-nearest neighbor methods for density estimation
716 that assign an equal weight to all k neighbors irrespective of their distances, the continuous covariance
717 function of the Gaussian process accounts for the distance between cells, smoothly adjusting the weight
718 of their contribution. The resulting impact on the local density estimate facilitates a more precise
719 representation of the cell-state landscape.

720 The ideal length scale strongly depends on the availability of data at different points of the cell-state
721 space and encompasses a specific amount of cells needed to support a reliable density estimate of a
722 given state. We therefore derived a heuristic for length scale as function of the mean nearest neighbor
723 distance between cells:

724
$$l = \exp\left(\lambda + \frac{1}{n} \sum_{j=1}^n \log \circ \text{dn}(x_j) \right)$$

725 Here, $\lambda = 3$ is a heuristic value inferred from an extensive cross-analysis of multiple datasets. The
726 derivation of the length-scale heuristic is described in **Supplementary Note 4**.

727

728 **Sparse Gaussian Process**

729 Gaussian process computation usually necessitates $O(n^3)$ operations, where n is the number of cells
730 and thus can be prohibitively expensive for large datasets. To address this computational challenge,
731 Mellon utilizes a sparse approximation of the GP. This approach substantially reduces the computational
732 complexity while maintaining the versatility and expressiveness of the full GP model.

733 The sparse GP in Mellon is constructed using a subset of data points, referred to as “landmark cell-
734 states,” that essentially act as inducing points. These landmark states are chosen to capture the essential
735 structure of the cell-state space, providing a representative skeleton for the full GP model. This sparse
736 GP approach translates to an efficient $O(nk^2)$ time complexity for inference, where k is the number of
737 landmark points, a substantial reduction from the cubic time complexity of the full GP.

738 The specifics of the sparse GP implementation play a crucial role in the overall performance of the Mellon
739 and are described in **Supplementary Note 2**.

740

741 *Landmark selection*

742 The choice of landmarks, akin to the “inducing points” in a Gaussian process, is essential to ensure
743 precise recovery of the approximated covariance structure. Previous studies have demonstrated that k-
744 means centroids are well suited for this purpose⁵⁸. We assessed the accuracy of this approach by
745 comparing the inferred density derived from the landmarks against the density function inferred from a
746 non-sparse, or “no-landmarks” version (**Supplementary Fig. 15**). This comparison showed a
747 convergence of the landmark-based model towards the non-sparse version, thereby confirming the
748 efficacy of the landmark selection. We therefore use k-means clustering as the default landmark selection
749 in our algorithm and initialize it with `kmeans++`⁵⁹ to ensure computational efficiency.

750

751 Full Bayesian Model

752 The full Bayesian model used in the Mellon algorithm is formally defined as follows:

$$\begin{aligned} X &= (x_i)_{i \in \{1 \dots n\}}, x_i \in \mathbb{R}^{d'} \\ l &= \exp \left(\lambda + \frac{1}{n} \sum_{j=1}^n \log \circ \text{dn}(x_j) \right) \\ m &= P_{1\%}[\hat{\rho}(\text{dn}(x_i)|d)_{i \in 1 \dots n}] - 10 \\ f(x_i) &\sim \text{GP}(m, \text{Matern52}(l)) \\ \rho(x_i) &= \exp \circ f(x_i) \\ \text{dn}(x_i) &\sim \text{NN}(\rho(x_i), d) \end{aligned}$$

753

754 Where

- 755 • X represents the cell states, where each cell state, x_i , is a vector in the d' -dimensional Euclidean
756 space. The cell states form the primary input data for the model.
- 757 • l is the length scale of the GP covariance function. l is calculated from the distances to the
758 nearest neighbors in the cell-state space. The logarithm of these distances is averaged and added
759 to a fixed parameter $\lambda = 3$. The sum is then exponentiated to produce the length scale.
- 760 • $\hat{\rho}$ is the heuristic maximum likelihood estimate for the density.
- 761 • $\text{dn}(x_i)$ is Nearest-Neighbor Distance of cell state x_i
- 762 • m is the GP mean function. m is calculated as the 1% percentile of the heuristic maximum
763 likelihood estimates of density subtracted by a constant (10 in this case). This mean function
764 represents the average behavior of the underlying cell-state densities.
- 765 • Matern52 is the Matern covariance function with $\nu = \frac{5}{2}$ and length scale l .
- 766 • $f(x_i)$ is a random function generated by a sparse Gaussian process (GP), where x_i is the input
767 cell-state vector. The GP is defined by the mean function m and the Matern covariance function.

768 The cell-state density function $\rho(x_i)$ is the random variable of interest and is calculated by exponentiating
769 the function $f(x_i)$. This ensures that the density is always positive. The final part of the model is the
770 Nearest Neighbor Distance distribution NN of the Nearest Neighbor distance $\text{dn}(x_i)$, which is calculated
771 as a function of the cell-state density $\rho(x_i)$ and the dimensionality d .

772

773 Initialization

774 An appropriate initialization y' for the density function $y = \rho(x_i)$ can improve the convergence of the
775 maximum a posterior estimation. We employ a regression approach to initialize density estimation using
776 the heuristic maximum likelihood estimates of the log-density $\hat{\rho}$ (**Supplementary Note 1**):

$$777 \quad y' = \operatorname{argmin}_y \|\hat{\rho} - Ly\|_2^2 + \|y\|_2^2$$

778 Where L represents the transformation matrix within the Gaussian process, facilitating the conversion of
779 the latent representation, y , into the log-density function. High values in y are penalized through a ridge
780 regression to simulate the additional smoothness of the true density over $\hat{\rho}$.

781

782 Density at single-cell resolution

783 The log-density function $f(x_i)$ is evaluated at each single cell $x_i \in \mathbb{R}^{d'}$, to estimate log cell-state density
784 at single-cell resolution. The estimated densities in cell-state space are visualized using techniques such
785 as UMAPs for convenience. Note that the density function can be evaluated at any point in the cell-state
786 space including states that are not measured in the dataset. Single-cell densities can be examined along
787 pseudo-time, individual diffusion components or between interconnected clusters to identify high- and
788 low-density regions.

789

790 Note on the number of landmarks for sparse Gaussian Process

791 The number of landmarks serves as a parameter to the sparse Gaussian Process within Mellon. It's
792 crucial that the number of landmarks is sufficiently large to accurately capture the intricate patterns and
793 variability within the cell state density function. However, it is important to consider the trade-off involved:
794 an increased number of landmarks enhances the model's capacity to encapsulate finer details, but it also
795 increases the computational demands.

796 Mellon employs a default selection of 5,000 landmarks, an empirical decision grounded in extensive
797 testing with a multitude of datasets with different properties (**Supplementary Table 1**). Our evaluation
798 underscores the robustness of Mellon's density estimates across all investigated datasets, consistently
799 demonstrating stability even when the number of landmarks is substantially altered (**Supplementary Fig.**
800 **15**).

801 Nevertheless, the optimal number of landmarks can be contingent on the complexity and volume of the
802 particular dataset under examination. To assist users in selecting a representative number of landmarks,
803 Mellon incorporates a test for approximating the rank of the covariance matrix. Should the complexity of
804 the function appear exhausted using the existing landmark skeleton, a warning will be issued. This serves
805 as an indication that the selected number of landmarks might be insufficient for the model to accurately
806 capture the density function of the cell-states.

807

808

809 **Scalability of Mellon**

810 The implementation of Mellon leverages modern advances in numerical computation libraries, specifically
811 the JAX library, to enable efficient calculations and seamless differentiation. JAX is particularly suited for
812 our purposes due to its unique capability of just-in-time (JIT) compilation using XLA (Accelerated Linear
813 Algebra), a linear algebra compiler developed by Google⁵⁶. This feature ensures efficient utilization of
814 hardware resources, especially for large-scale computations and vectorized operations, which are
815 intrinsic to our method.

816 Mellon's scalability to large single-cell datasets is ensured through the use of a Sparse Gaussian Process
817 (GP). The sparse GP allows us to approximate the full GP model, significantly reducing the computational
818 demands while retaining the essence of GP's expressiveness. This scalability (**Fig. 6**) makes Mellon
819 practical for atlas-scale single-cell data sets, which often involve millions of cells.

820 Finally, model tractability in Mellon is achieved through the adoption of a length-scale heuristic for the GP
821 covariance function. The covariance function is crucial in a Gaussian Process as it dictates how many
822 nearby points in the input space influence each other in the output space. Typically, the length scale of
823 this function is subject to inference or optimization, often involving computationally intensive iterative
824 processes that require repeated updates of the covariance matrix and its Cholesky decomposition. In
825 Mellon, we sidestep this computational demand by deriving an appropriate length scale with a data driven
826 approach designed to adapt to the varying local densities in the high-dimensional cell-state space
827 (**Supplementary Note 4**). This not only streamlines the computation but also assists in avoiding
828 overfitting to dense regions, resulting in a smooth and accurate portrayal of cell-state density
829 relationships.

830 Together, these components create a balance between computational efficiency and model
831 expressiveness, making Mellon an effective and practical tool for cell-state density estimation from large,
832 high-dimensional single-cell data.

833

834 **Inference**

835 Mellon, by default, employs the L-BFGS-B optimization algorithm to infer the maximum a posteriori (MAP)
836 estimates of the posterior likelihood. Notably, our implementation provides direct access to the posterior
837 distribution of the density function. This is realized through a JAX function with automatic differentiation,
838 thus facilitating the use of any preferred inference scheme while retaining computational simplicity.

839 This flexibility is crucial because, in Bayesian inference, the MAP estimate can be subject to the
840 transformation of the latent representation and might not necessarily represent the "true" underlying cell-
841 state density. In fact, empirical evidence (**Supplementary Fig. 30**) indicates that the MAP estimate
842 strongly coincides with the posterior mean. However, without a definitive ground truth, it is challenging to
843 ascertain which estimate more closely resembles the true cell-state density.

844 In essence, Mellon's versatile implementation provides a robust framework for density estimation that
845 can adapt to diverse inference schemes, offering users the freedom to employ the technique best suited
846 to their specific study.

847

848

849 **Cell-state Representation**

850 Mellon utilizes diffusion components¹¹, as implemented in Palantir¹⁴, as the representation of cell-state
851 space. Diffusion maps have been widely used in single-cell data analysis owing to their reliable and
852 robust representation of cell-states^{12,14}. Cellular states in phenotypic landscapes reside in substantially
853 lower dimensions compared to measured gene expression owing to gene regulatory networks inducing
854 a strong covariate structure amongst genes. Therefore, biological similarity between cell states is more
855 closely linked to the distance they can traverse along the phenotypic landscape, rather than solely their
856 direct proximity in gene expression space. Diffusion maps identify the intrinsic structure in single-cell
857 data, mitigating noise by treating the data as realizations of a stochastic process. They not only efficiently
858 reduce noise in single-cell data but also extract a faithful representation of the underlying cell-state
859 manifold.

860 Further, the distances computed using diffusion maps, termed diffusion distances, are a measure that
861 reflects the interconnectedness of data points along the phenotypic manifold. Importantly, diffusion
862 distance operates along this manifold, which is constructed from the observed cell states, thereby
863 providing a meaningful indicator of biological similarity between cells. Therefore, the use of diffusion
864 distance in the estimation of cell-state density leads to a biologically relevant quantification of cells sharing
865 a similar state.

866 Diffusion maps can be constructed for different single-cell data modalities with appropriate preprocessing.
867 We recommend the use of PCA for RNA and SVD for ATAC and histone modification data. “Data
868 preprocessing” section provides more details on preprocessing of single-cell datasets. Diffusion maps
869 can also be constructed using other latent representations³⁶ or multimodal representations⁵¹.

870 **Number of Diffusion Components**

871 The dimensionality of the subspace where the data is represented is determined by the number of
872 diffusion components utilized in Mellon. Mellon results are robust to the number of diffusion components
873 indicating that pinpoint precision in their selection isn't strictly necessary (**Supplementary Fig. 13**).
874 However, some considerations are important while choosing the number of diffusion components:
875 Selecting a high number of diffusion components might lead to the inclusion of unnecessary noise within
876 the state representation, reducing the granularity of the resulting density model. Conversely, choosing a
877 small number of diffusion components might under-represent the complexity of the data, thus also leading
878 to a less detailed density model. The optimal number of diffusion components is therefore largely data-
879 specific and should be chosen to best capture the inherent structure and complexity of the cell data,
880 without unnecessarily increasing noise or forfeiting essential information. For example, Eigen gap statistic
881 has been previously employed to choose the number of diffusion components¹⁴.

883

884

885 **Genes Driving Low-Density Cell-State Transitions**

886 Low-density regions representing rare transitory cells are critical for diverse biological processes. We
887 devised a gene change analysis procedure to identify genes that drive cell-state transitions in low-density
888 regions and thus can be used to describe the dynamic behavior of the biological system. The input is a
889 relevant set of cells $S \subset \{1, \dots, n\}$, such as those representing a transition of interest. These could include
890 a branch in the cell-differentiation landscape or clusters interconnected by transitory cells. The output is
891 a ranking of genes ordered by their change scores representing their association with the low-density
892 regions in the selected set of cells. Top genes in this ranking can be interpreted as driving the transitions
893 in low-density regions.

894 We first compute a measure of local variability of a gene for each cell-state: We compute the expression
895 change from a cell when transitioning to each of its neighbors and normalize the change by the distance
896 between the cells in state space to account for the magnitude of the state transition. The maximal
897 normalized change amongst the neighbors of the cell is nominated as the local variability of the gene for
898 the corresponding state. Formally, the local variability for gene j in cell-state i is defined as:

$$899 \quad d_i^j := \max_{l \in N_i} \sqrt{\frac{(m_i^j - m_l^j)^2}{\|x_i - x_l\|_2}}$$

900 where, m_i^j denotes the MAGIC imputed expression of gene j in cell i , N_i is the set of k nearest neighbors
901 of cell i .

902 We next compute, a low-density change score s_j for each gene j , as the sum of the gene change rates
903 d_i^j across the selected cells, inversely weighted by the cell-state densities, $\rho(x_i)$:

$$904 \quad s_j := \sum_{i \in S} \frac{d_i^j}{\rho(x_i)}$$

905 This scoring approach encapsulates the hypothesis that genes with high change score in low cell-state
906 density regions may be driving transitions. Genes are ordered by the change score and genes with scores
907 $> 95^{\text{th}}$ percentile are considered to be driving low-density changes (**Supplementary Fig. 7**).

908

909

910 **Primed and lineage-specific accessibility scores from scATAC-seq data**

911 Gene scores from scATAC-seq are typically computed by summarizing the accessibility of peaks in the
912 body of the gene and its vicinity³⁶. This, however, does not consider the history and temporal dynamics
913 of peak accessibility. Enhancer priming, where open chromatin peaks are pre-established in stem cells
914 without turning on gene expression but maintain the gene locus in an open state for lineage-specific
915 upregulation, is an important mechanism through which stem cells encode high differentiation potential³²⁻
916 ³⁴. To investigate the establishment of peak accessibility, we devised a procedure to disentangle primed
917 and lineage-specific peaks in the context of cell-fate specification. As a result of the sparsity and noise of
918 scATAC-seq data, our approach utilizes several abstractions and consists of the following steps:

- 919 1. Identification of peaks with accessibility strongly correlated with gene expression at metacell
920 resolution
- 921 2. Determination of peaks with higher accessibility in the lineage under consideration compared to
922 other lineages using differential accessibility testing between metacells
- 923 3. Classification of peaks as primed or lineage-specific based on accessibility patterns in stem cells
- 924 4. Determination of primed and lineage-specific accessibility scores for each gene at single-cell
925 resolution.

926 We developed this approach to identify primed and lineage-specific peaks in the transition from
927 hematopoietic stem cells (HSCs) to B-cell fate committed cells (proB) (**Fig. 3**). We used the monocyte
928 and erythroid lineages as the alternative lineages to test for B-cell lineage specificity.

929 930 **Determination of primed and lineage-specific peaks**

931 **Metacells and gene-peak correlations using SEACells**

932 We used our SEACells algorithm²¹ to identify metacells from the T-cell depleted bone marrow. SEACells
933 aggregates highly related cells into metacells overcoming the sparsity in single-cell data while retaining
934 heterogeneity. We used the ATAC modality of the multiome data to identify metacells. We used metacells
935 to compare the expression of a gene with the accessibility of each peak in a window of 100kb around the
936 gene to identify the subset of peaks that significantly correlate with expression of the gene (correlation
937 ≥ 0.1 , p -value ≤ 0.1 , Empirical null) (**Supplementary Fig. 19A**).

938

939 **Peaks relevant to particular lineages**

940 Metacells and gene-peak correlations were computed using all hematopoietic lineages in our dataset.
941 We performed differential accessibility analysis to identify the subset of peaks with greater accessibility
942 in the lineage under consideration. We used edgeR⁶⁰ to perform differential accessibility with metacell
943 counts as input. The use of metacells rather than single-cell data for differential accessibility has been
944 demonstrated to provide better sensitivity and specificity³⁷. To identify peaks that are relevant to the B-
945 cell lineage, we compared accessibility in pro B-cell metacells and metacells of the erythroid (EryPre1)
946 or monocyte (Monocyte) lineages and retained peaks with the accessibility fold-change $\log_2FC > 0$ in
947 either comparison. While this ensures that the selected peaks have greater accessibility compared to
948 other lineages, it does not exclude ubiquitously accessible peaks. We therefore excluded peaks with
949 $\log_2FC < 0.25$ in the comparison between stem-cells (HSCs) and erythroid and monocyte lineages.

950 The resulting set of peaks demonstrate substantially greater accessibility in B-cell lineages compared to
951 all other cell-types (**Supplementary Fig. 19C-D**)

952

953 Classification of primed and lineage-specific peaks

954 After identifying peaks with greater accessibility in the B-cell lineage, we assigned primed or
955 lineage-specific status to each peak with a simple logic: A peak is annotated as primed if it is accessible
956 in HSCs and lineage-specific if it is not. Accessibility in HSCs was determined using Poisson statistics as
957 described in SEACells²¹. The mean of the Poisson distribution for a cell-type c is estimated using

$$958 \lambda = \frac{\text{Total fragments in } c}{\text{Effective genome length}}$$

959 Where *effective genome length* is set to *num of peaks* * 5000. For a peak p in cell-type c with n
960 fragments, λ is used to estimate the P value of observing more than n fragments, and p is considered
961 open in c if $P < 1e - 2$.

962 Primed and lineage-specific scores

963 We utilized the primed and lineage-specific peaks to derive primed and lineage-specific scores for the
964 associated genes at single-cell resolution. For each gene g and cell i , the primed accessibility score
965 s_{ig}^{primed} is computed as

$$966 s_{ig}^{primed} = \frac{\sum_{p \in g_{primed}} a_{ip} c_{gp}}{\sum_{p \in g_{primed}} c_{gp}}$$

967
968 Where g_{primed} is the set of primed peaks that significantly correlate with gene g , a_{ip} is the accessibility
969 of peak p in cell i , and c_{gp} is the correlation between peak p and expression of gene g computed using
970 metacells. Therefore, the primed score is a weighted average of the accessibility of primed peaks that
971 correlate with the gene. The lineage-specific score s_{ig}^{lin} is computed in an analogous manner where g_{lin}
972 is the set of lineage-specific peaks that significantly correlate with gene g :

$$973 s_{ig}^{lin} = \frac{\sum_{p \in g_{lin}} a_{ip} c_{gp}}{\sum_{p \in g_{lin}} c_{gp}}$$

974
975 Given the sparsity of the scATAC data, we used imputed peak accessibility for computing scores. The
976 peak counts dataset was TF-IDF normalized⁶¹ to preferentially weight peaks which are highly accessible
977 in a small proportion of cells. The MAGIC algorithm² was then used to perform imputation using
978 normalized accessibility as the input.

979 Data visualization

980 Accessibility trends along pseudo-time were computed using Mellon. Trends are visualized as a
981 percentage of the maximum value of each trend, to allow for better comparison across genes.

982 Application to T-cell depleted bone marrow data

983 We applied primed and lineage-specific accessibility scores to characterize commitment of hematopoietic
984 stem cells to B-cells using the T-cell depleted bone marrow multiome data. We used hematopoietic stem
985 cells (HSC), hematopoietic multipotent cells (HMP), common lymphoid progenitor (CLP) and pro B-cells
986 along the B-cell lineage to investigate the open chromatin landscape (**Fig. 3B**). The cells were chosen to

989 span the commitment of stem cells to the B-cell lineage. The high- and low-density regions were manually
990 assigned by comparing pseudotime and log-density of the selected subset of cells (**Fig. 3B**).

991 **Primed and lineage-specific accessibility scores in B-cell specification**

992 We applied the SEACells algorithm²¹ to identify metacells using the ATAC modality of the T-cell depleted
993 bone marrow data. Metacells were identified using all cells, resulting in 115 metacells according to
994 recommended heuristic for selecting the number of metacells. Peak accessibility and gene expression
995 correlations were determined using all metacells and the subset of genes with at least 5 peaks were
996 selected for downstream analysis (**Supplementary Fig. 19A**). We computed gene change scores using
997 Mellon using the subset of cells that define B-cell lineage commitment. Genes in the 95th percentile of
998 gene change scores with B-cell specific upregulation in the low-density regions were used to characterize
999 the role of enhancer priming (**Fig. 3**). Primed and lineage-specific accessibility scores were computed for
1000 the subset of these genes with at least one lineage-specific and primed peak each.

1001 **In silico ChIP**

1002 We used in silico ChIP-seq³⁷, a recently published approach to identify predicted targets of master
1003 regulators of B cell lineage commitment, specifically EBF1 and PAX5. Approaches like FIMO⁶² can
1004 determine enrichment scores for TF motifs in ATAC-seq peak sequences but the scores alone are not
1005 sufficiently reliable to predict TF targets. In silico ChIP-seq provides a framework for predicting TF targets
1006 by using single-cell multiome (scRNA-seq and scATAC-seq) data in addition to motif enrichment by
1007 correlating the expression of a TF to the accessibility of a peak. A combination of a high gene-peak
1008 correlation and high motif score is more indicative of potential TF binding compared to a peak with only
1009 a high motif score³⁷. We used our Python adaptation of in silico ChIP-seq using the SEACells metacells
1010 as input (github.com/settylab/atac-metacell-utilities). FIMO⁶² was used to associate TF motifs with ATAC-
1011 seq peaks, resulting in a peak by TF matrix of scores indicating the strength of match of the TF motif in
1012 the peak sequence. In silico TF binding scores are computed as product of correlation between TF
1013 expression and peak accessibility and FIMO motif scores as follows:

$$1014 \quad x_{ij} = \rho_{ij} * \text{minmax} \left(\frac{s_{ij}}{\max(s_j)} * \max(a_i) \right)$$

1015 Where i is the a ATAC-seq peak and j is the a TF of interest, ρ_{ij} is the Spearman rank correlation
1016 coefficient of accessibility of i and expression of j computed across all metacells, s_{ij} is the FIMO motif
1017 enrichment score for TF motif j binding in sequence of peak i , $\max(s_j)$ is the maximum FIMO score for
1018 TF j across all peaks, and a_i is the maximum accessibility of peak i across all cell type metacells.

1019 Minmax normalization is performed as follows:

$$1020 \quad \text{minmax}(x) = \frac{x - \min(x)}{\max(x) - \min(x)}$$

1021 The final in silico ChIP-seq output is a peak by TF matrix, containing a value between -1 and 1 indicating
1022 how likely a TF is to bind at a given peak and whether it has a repressive (negative) or activating (positive)
1023 effect, or 0 if a peak does not meet the minimum in silico ChIP-seq score (0.15).

1024

1025 **Regulation of EBF1**

1026 Peaks correlated with EBF1 expression were ordered using the procedure outlined in the section “Genes
1027 Driving Low-Density Cell-State Transitions” using imputed peak accessibility to compute accessibility
1028 change scores instead of gene change scores. In silico-ChIP was to identify the transcription factors with
1029 predicted binding sites in the top peak.

1030 Time-Continuous Density

1031 Time-series single-cell datasets provide snapshots of the changing cell-state densities at discrete time
1032 intervals. Our goal is to compute a time-continuous density function to interpolate cell-state densities at
1033 any time between the measured timepoints.

1034 We therefore incorporated a time coordinate into the Gaussian process used to generate the log density
1035 function and use the covariance of the Gaussian process to link temporally similar cell-states. Effectively
1036 the covariance function of time-continuous density has two components: (i) similarity between cells in the
1037 cell-state space and (ii) similarity between cells based on their measurement times. Similarity in cell-state
1038 space is encoded through the Matern52 kernel with the length-scale parameter as described in
1039 **Supplementary Note 5**. We now describe the Matern52 length-scale parameter for the temporal
1040 similarity component.

1041 The length scale should be designed such that the covariance between cells from different timepoints
1042 reflects the covariance of densities between those timepoints. Therefore, we optimized the length scale
1043 to reflect the empirically observed covariance of density functions between different time points.
1044 Specifically, we employ Mellon to compute first time-point specific density functions ρ_t using only the cells
1045 from the corresponding time point t . We next evaluated these functions on *all cells from all timepoints*,
1046 and computed a correlation of cell-state density between timepoints:

$$1047 \quad p_{t,t'} := \text{Corr}[\rho_t(x_i)_{i \in \{1, \dots, n\}}, \rho_{t'}(x_i)_{i \in \{1, \dots, n\}}]$$

1048 Where t and t' represent two time points, and $\text{Corr}[\cdot, \cdot]$ denotes the Pearson correlation. This is used to
1049 derive a correlation matrix between all measured timepoints T (**Supplementary Fig. 23A-D**):

$$1050 \quad P := (p_{t,t'})_{t,t' \in T}$$

1051 This matrix P is then compared to the covariance matrix of time points using the Matern52 kernel. Given
1052 the isotropy of the kernel function, it maps a scalar temporal difference $t - t'$ to a covariance value. The
1053 kernel-based covariance matrix is defined as:

$$1054 \quad K_L := \text{Matern52}(l_t)(t - t')_{t,t' \in T}$$

1055 Where l_t is the length scale parameter for the time coordinate. We thus select the l_t by optimizing:

$$1056 \quad l_t := \arg \min_{L' \in \mathbb{R}^+} \| P - K_{L'} \|_2.$$

1057 The optimized length scale is used for the Matern52 covariance kernel for the time coordinate, denoted
1058 as $\text{Matern52}(l_t)$ (**Supplementary Fig. 23E-F**).

1059 The resulting covariance kernel for cells i and j , situated at their respective states x_i, x_j , and
1060 measurement times t_i, t_j , is then given as:

$$1061 \quad k(i, j) = \text{Matern52}(l_t)(t_i - t_j) \cdot \text{Matern52}(l)(x_i - x_j)$$

1062 Where $\text{Matern52}(l)$ designates the Matern52 covariance kernel for cell-state coordinates and
1063 $\text{Matern52}(l_t)$ designates the Matern52 covariance kernel for time coordinates.

1064 This construction is easily implemented with Mellon, since it is designed to support any combination of
1065 covariance functions, each operating in distinct active dimensions – in this case, either time or cell-state
1066 coordinates.

1067 Using this covariance function, Mellon can compute a continuous density function over time and state
1068 space using all samples, and thus can be used to interpolate cell-state densities at unmeasured
1069 timepoints. This function is also differentiable in time and state space, and the change in density over
1070 time can be determined using the first derivative (**Supplementary Video 1, Figure 25**).

1071

1072 Leave-one-out Cross Validation

1073 We validated the effectiveness of the time-continuous density function using a leave-one-out cross-
1074 validation strategy (**Supplementary Fig. 24**). We computed a time-continuous density function after
1075 excluding cells from a particular timepoint and evaluated the densities at the excluded timepoint using
1076 this density function. We compared these densities with a time-agnostic density, which was computed
1077 exclusively using cells from the excluded time point and then evaluated across all states. Note that these
1078 two density functions were derived from mutually exclusive training datasets.

1079 Density along Trajectory

1080 Time-continuous density provides a platform to decipher the dynamics of cell-type proportions and fate
1081 choices in true temporal order. As proof of principle, we investigated the cell-type proportion dynamics
1082 along the trajectory of a particular lineage. We first used Palantir¹⁴ to derive fate propensities for all cells
1083 and selected the subset of cells with high propensity towards a particular fate.

1084 In the mouse gastrulation data, we applied Palantir using all cells across all timepoints and selected cells
1085 which specify the erythroid cells (**Supplementary Fig. 26**). Palantir was also used to derive a pseudo-
1086 temporal order of progression of cells in the erythroid trajectory (**Supplementary Fig. 26**). Note that the
1087 pseudo-time order does not take measurement time into consideration and represents the potential
1088 journey of a cell through the cell-state space as it acquires erythroid fate. Further, cells measured at any
1089 timepoint can span a range of pseudo-time depending on the developmental stage.

1090 The Palantir fate probability of cell state x_i reaching fate F is represented by the function $f(x_i, F)$.
1091 Accordingly, we define our threshold function for fate F at pseudotime t as:

$$1092 \quad T_F(t) = \max_{s < t} P_{99\%}[f(x_k, F)_{k \in I_s}].$$

1093 In this equation, $P_{99\%}$ is the 99% percentile function and I_s is the set of all cells whose Palantir pseudotime
1094 is less than or equal to s . We then identify the subset of cells that are part of the branch leading to fate F
1095 as:

$$1096 \quad N_F = \{i \in \{1, \dots, n\} | f(x_i, F) > T_F(t_i) - \epsilon\}.$$

1097 In the above equation, t_i is the pseudotime of the i^{th} cell, and ϵ is a small chosen value (in our case,
1098 0.01), which manages how much a cell can fall below the threshold while still being accepted as part of
1099 the branch. To simplify computation, we only calculate T_F for 500 specific pseudotime points along the
1100 trajectory, using the next larger pseudotime relative to t_i in this range to evaluate $T_F(t_i)$. This algorithm
1101 has been incorporated into the existing Palantir python package.

1102 We next determined the joint cell-state density between pseudo-time and real time leveraging the time-
1103 continuous density function. We first used Gaussian process as implemented in Mellon to map pseudo-
1104 time to each coordinate of the cell-state space. This effectively generates a trajectory traversing the cell-
1105 state space by mapping the 1-dimensional pseudotime to high-dimensional cell-state space. Formally,
1106 the trajectory for each dimension $m \in (1, \dots, d')$ is defined via the mean of the posterior distribution of T_F^m
1107 in the Bayesian model:

$$\begin{aligned} T_F^m &\sim GP\left(\overline{x_{j \in N_F}^m}, \text{Matern52}(1)\right) \\ x_i^m &\sim N\left(\hat{T}_F^m(s_i), 0.01\right), \quad x_i \in N_F \end{aligned} \quad (2)$$

1108

1109 Where $\overline{x_{j \in N_F}^m}$ represents the average of this coordinate across all cells in branch F . The trajectory can
1110 then be denoted by

1111

$$\begin{aligned} T_F &: [0,1] \rightarrow \mathbb{R}^{d'} \\ T_F &= (\hat{T}_F^m)_{m \in (1, \dots, d')} \end{aligned}$$

1112 where \hat{T}_F^m is the mean of the posterior of (2). The length scale of 1 and variance of 0.01 were selected
1113 by examination of a range of values for compatibility with cell states represented via Palantir diffusion
1114 maps.

1115 Finally, the time-continuous density function $\rho: \mathbb{R}^{d'} \times [0,1] \rightarrow \mathbb{R}^+$ can be evaluated along the trajectory to
1116 calculate joint cell-state density $\rho(T_F(s), t)$ for any given pseudotime s and actual time t (**Fig. 4F**).

1117

1118 Marginal Cell Type Proportions over Time

1119 We used the joint cell-state density $\rho(T_F(s), t)$ to determine the dynamics of cell-type proportions over
1120 real time. We first assign a cell type to each section of the pseudo-temporal trajectory T_F . This is achieved
1121 by computing a density function ρ_H for each annotated cell type H using Mellon. The cell type annotation
1122 $h(s)$, for a given pseudotime s is then given by the largest cell type density for this point on the trajectory
1123 as follows

1124

$$h(s) := \operatorname{argmax}_{H'} \rho_{H'} \circ T_F(s).$$

1125 The cell-type annotation pseudotime s can then be represented as an indicator function:

1126

$$\mathbb{I}_H(s) = \begin{cases} 1 & , h(s) = H \\ 0 & , \text{otherwise} \end{cases}$$

1127 We next marginalized the joint cell-state density $\rho(T_F(s), t)$ over pseudo-time to determine the total mass
1128 of a cell type. Specifically, the mass of cell type H along the trajectory of fate F at a real time point t is
1129 determined as

1130

$$m_F^H(t) = \int_0^1 \rho(T_F(s), t) \cdot \mathbb{I}_H(s) ds.$$

1131 Finally, the relative proportion of cell type H at a real time t is given by normalizing the masses across all
1132 cell-type as follows:

1133

$$a_F^H(t) = \frac{m_F^H(t)}{\sum_{H'} m_F^{H'}(t)}.$$

1134 This provides a quantifiable measure of cell type proportions over time, offering valuable insights into the
1135 temporal evolution of cell types in a given biological system.

1136

1137

1138 **Application to mouse gastrulation data**

1139 We applied Mellon to determine the time-continuous cell-state density for the mouse gastrulation data⁴⁴
1140 across all measured time points: E6.5, E6.75, E7.0, E7.25, E7.5, E7.75, E8.0, E8.25 and E8.5. Data was
1141 preprocessed as described in section “Mouse gastrulation data in Data preprocessing”. Diffusion maps
1142 were constructed using batch corrected PCs across all cells using the Palantir package¹⁴. We selected
1143 25 components, as they encapsulated all significant biological variations. Density results remained stable
1144 beyond this point with respect to the number of components (**Supplementary Fig. 13**). Time-continuous
1145 densities were computed following the procedure described above with default parameters.

1146 Palantir¹⁴ was to derive pseudo-temporal order and cell-fate propensities. Palantir was run with default
1147 parameters by using an Epiblast cell as the start and manually setting the following cell-types as
1148 terminals: Cardiomyocytes, Erythroid, Endothelial, Neural crest, Brain, Notochord, Allantois, ExE
1149 endoderm. Since our goal was to identify cells with high fate propensity to erythroid lineage, a finer
1150 resolution terminal state identification was not necessary. Erythroid lineage cells were identified using
1151 Equation (1). Joint cell-state density over pseudo-time and real-time were visualized using 200 points
1152 along pseudo-time and 500 points between every pair of measured timepoints.

1153

1154

1155 **Data preprocessing**

1156

1157 scRNA-seq data preprocessing and analysis

1158 The following procedure was used for preprocessing scRNA-seq data across datasets unless specified
1159 otherwise: Raw counts were normalized by dividing the counts by the total counts per cell. The normalized
1160 data was multiplied by the median of total counts across cells to avoid numerical issues and then log-
1161 transformed with a pseudocount of 0.1. Feature selection was then performed to select the top 2500 most
1162 highly variable genes, which was used as input for principal component analysis with 50 components.
1163 PCs were used as inputs for leiden clustering and UMAP visualizations. Preprocessing and analysis was
1164 performed using the scanpy⁶³ package.

1165 Diffusion maps were computed using the Palantir¹⁴ package with default parameters and PCs as the
1166 inputs. The diffusion kernel was also used for MAGIC² gene expression imputation.

1167 Batch correction where applicable was performed using Harmony with default parameters⁶⁴. Batch
1168 corrected PCs if applicable were used as inputs for UMAPs, diffusion maps, and imputation.

1169

1170 T-cell depleted bone marrow single-cell multiome data

1171 Raw gene counts, ATAC fragment files and cell metadata were downloaded from⁶⁵.

1172 **RNA modality**

1173 scRNA-seq data was processed using the procedure described in section “scRNA-seq data
1174 preprocessing and analysis”, which mimics the analysis in²¹.

1175 *Cell-type annotation*

1176 All hematopoietic stem and progenitor cells (HSPCs) were grouped as one cell-type in the T-cell depleted
1177 bone marrow. To achieve higher granularity among the stem and progenitor cells, we integrated this data
1178 with a dataset of CD34+ bone marrow cells using Harmony⁶⁴. This dataset is enriched for stem and
1179 progenitor cells and thus the associated cell-type information can be utilized to better resolve the cell-
1180 types within the HSPC cluster of the T-cell depleted bone marrow data. Batch corrected PCs were used
1181 for leiden clustering, and the HSPC cluster of the T-cell depleted data were assigned to different stem
1182 and progenitor cell-types based on their clustering with the CD34+ bone marrow data. Clusters
1183 associated with the B-cell trajectory were annotated using the markers described in⁶⁶.

1184 *Mellon cell-state density*

1185 Mellon was applied with default parameters using 20 diffusion components to compute cell-state density.
1186 Gene change scores, primed accessibility scores, and lineage-specific accessibility scores were
1187 computed as described above. IL7R signaling targets were downloaded from NicheNet⁶⁷ and signature
1188 scores were computed by averaging the z-scored imputed gene expression.

1189

1190 **ATAC modality**

1191 ArchR³⁶ pipeline was used for analysis of the ATAC modality. In ArchR, data was normalized using
1192 IterativeLSI and SVD to determine a lower-dimensional representation of the sparse data. The first SVD
1193 component showed greater than 0.97 correlation with log library size and was excluded from downstream
1194 analysis. SVD was used as input to cluster the data with leiden and visualization using UMAPs. SVD also

1195 served as input for computing diffusion and MAGIC imputation of peak accessibilities and gene scores.
1196 Peak calling was performed within ArchR using only the nucleosome free fragments as described in ²¹.

1197 A handful of cells which passed the RNA QC thresholds did not clear the thresholds in the ATAC modality.
1198 RNA preprocessing and analysis was repeated after excluding these cells. Mellon was applied with
1199 default parameters using 20 diffusion components to compute cell-state density of the ATAC modality.

1200 **Palantir trajectories**

1201 Palantir¹⁴ was used to infer pseudo-temporal trajectories of hematopoietic differentiation. Palantir was
1202 applied to the RNA modality using default parameters with the number of diffusion components (n=8)
1203 chosen by the Eigen gap statistic. A CD34+ hematopoietic stem cell was used as the start. Terminal cells
1204 were manually specified for erythroid, monocyte, B-cells, plasmacytoid dendritic cells. Note that the pre-
1205 pro B state of the B-cell trajectory is almost exclusively defined by cell-cycle⁶⁶ and hence Palantir was
1206 run with pre-pro B and naïve B as the terminals. The B-cell fate probability was then computed as the
1207 sum of pre-pro B and naïve B probabilities.

1208 Cells with increasing probability towards each lineage were selected as the lineage cells highlighted in
1209 **Fig. 1D**. B-cell lineage cells were comprised of Hematopoietic stem cells (HSCs), Hemopoietic
1210 multipotent progenitors (HMPs), Common Lymphoid progenitors (CLPs), prepro B-cells, pre B-cells, pro
1211 B-cells and Naïve B-cells. pDC lineage cells were comprised of HSCs, HMPs, Myeloid precursors, and
1212 pDCs. Erythroid lineage cells were comprised of HSCs, Megakaryocyte erythroid precursors (MEPs) and
1213 erythroid precursors. Monocyte lineage cells were comprised of HSCs, HMPs, Myeloid precursors,
1214 monocyte precursors and monocytes.

1215 Cells involved in lineage specification (highlighted cells in **Fig. 1D**) where chosen as the subset of the
1216 lineage cells spanning from HSCs to the cell-type where the fate propensity reached 1. B-cells: HSCs,
1217 HMPs, CLPs, prepro B-cells, pro B-cells. pDCs: HSCs, HMPs, MyeloidPre, pDCs. Erythroid lineage:
1218 HSCs, MEPs. Monocytes: HSCs, HMPs, Myeloid precursors, monocyte precursors and monocytes.

1219

1220 **HCA bone marrow**

1221 The processed annData was downloaded from²⁷. The downloaded data was pre-batch corrected across
1222 all donors. Cell types that do not differentiate in the bone marrow such as T-cells, NK cells and plasma
1223 cells were excluded from the analysis. Following the cell filtering, each donor was analyzed separately
1224 using the steps outlined in the section “scRNA-seq data preprocessing and analysis”.

1225 Palantir¹⁴ was applied separately for each donor using the same procedure that was described for the T-
1226 cell depleted bone marrow dataset. Mellon was applied with default parameters using 20 diffusion
1227 components to compute cell-state density.

1228

1229 **Pancreatic development**

1230 Processed anndata was downloaded from ¹⁷ and the data was generated by ²⁹. The pre-computed
1231 UMAPs, cell-type annotations and diffusion maps were used for analysis. Mellon was applied with default
1232 parameters to compute cell-state density.

1233

1234 **In-vitro endoderm differentiation**

1235 Raw counts and cell metadata was downloaded from³⁰. Wild-type cells were used for all analysis. Data
1236 analysis was performed using the steps outlined in the section “scRNA-seq data preprocessing and

1237 analysis”, batch correction was used to correct technical differences between two batches. Mellon was
1238 applied with default parameters to compute cell-state density.

1239

1240 Spatial organization of intestinal tissue

1241 Raw counts and zone information were downloaded from³¹ and processed using the steps outlined in the
1242 section “scRNA-seq data preprocessing and analysis”. Mellon was applied with default parameters to
1243 compute cell-state density.

1244

1245 Lung regeneration

1246 Processed anndata was downloaded from²⁸. The pre-computed UMAPs, cell-type annotations and
1247 diffusion maps were used for analysis. Mellon density functions were computed for each timepoint
1248 separately and evaluated across all cells.

1249

1250 scRNA-seq of murine models of lung adenocarcinoma

1251 Processed anndata object containing counts, visualization and cell-metadata were downloaded from⁹.
1252 scVI⁶⁸ was used in the publication for data integration and to derive a latent representation. scVI latent
1253 space was used as input for computing force directed layouts and diffusion maps instead of PCs like
1254 other datasets.

1255

1256 Mouse gastrulation atlas

1257 Processed data including batch corrected principal components and cell metadata were downloaded
1258 from⁴⁴. Batch corrected PCs were used as input for computing diffusion maps. Cells from the
1259 “mixed_gastrulation” samples were excluded since the timepoints are not well-defined. Further, ExE
1260 ectoderm, ExE endoderm and Parietal endoderm cells were excluded since their parental cells are not
1261 measured in the dataset. Given the complexity of the data, 25 diffusion components for computing time-
1262 continuous cell-state densities using Mellon.

1263

1264 iPS reprogramming dataset

1265 Raw counts and cell metadata were downloaded from⁸. The dataset contains reprogramming in two
1266 culture conditions: Serum and 2i. Cells cultured in 2i media were used for the analysis. Highly variable
1267 genes computed in the publication were used for the analysis using the steps outlined in the section
1268 “scRNA-seq data preprocessing and analysis”. iPS data was used for robustness analysis and
1269 benchmarking performance.

1270

1271 scATAC-seq of murine models of lung adenocarcinoma

1272 Raw peak counts and cell metadata were downloaded from⁴⁹. Immune and stromal cells were excluded
1273 from the analysis. Following cell filtering, peak counts were normalized using TFIDF following the
1274 procedure in²¹. SVD was to determine a lower-dimensional representation using normalized data as
1275 input. The first SVD component showed greater than 0.97 correlation with log library size and was
1276 excluded from downstream analysis. SVD was used as input for visualization using force directed layouts
1277 and diffusion maps. Mellon was applied with default parameters to compute cell-state density.

1278

1279 sortChIC data profiling histone modifications in murine hematopoiesis
1280 Raw peak counts and cell metadata were downloaded from⁴⁸ for all available histone modifications:
1281 H3K4me1, H3K4me3, H3K27me3, H3K9me3. Each modification was analyzed separately following the
1282 procedure described in the section “scATAC-seq of murine models of lung adenocarcinoma”: Data was
1283 normalized using TF-IDF, and then SVD was used to derive a low-dimensional representation. The first
1284 component of SVD was excluded due to high correlation with log library size and was excluded from
1285 downstream analysis. Mellon was applied with default parameters to compute cell-state density.

1286

1287 Skin differentiation Share-seq data

1288 The processed annData was downloaded from⁵¹ using the data generated by³². The pre-computed
1289 UMAPs, cell-type annotations and diffusion maps were used for analysis. Note that the diffusion
1290 components were derived using the MIRA multimodal representation which uses both RNA and ATAC
1291 modalities. Mellon was applied with default parameters to compute cell-state density.

1292

1293

1294

1295

1296

1297 **Robustness analysis**

1298

1299 The robustness of Mellon was evaluated by recalculating density estimations across a broad spectrum
1300 of parameter settings on multiple datasets. We carried out full density inference for an extensive range
1301 of length scales, numbers of landmarks, and numbers of diffusion components in the following scRNA-
1302 seq datasets : T-cell depleted bone marrow of human hematopoiesis (BM)²¹; CD34+ human bone marrow
1303 cells, a dataset of hematopoietic stem and precursor cells (CD34)²¹; COVID-19 atlas of peripheral blood
1304 mono nuclear cells (PBMCs) from healthy donors and critical patients (Covid)⁶⁹; iPS reprogramming
1305 dataset (ips)⁸ and the mouse gastrulation atlas (mgast)⁴⁴. These datasets cover a broad spectrum of
1306 systems with different complexities, cell numbers and contain discrete and continuous cell-states and
1307 cell-types. We compared the densities using Spearman correlation between results obtained from
1308 different parameter settings. As shown in **Supplementary Fig. 13-16**, Mellon results exhibited a high
1309 level of consistency in the results even when the parameters are varied orders of magnitude beyond the
1310 defaults.

1311 We further evaluated Mellon's robustness to down sampling the cells in the dataset. Starting with the full
1312 dataset, we serially removed 10% of cells until at least 100 cells were retained. We next computed
1313 densities for independently for each subsample by recomputing the principal components and diffusion
1314 components using only the cells in the subset. We then compared the density between all pairs of
1315 subsamples using the intersection of cells between the two samples (**Supplementary Figure 11**). The
1316 consistency is retained even when cells in the bottom 10th percentile of the average density between the
1317 pair of runs are used for comparison (**Supplementary Fig. 12**).

1318 This robustness evaluation provides empirical evidence of Mellon's ability to perform consistently under
1319 a wide range of parameters and under the condition of subsampling, which underscores its utility for
1320 accurate density estimation from high-dimensional single-cell data.

1321

1322 **Simulated datasets with ground-truth densities**

1323

1324 In order to validate the accuracy and precision of Mellon, we generated three datasets mirroring single-
1325 cell datasets of either continuum of cell-states or discrete clusters. Each dataset is accompanied by a
1326 predefined 'ground truth' density serving as a performance benchmark for Mellon.

1327 The datasets with continuum of cell-states were generated using a large Gaussian Mixture Model (GMM)
1328 designed to emulate a cellular differentiation tree. This tree was conceptualized as a series of velocity
1329 vectors, each connecting branching points and each being a slightly perturbed version of the vector of its
1330 parent node. For each node in this tree, a unique Gaussian was defined. The Gaussian's covariance
1331 matrix and mean were designed to create a distribution aligning with the velocity vector. Considering the
1332 inherent low dimensionality typically exhibited by a cell-state manifold, we adjusted the principal
1333 components of these Gaussians using an exponential decay scalar. The two continuous styles mimic the
1334 structure of CD34+ bone marrow RNA-seq and T-Cell depleted bone marrow RNA-seq datasets.

1335 The synthetic datasets representing single-cell datasets of discrete clusters was also generated using a
1336 GMM but with a different configuration. In this setup, we randomly sampled mean and covariance
1337 matrices to create an arbitrary GMM, resulting in mostly isolated clusters of simulated cells. This approach
1338 provided an alternative, contrasting framework for testing robustness of Mellon.

1339 The GMM allowed us to easily sample simulated cell states from both dataset types and to define
1340 corresponding ground truth probability density functions. We then utilized Mellon to compute the log-
1341 density of these simulated datasets. Ground-truth densities were compared with Mellon densities using
1342 Spearman correlations. As shown in **Supplementary Fig. 4**, this comparison effectively quantified
1343 Mellon's ability to infer cell densities from high-dimensional single-cell data, with Mellon exhibiting high
1344 consistency with the ground truth for both synthetic datasets.

1345 See **Supplementary Note 5** for further details and parameter choices for dataset simulations.

1346

1347

1348

1349 Comparison to density estimation approaches

1350

1351 A commonly used approach for density estimation with single-cell data is to calculate the reciprocal of
1352 the distance to the k th nearest neighbor, treating this value as a proxy for density². While straightforward,
1353 this method tends to produce a noisy density estimation and frequently fails to capture meaningful global
1354 trends (**Supplementary Fig. 2**).

1355 Another prevalent approach involves application of kernel density estimation (KDE) to the low-
1356 dimensional embeddings generated by tools like UMAPs or Force-Directed Layouts. While these
1357 visualization tools are powerful, their main design is not for density inference. They can produce unstable
1358 embeddings, and when KDE is applied, the instability in the embeddings directly translates into the
1359 density inference, resulting in less reliable outputs. Furthermore, the high compression involved in
1360 generating these low-dimensional representations means that they cannot capture all the relevant
1361 biological variability inherent in the data. Consequently, these methods often fail to depict all the nuanced
1362 details of the underlying cell-state density function (**Supplementary Fig. 2**).

1363

1364

1365 Efficient Pseudotime Trend Computation with Mellon

1366

1367 The versatility of Mellon extends beyond density inference, showcasing its robust capability in the swift
1368 computation of gene trends, defined as continuous, smooth functions that trace the trajectory of gene
1369 expression over pseudotime. Our Gaussian process (GP) regression-centric design not only serves as
1370 the backbone for Mellon's primary application but also efficiently caters to general GP regression, due to
1371 scalable features such as the fixed length scale for the Matern52 covariance kernel and landmarks for
1372 Sparse Gaussian process regression.

1373

1374 Gaussian processes shine in their adeptness at handling high-noise scenarios, for instance, non-imputed
1375 gene expression values. This strength enables Mellon to generate smooth gene trends from a selected
1376 cellular branch's temporal ordering using unimputed gene expression values, effectively capturing the
1377 dynamics of gene expression as cellular differentiation unfolds (**Supplementary Fig. 21**).

1378

1379 Mellon's implementation harnesses the power of the JAX library's vectorization capabilities and low-
1380 dimensional latent representations of functions within the GP framework, enabling efficient gene trend
1381 computations across a substantial quantity of genes. In tests using a 36-core CPU, Mellon was able to
1382 generate gene trends for up to 10,000 genes and 1,500 cells at 500 pseudotime points in about one
1383 second. This efficient computation allows high-throughput exploration of gene expression dynamics
1384 during cellular differentiation from large-scale single-cell datasets.

1385

1386

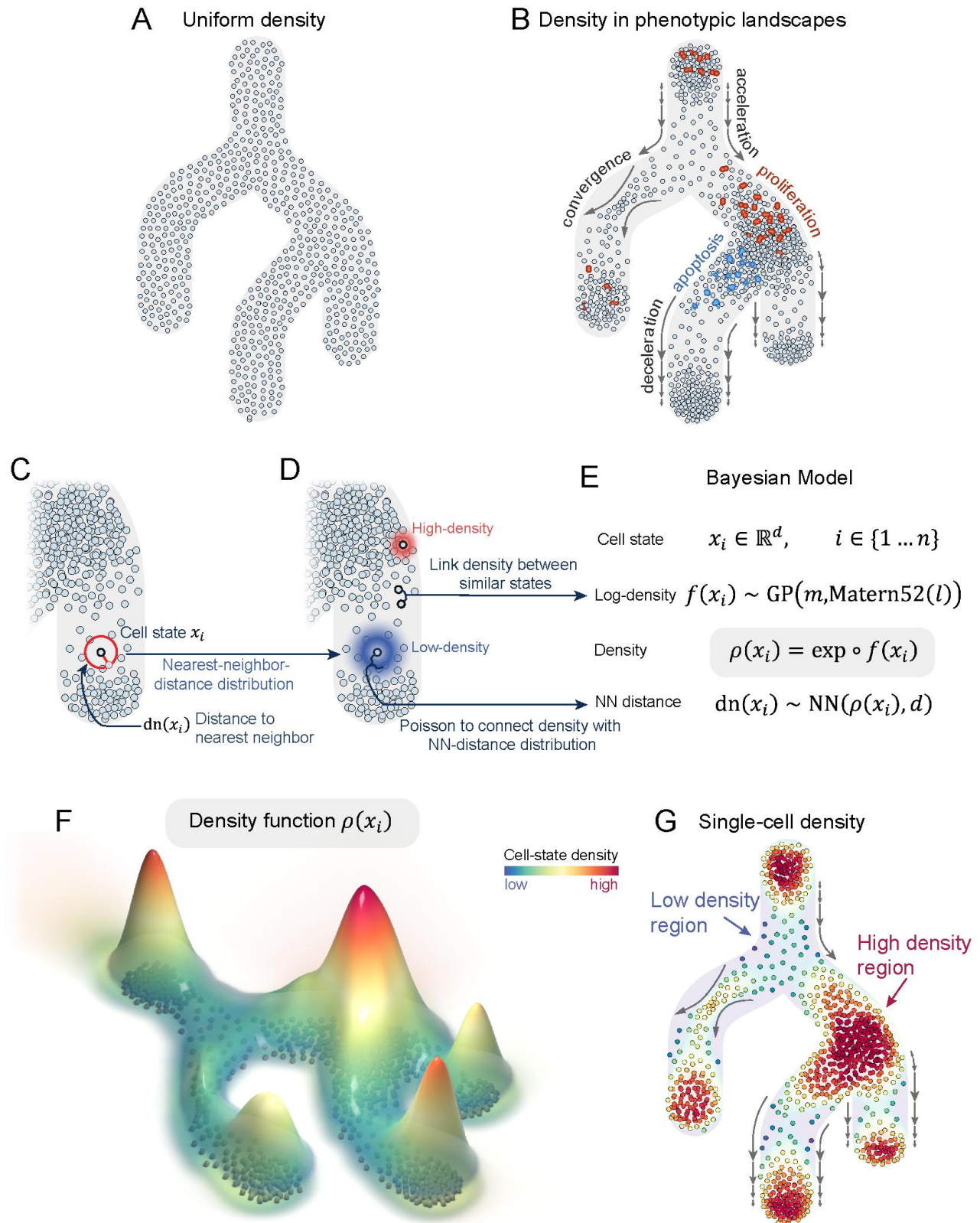
1387

1388

1389

1390

1391 **Figures**



1392

1393

1394 **Figure 1: Illustrative diagram detailing the principles and processes of Mellon.**

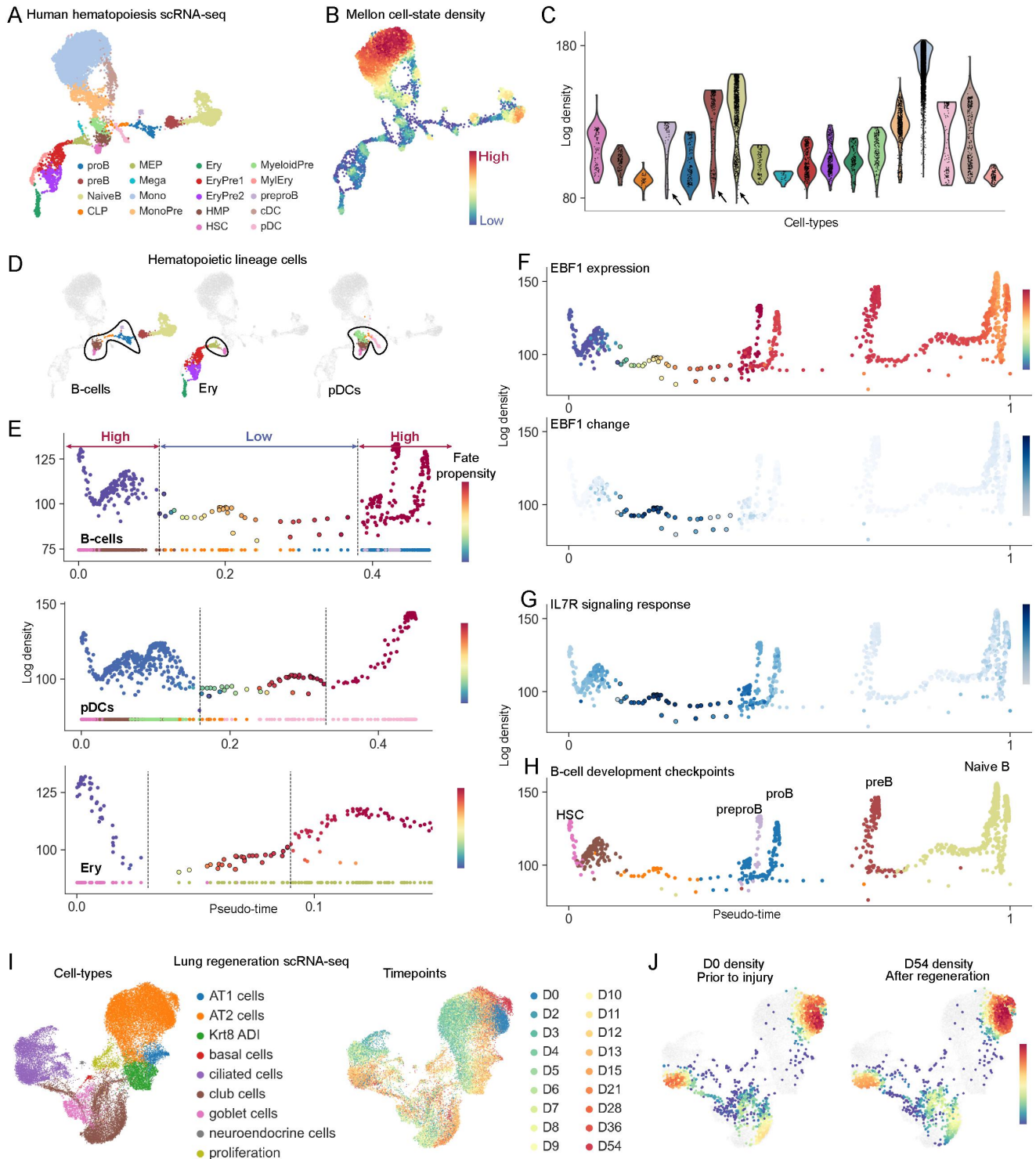
1395 A-B. An abstract depiction of a cellular differentiation landscape with cells uniformly distributed along its
1396 branches, representing a scenario not commonly found in biological systems. Diverse biological
1397 phenomena, as depicted in (B), impact cell-state density: apoptosis, acceleration, and divergence of cell-
1398 state changes lead to a decrease in density, while proliferation, deceleration, and convergence of cell-
1399 state changes increase density. Therefore, heterogeneity in cell-state densities is a norm rather than an
1400 exception in differentiation landscapes.

1401 C. Subset of cells with heterogeneous density are highlighted to illustrate the influence of biological
1402 factors in (B). Color gradient signifies the nearest-neighbor distribution around two example cells - one
1403 in a high-density state with a tighter distribution (red gradient) and another in a lower-density state with a
1404 broader distribution (blue gradient).

1405 D. Bayesian model employed by Mellon for density inference, underpinning the connection between the
1406 density estimation between neighboring cells using a Gaussian process and the log-density function as
1407 its random variable. Arrows relate the examples in panel C with their corresponding equations in D.

1408 E-F. Depict the resulting continuous density function from Mellon's inference process over the set of cells
1409 in B. E: Density function is visualized as a 3D landscape, where the z-axis represents density, and
1410 individual cell states are illustrated as spheres at the base. F color-codes the cell states from B according
1411 to their inferred densities, overlaying these with a translucent representation of the continuous density in
1412 the background. Examples of high- and low-density regions are highlighted.

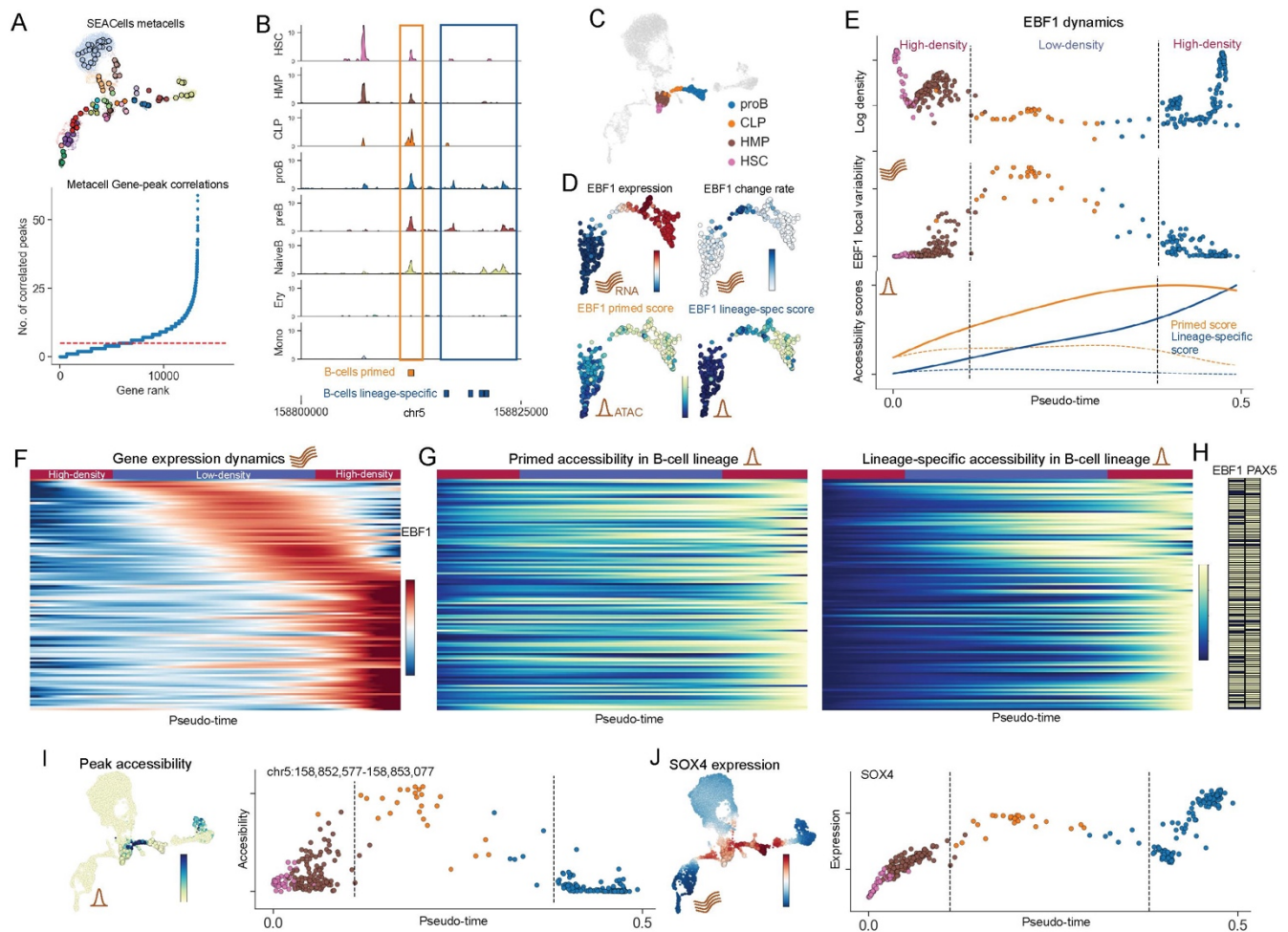
1413



1414

1415

1416 **Figure 2: Mellon reveals the density landscape of human hematopoietic differentiation**
1417 A. UMAP representation of the scRNA-seq dataset of T-cell depleted bone marrow²¹ colored by cell-
1418 types.
1419 B. Same UMAP as (A), colored by Mellon cell-state density
1420 C. Violin plots to compare cell-state densities among different hematopoietic cell-types. Arrowheads
1421 indicate example cell-types with high variability in density.
1422 D. UMAPs as in (A), highlighted by cells of the different lineages, left to right: B-cells, Erythroid lineage
1423 cells and plasmacytoid dendritic cells (pDCs). Lineage cells were selected based on cell-fate
1424 propensities. Cells spanning hematopoietic stem-cells to fate committed cells along each lineage.
1425 E. Plots comparing pseudotime ordering and log density during the fate specification of each lineage.
1426 Top to bottom: B-cells, pDCs and Erythroid lineages. Cells are colored by Palantir fate propensities, which
1427 represent the probability of each cell differentiating to the corresponding lineage. Points at the bottom of
1428 each plot are colored by cell- type. Subset of cells along each lineage spanning hematopoietic stem cells
1429 to fate committed cells are shown. Dotted lines indicate the low-density region within which fate
1430 specification takes place and were added manually.
1431 F. Plots comparing pseudotime and log density for all cells of the B-cell trajectory colored by EBF1
1432 MAGIC² imputed expression (top) and EBF1 local variability in gene expression (bottom).
1433 G. Same as (F), with cells colored by signature scores for IL7R response genes.
1434 H. Same as (F), with cells colored by cell-types. Density peaks correspond to well-characterized
1435 checkpoints during B-cell differentiation.
1436 I. UMAP representation of the scRNA-seq dataset of lung regeneration²⁸. Cells are colored by cell-type
1437 (left) and by timepoint of measurement (right). D0 is prior to injury and all subsequent timepoints show
1438 recovery from injury.
1439 J. UMAPs colored by density at D0 (left) and density at D54 (right). Cells from D0 and D54 are colored
1440 by density with cells from other timepoints in grey.
1441



1442

1443 **Figure 3: Dynamics of chromatin accessibility and gene expression during B-cell fate specification.**

1444
1445 A. Top: UMAP colored by cell-types and highlighted by SEACells²¹ metacells. Bottom: Plots showing the
1446 number of peaks significantly correlated with each gene. The correlations were computed using
1447 SEACells²¹ metacells.

1448 B. Coverage plots highlighting examples of B-cell primed (in orange) and B-cell lineage specific peaks
1449 (in blue). The genomic region is part of the EBF1 gene locus

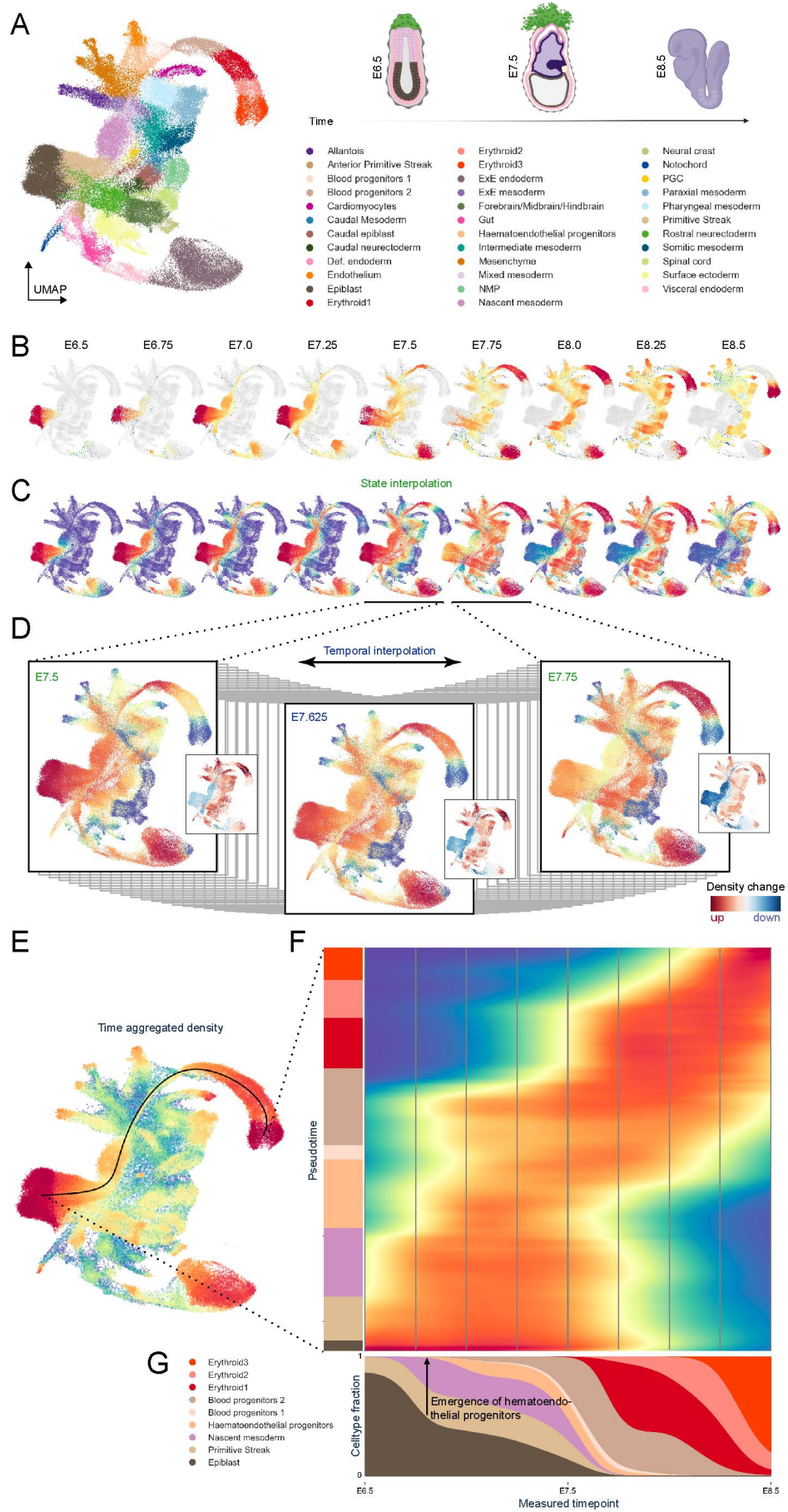
1450 C. UMAP colored by cell types included in B-cell specification. The full dataset is shown in grey.

1451 D. UMAP colored by EBF1 MAGIC imputed expression, EBF1 local variability, EBF1 primed accessibility
1452 scores and EBF1 lineage-specific accessibility scores. The subset of cells involved in B-cell specification
1453 (C) are shown.

1454 E. Top: Plots comparing pseudotime and Mellon density for the B-lineage cells, colored by cell-type.
1455 Middle: Plots comparing pseudotime and EBF1 local change for the B-lineage cells, colored by cell-type.
1456 Bottom: Solid lines show the trend of primed and lineage-specific accessibility scores for EBF1 in B-cell
1457 lineage. Dotted lines show the corresponding trends in the erythroid lineage. Vertical dotted lines show
1458 high- and low-density regions selected manually.

1459 F. Heatmaps with z-score expression of genes with high change scores and upregulation during B-cell
1460 specification. Genes are sorted based on their expression along pseudotime. Genes with at least 1
1461 primed and at least 1 lineage-specific peak from **Supplementary Fig. 20A** were used.

- 1462 G. Heatmaps of primed (left) and lineage-specific (right) accessibility scores for genes in (F) in the same
1463 order. Scores were scaled to maximum of 1 along the trend.
- 1464 H. Matrix indicating whether the genes in (F) are predicted targets of EBF1 or PAX5 using Insilico-ChIP³⁷.
- 1465 I. Left: UMAP colored by MAGIC imputed accessibility of the single ATAC peak (chr5:158,852,577-
1466 158,853,077) with highest change score in EBF1 correlated peaks. Right: Plot comparing pseudotime to
1467 peak accessibility for cells during B-cell specification in (C).
- 1468 J. Left: UMAP colored by SOX4 MAGIC imputed expression. Right: Plot comparing pseudotime to gene
1469 expression for cells during B-cell specification in (C).
- 1470



1472 **Figure 4: Depiction of time-continuous cell-state density estimation during mouse gastrulation**
1473 **using Mellon.**

1474 A. UMAP representation of the mouse gastrulation dataset⁴⁴. Illustrations on the right show a
1475 diagrammatic overview of the mouse embryo during gastrulation from E6.5 to E8.5, providing context to
1476 the developmental progression. *Created using BioRender.*

1477 B. UMAPs colored by Mellon cell-state density at each measured timepoint, demonstrating variability in
1478 cell-state densities within each observed timepoint.

1479 C. UMAPs colored by state-interpolated densities, derived from densities from (B), but evaluated across
1480 all cells. This showcases the potential of Mellon for extrapolating cell-state densities beyond directly
1481 sampled cell states.

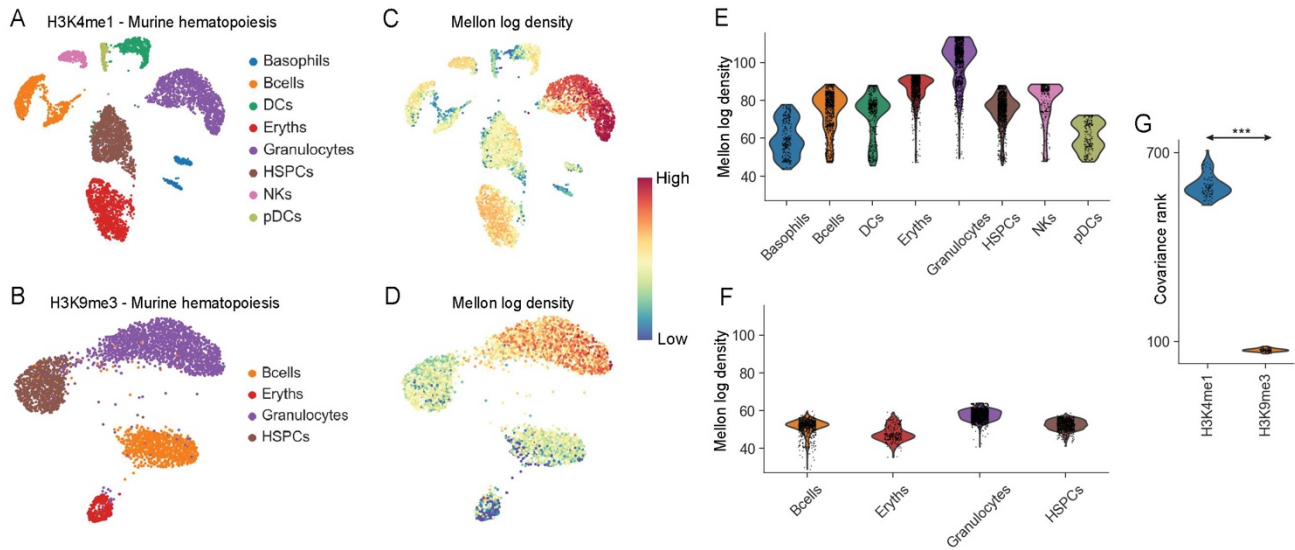
1482 D. Illustration of time-continuous density on UMAP for measured (E7.5, E7.75) and interpolated (E7.25)
1483 timepoints, further demonstrating the application of Mellon in interpolating cell-state densities beyond
1484 measured timepoints. Smaller accompanying UMAPs denote the temporal rate of change in cell-state
1485 density, with red signifying increasing density (enrichment) and blue indicating decreasing density
1486 (depletion).

1487 E. UMAP colored by cell-state density inferred using all-cells without using temporal information. Trend
1488 highlights the erythroid trajectory.

1489 F. Heatmap displays the time-dependent cell-state densities along the trajectory (pseudotime on the y-
1490 axis and real-time on the x-axis), with vertical grey lines signifying the measured timepoints.

1491 G. Marginal plot illustrating the proportional composition of cell-types along the erythroid trajectory at
1492 each timepoint, derived by integrating density in F across the trajectory segment associated with each
1493 specific cell type.

1494



1495

1496

1497 **Figure 5: Application of Mellon density estimation to single-cell chromatin data modalities.**

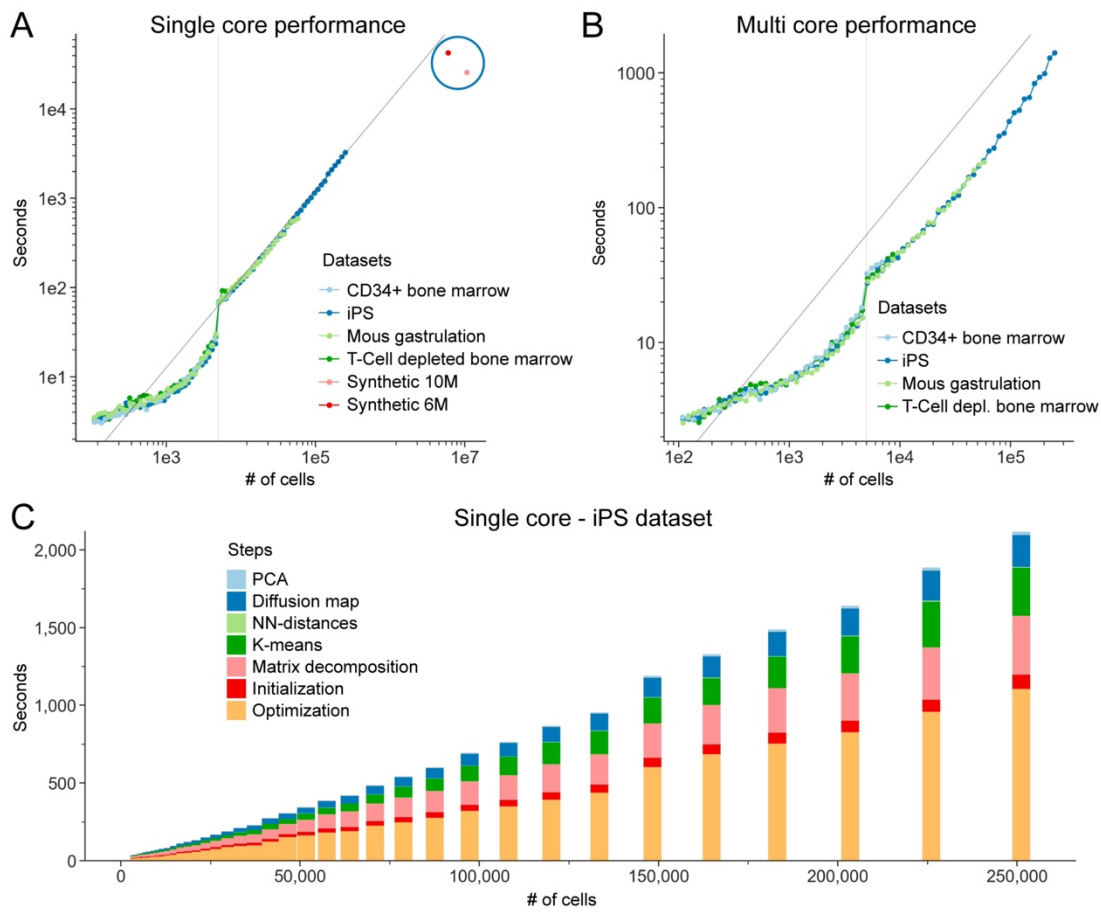
1498 A-B. UMAPs of H3K4me1 (A) and H3K9me3 (B) mouse bone marrow sort-ChIC dataset⁴⁸ colored by
1499 cell-type.

1500 C-D. Same as (A-B), with UMAPs colored by Mellon log density

1501 E-F. Violin plots to compare cell-state densities among different hematopoietic cell-types. Top: H3K4me1,
1502 Bottom: H3K9me3

1503 G. Violin plot of covariance matrix rank for each sort-ChIC dataset for 100 runs of Mellon by repeatedly
1504 subsampling 80% of the dataset. (***) p-value < 1e-30, Wilcoxon rank-sum test)

1505



1506

1507 **Figure 6: Performance benchmarking of Mellon for demonstrating its scalability and linear time**
 1508 **complexity.**

1509 A. Demonstrates the CPU time required for Mellon's density inference on a single core across various
 1510 dataset sizes from four distinct datasets. Each dataset is successively downsized by randomly removing
 1511 10% of cells. The data points in this log-log plot align closely with the diagonal line that has a slope of 1,
 1512 indicating a linear relationship between the number of cells and the CPU time required, which suggests
 1513 a linear time complexity of Mellon's algorithm, particularly for large datasets. Notably, statistics for the
 1514 two large synthetic datasets (6 million and 10 million cells), marked by a blue circle, fall below the
 1515 diagonal. This emphasizes that a nonlinear increase in compute time does not dominate, even for these
 1516 larger datasets. For these two synthetic datasets, the computation of diffusion components was omitted,
 1517 and the larger dataset (10 million cells) uses only 1,000 landmarks, instead of the usual 5,000. The
 1518 vertical line at 5,000 cells marks the point where the Gaussian process changes from a full process to a
 1519 sparse one, demonstrating how Mellon adapts to larger datasets by computing the density based on a
 1520 subset of 'landmark' cell states.

1521 B. Same as (A) but using 36 CPU cores, showcasing the computational efficiency achieved through
 1522 parallel processing. The data points, situated below the slope-1 diagonal, represent a decrease in CPU
 1523 time due to the parallelization of tasks.

1524 C. Breakdown of the total single-core CPU time for the iPS dataset into individual computational stages,
 1525 offering insights into the contribution of each stage to the overall density inference process.

1526

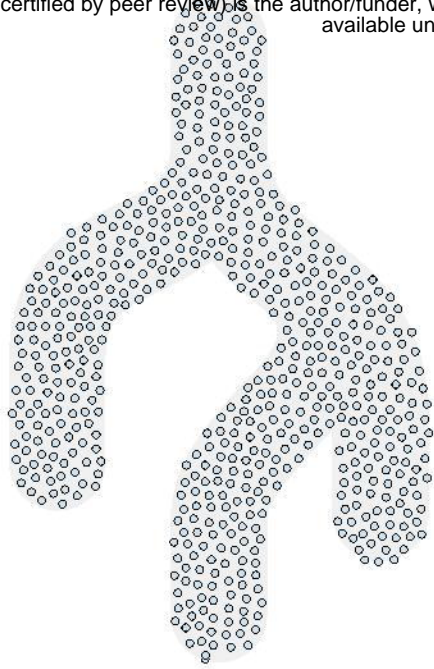
1527 References

- 1528
- 1529 1 Bendall, S. C. *et al.* Single-cell trajectory detection uncovers progression and regulatory
1530 coordination in human B cell development. *Cell* **157**, 714-725 (2014).
1531 <https://doi.org/10.1016/j.cell.2014.04.005>
- 1532 2 van Dijk, D. *et al.* Recovering Gene Interactions from Single-Cell Data Using Data Diffusion. *Cell*
1533 **174**, 716-729 e727 (2018). <https://doi.org/10.1016/j.cell.2018.05.061>
- 1534 3 Burkhardt, D. B. *et al.* Quantifying the effect of experimental perturbations at single-cell resolution.
1535 *Nat Biotechnol* **39**, 619-629 (2021). <https://doi.org/10.1038/s41587-020-00803-5>
- 1536 4 Antolovic, V., Lenn, T., Miermont, A. & Chubb, J. R. Transition state dynamics during a stochastic
1537 fate choice. *Development* **146** (2019). <https://doi.org/10.1242/dev.173740>
- 1538 5 Westbrook, E. R., Lenn, T., Chubb, J. R. & Antolović, V. Collective signalling drives rapid jumping
1539 between cell states. *bioRxiv*, 2023.2005.2003.539233 (2023).
1540 <https://doi.org/10.1101/2023.05.03.539233>
- 1541 6 Rukhlenko, O. S. *et al.* Control of cell state transitions. *Nature* **609**, 975-985 (2022).
1542 <https://doi.org/10.1038/s41586-022-05194-y>
- 1543 7 Nelms, B. & Walbot, V. Defining the developmental program leading to meiosis in maize. *Science*
1544 **364**, 52-56 (2019). <https://doi.org/10.1126/science.aav6428>
- 1545 8 Schiebinger, G. *et al.* Optimal-Transport Analysis of Single-Cell Gene Expression Identifies
1546 Developmental Trajectories in Reprogramming. *Cell* **176**, 928-943 e922 (2019).
1547 <https://doi.org/10.1016/j.cell.2019.01.006>
- 1548 9 Yang, D. *et al.* Lineage tracing reveals the phylogenetics, plasticity, and paths of tumor evolution.
1549 *Cell* **185**, 1905-1923 e1925 (2022). <https://doi.org/10.1016/j.cell.2022.04.015>
- 1550 10 Burdziak, C. *et al.* Epigenetic plasticity cooperates with cell-cell interactions to direct pancreatic
1551 tumorigenesis. *Science* **380**, eadd5327 (2023). <https://doi.org/10.1126/science.add5327>
- 1552 11 Coifman, R. R. *et al.* Geometric diffusions as a tool for harmonic analysis and structure definition
1553 of data: diffusion maps. *Proc Natl Acad Sci U S A* **102**, 7426-7431 (2005).
1554 <https://doi.org/10.1073/pnas.0500334102>
- 1555 12 Haghverdi, L., Buettner, F. & Theis, F. J. Diffusion maps for high-dimensional single-cell analysis
1556 of differentiation data. *Bioinformatics* **31**, 2989-2998 (2015).
1557 <https://doi.org/10.1093/bioinformatics/btv325>
- 1558 13 Setty, M. *et al.* Wishbone identifies bifurcating developmental trajectories from single-cell data.
1559 *Nat Biotechnol* **34**, 637-645 (2016). <https://doi.org/10.1038/nbt.3569>
- 1560 14 Setty, M. *et al.* Characterization of cell fate probabilities in single-cell data with Palantir. *Nat*
1561 *Biotechnol* **37**, 451-460 (2019). <https://doi.org/10.1038/s41587-019-0068-4>
- 1562 15 Trapnell, C. *et al.* The dynamics and regulators of cell fate decisions are revealed by
1563 pseudotemporal ordering of single cells. *Nat Biotechnol* **32**, 381-386 (2014).
1564 <https://doi.org/10.1038/nbt.2859>
- 1565 16 Buenrostro, J. D. *et al.* Integrated Single-Cell Analysis Maps the Continuous Regulatory
1566 Landscape of Human Hematopoietic Differentiation. *Cell* **173**, 1535-1548 e1516 (2018).
1567 <https://doi.org/10.1016/j.cell.2018.03.074>
- 1568 17 Lange, M. *et al.* CellRank for directed single-cell fate mapping. *Nature Methods* **19**, 159-170
1569 (2022). <https://doi.org/10.1038/s41592-021-01346-6>
- 1570 18 Rasmussen, C. E. & Williams, C. K. I. *Gaussian Processes for Machine Learning*. (2006).
- 1571 19 Snoek, J., Larochelle, H. & Adams, R. P. in *Advances in Neural Information Processing Systems*.
- 1572 20 Orkin, S. H. & Zon, L. I. Hematopoiesis: an evolving paradigm for stem cell biology. *Cell* **132**, 631-
1573 644 (2008). <https://doi.org/10.1016/j.cell.2008.01.025>
- 1574 21 Persad, S. *et al.* SEACells infers transcriptional and epigenomic cellular states from single-cell
1575 genomics data. *Nat Biotechnol* (2023). <https://doi.org/10.1038/s41587-023-01716-9>
- 1576 22 Oetjen, K. A. *et al.* Human bone marrow assessment by single-cell RNA sequencing, mass
1577 cytometry, and flow cytometry. *JCI Insight* **3** (2018). <https://doi.org/10.1172/jci.insight.124928>
- 1578 23 Pietras, E. M., Warr, M. R. & Passegue, E. Cell cycle regulation in hematopoietic stem cells. *J*
1579 *Cell Biol* **195**, 709-720 (2011). <https://doi.org/10.1083/jcb.201102131>

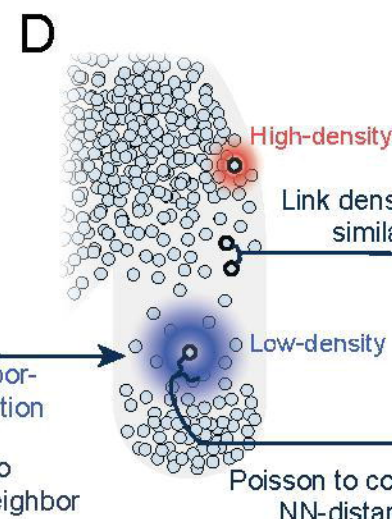
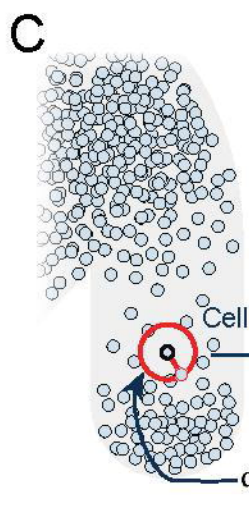
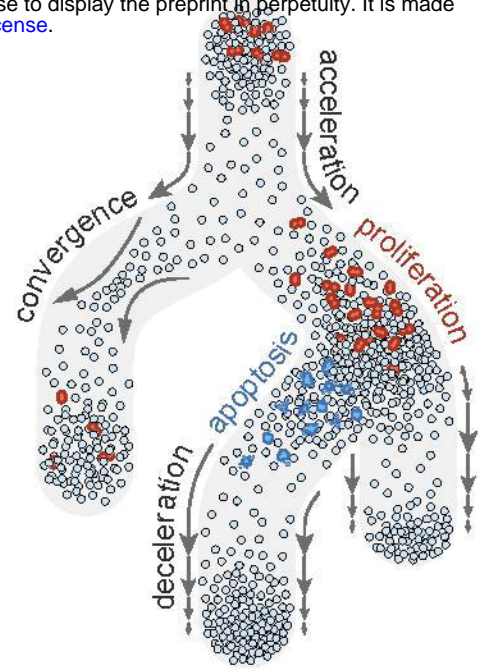
- 1580 24 Boller, S. & Grosschedl, R. The regulatory network of B-cell differentiation: a focused view of early
1581 B-cell factor 1 function. *Immunol Rev* **261**, 102-115 (2014). <https://doi.org/10.1111/imr.12206>
- 1582 25 Kim, H., Hwang, J. S., Lee, B., Hong, J. & Lee, S. Newly Identified Cancer-Associated Role of
1583 Human Neuronal Growth Regulator 1 (NEGR1). *J Cancer* **5**, 598-608 (2014).
1584 <https://doi.org/10.7150/jca.8052>
- 1585 26 Melchers, F. Checkpoints that control B cell development. *J Clin Invest* **125**, 2203-2210 (2015).
1586 <https://doi.org/10.1172/JCI78083>
- 1587 27 Atlas, H. C. (2020).
- 1588 28 Strunz, M. *et al.* Alveolar regeneration through a Krt8+ transitional stem cell state that persists in
1589 human lung fibrosis. *Nat Commun* **11**, 3559 (2020). <https://doi.org/10.1038/s41467-020-17358-3>
- 1590 29 Bastidas-Ponce, A. *et al.* Comprehensive single cell mRNA profiling reveals a detailed roadmap
1591 for pancreatic endocrinogenesis. *Development* **146** (2019). <https://doi.org/10.1242/dev.173849>
- 1592 30 Yang, D. *et al.* CRISPR screening uncovers a central requirement for HHEX in pancreatic lineage
1593 commitment and plasticity restriction. *Nat Cell Biol* **24**, 1064-1076 (2022).
1594 <https://doi.org/10.1038/s41556-022-00946-4>
- 1595 31 Moor, A. E. *et al.* Spatial Reconstruction of Single Enterocytes Uncovers Broad Zonation along
1596 the Intestinal Villus Axis. *Cell* **175**, 1156-1167 e1115 (2018).
1597 <https://doi.org/10.1016/j.cell.2018.08.063>
- 1598 32 Ma, S. *et al.* Chromatin Potential Identified by Shared Single-Cell Profiling of RNA and Chromatin.
1599 *Cell* **183**, 1103-1116 e1120 (2020). <https://doi.org/10.1016/j.cell.2020.09.056>
- 1600 33 Gonzalez, A. J., Setty, M. & Leslie, C. S. Early enhancer establishment and regulatory locus
1601 complexity shape transcriptional programs in hematopoietic differentiation. *Nat Genet* **47**, 1249-
1602 1259 (2015). <https://doi.org/10.1038/ng.3402>
- 1603 34 Lara-Astiaso, D. *et al.* Immunogenetics. Chromatin state dynamics during blood formation.
1604 *Science* **345**, 943-949 (2014). <https://doi.org/10.1126/science.1256271>
- 1605 35 Kaikkonen, M. U. *et al.* Remodeling of the enhancer landscape during macrophage activation is
1606 coupled to enhancer transcription. *Mol Cell* **51**, 310-325 (2013).
1607 <https://doi.org/10.1016/j.molcel.2013.07.010>
- 1608 36 Granja, J. M. *et al.* ArchR is a scalable software package for integrative single-cell chromatin
1609 accessibility analysis. *Nat Genet* **53**, 403-411 (2021). <https://doi.org/10.1038/s41588-021-00790-6>
- 1610 37 Argelaguet, R. *et al.* Decoding gene regulation in the mouse embryo using single-cell multi-omics.
1611 *bioRxiv*, 2022.2006.2015.496239 (2022). <https://doi.org/10.1101/2022.06.15.496239>
- 1612 38 Murre, C. 'Big bang' of B-cell development revealed. *Genes Dev* **32**, 93-95 (2018).
1613 <https://doi.org/10.1101/gad.311357.118>
- 1614 39 Sun, B. *et al.* Sox4 is required for the survival of pro-B cells. *J Immunol* **190**, 2080-2089 (2013).
1615 <https://doi.org/10.4049/jimmunol.1202736>
- 1616 40 Macnair, W., Gupta, R. & Claassen, M. psupertime: supervised pseudotime analysis for time-
1617 series single-cell RNA-seq data. *Bioinformatics* **38**, i290-i298 (2022).
1618 <https://doi.org/10.1093/bioinformatics/btac227>
- 1619 41 Tran, T. N. & Bader, G. D. Tempora: Cell trajectory inference using time-series single-cell RNA
1620 sequencing data. *PLoS Comput Biol* **16**, e1008205 (2020).
1621 <https://doi.org/10.1371/journal.pcbi.1008205>
- 1622 42 Mittnenzweig, M. *et al.* A single-embryo, single-cell time-resolved model for mouse gastrulation.
1623 *Cell* **184**, 2825-2842 e2822 (2021). <https://doi.org/10.1016/j.cell.2021.04.004>
- 1624 43 Klein, D. *et al.* Mapping cells through time and space with moscot. *bioRxiv*,
1625 2023.2005.2011.540374 (2023). <https://doi.org/10.1101/2023.05.11.540374>
- 1626 44 Pijuan-Sala, B. *et al.* A single-cell molecular map of mouse gastrulation and early organogenesis.
1627 *Nature* **566**, 490-495 (2019). <https://doi.org/10.1038/s41586-019-0933-9>
- 1628 45 McDole, K. *et al.* In Toto Imaging and Reconstruction of Post-Implantation Mouse Development
1629 at the Single-Cell Level. *Cell* **175**, 859-876 e833 (2018). <https://doi.org/10.1016/j.cell.2018.09.031>
- 1630 46 Wu, S. J. *et al.* Single-cell CUT&Tag analysis of chromatin modifications in differentiation and
1631 tumor progression. *Nat Biotechnol* **39**, 819-824 (2021). <https://doi.org/10.1038/s41587-021-00865-z>
- 1632
- 1633

- 1634 47 Bartosovic, M., Kabbe, M. & Castelo-Branco, G. Single-cell CUT&Tag profiles histone
1635 modifications and transcription factors in complex tissues. *Nat Biotechnol* **39**, 825-835 (2021).
1636 <https://doi.org/10.1038/s41587-021-00869-9>
- 1637 48 Zeller, P. *et al.* Single-cell sortChIC identifies hierarchical chromatin dynamics during
1638 hematopoiesis. *Nat Genet* **55**, 333-345 (2023). <https://doi.org/10.1038/s41588-022-01260-3>
- 1639 49 LaFave, L. M. *et al.* Epigenomic State Transitions Characterize Tumor Progression in Mouse Lung
1640 Adenocarcinoma. *Cancer Cell* **38**, 212-228 e213 (2020).
1641 <https://doi.org/10.1016/j.ccell.2020.06.006>
- 1642 50 Massague, J. & Ganesh, K. Metastasis-Initiating Cells and Ecosystems. *Cancer Discov* **11**, 971-
1643 994 (2021). <https://doi.org/10.1158/2159-8290.CD-21-0010>
- 1644 51 Lynch, A. W. *et al.* MIRA: joint regulatory modeling of multimodal expression and chromatin
1645 accessibility in single cells. *Nat Methods* **19**, 1097-1108 (2022). <https://doi.org/10.1038/s41592-022-01595-z>
- 1646 52 Meers, M. P., Llagas, G., Janssens, D. H., Codomo, C. A. & Henikoff, S. Multifactorial profiling of
1647 epigenetic landscapes at single-cell resolution using Multi-Tag. *Nat Biotechnol* **41**, 708-716
1648 (2023). <https://doi.org/10.1038/s41587-022-01522-9>
- 1649 53 Stuart, T. *et al.* Nanobody-tethered transposition enables multifactorial chromatin profiling at
1650 single-cell resolution. *Nat Biotechnol* **41**, 806-812 (2023). <https://doi.org/10.1038/s41587-022-01588-5>
- 1651 54 Regev, A. *et al.* The Human Cell Atlas. *Elife* **6** (2017). <https://doi.org/10.7554/eLife.27041>
- 1652 55 Rozenblatt-Rosen, O. *et al.* The Human Tumor Atlas Network: Charting Tumor Transitions across
1653 Space and Time at Single-Cell Resolution. *Cell* **181**, 236-249 (2020).
1654 <https://doi.org/10.1016/j.cell.2020.03.053>
- 1655 56 Google. JAX: composable transformations of Python+NumPy programs,
1656 <<http://github.com/google/jax>> (2018).
- 1657 57 Kumaraswamy, K. Fractal dimension for data mining. *Center for Automated Learning and*
1658 *Discovery School of Computer Science Carnegie Mellon University* **5000** (2003).
- 1659 58 Zhang, K., Tsang, I. W. & Kwok, J. T. in *Proceedings of the 25th international conference on*
1660 *Machine learning*. 1232-1239.
- 1661 59 Arthur, D. & Vassilvitskii, S. in *Proceedings of the eighteenth annual ACM-SIAM symposium on*
1662 *Discrete algorithms*. 1027-1035.
- 1663 60 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for differential
1664 expression analysis of digital gene expression data. *bioinformatics* **26**, 139-140 (2010).
- 1665 61 Cusanovich, D. A. *et al.* The cis-regulatory dynamics of embryonic development at single-cell
1666 resolution. *Nature* **555**, 538-542 (2018). <https://doi.org/10.1038/nature25981>
- 1667 62 Grant, C. E., Bailey, T. L. & Noble, W. S. FIMO: scanning for occurrences of a given motif.
1668 *Bioinformatics* **27**, 1017-1018 (2011). <https://doi.org/10.1093/bioinformatics/btr064>
- 1669 63 Wolf, F. A., Angerer, P. & Theis, F. J. SCANPY: large-scale single-cell gene expression data
1670 analysis. *Genome Biol* **19**, 15 (2018). <https://doi.org/10.1186/s13059-017-1382-0>
- 1671 64 Korsunsky, I. *et al.* Fast, sensitive and accurate integration of single-cell data with Harmony. *Nat*
1672 *Methods* **16**, 1289-1296 (2019). <https://doi.org/10.1038/s41592-019-0619-0>
- 1673 65 Persad, S. *et al.* (2022).
- 1674 66 Burrows, N. *et al.* Dynamic regulation of hypoxia-inducible factor-1alpha activity is essential for
1675 normal B cell development. *Nat Immunol* **21**, 1408-1420 (2020). <https://doi.org/10.1038/s41590-020-0772-8>
- 1676 67 Browaeys, R., Saelens, W. & Saeys, Y. NicheNet: modeling intercellular communication by linking
1677 ligands to target genes. *Nat Methods* **17**, 159-162 (2020). <https://doi.org/10.1038/s41592-019-0667-5>
- 1678 68 Lopez, R., Regier, J., Cole, M. B., Jordan, M. I. & Yosef, N. Deep generative modeling for single-
1679 cell transcriptomics. *Nat Methods* **15**, 1053-1058 (2018). <https://doi.org/10.1038/s41592-018-0229-2>
- 1680 69 Stephenson, E. *et al.* Single-cell multi-omics analysis of the immune response in COVID-19. *Nat*
1681 *Med* **27**, 904-916 (2021). <https://doi.org/10.1038/s41591-021-01329-2>
- 1682
- 1683
- 1684
- 1685
- 1686

A Uniform density



B Density in phenotypic landscapes



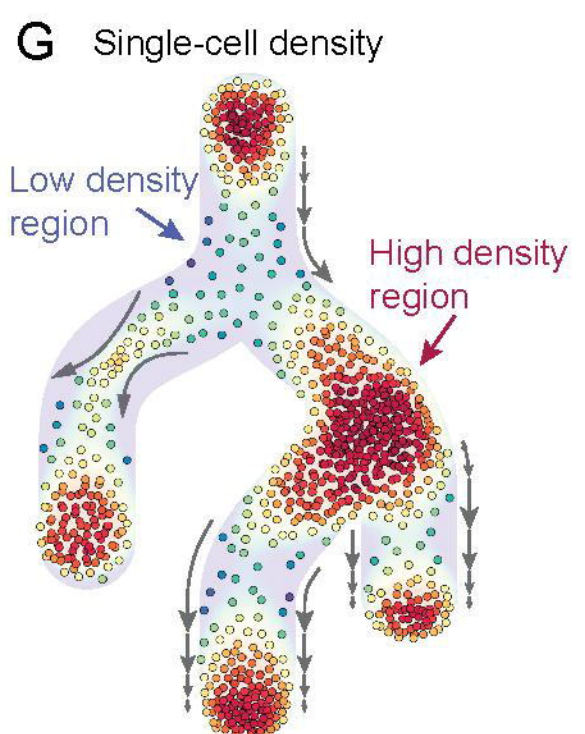
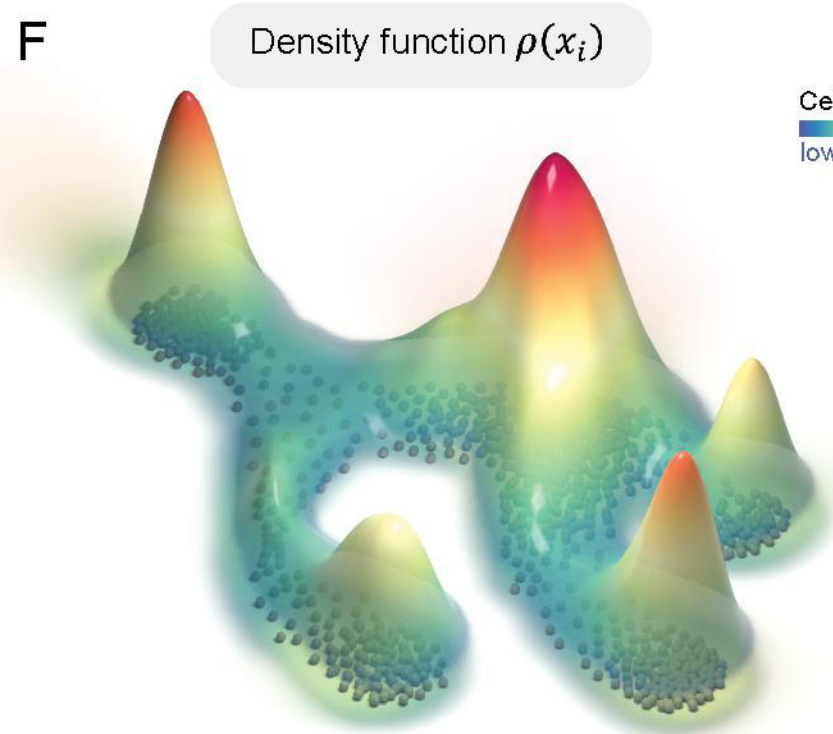
E Bayesian Model

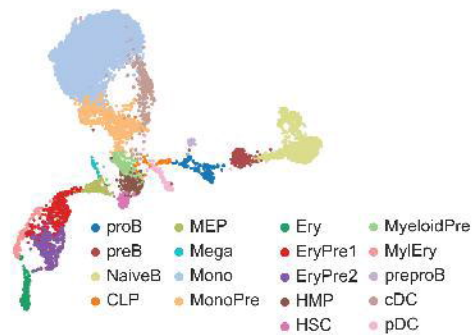
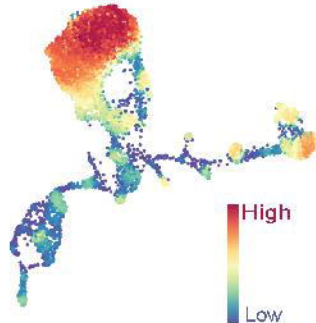
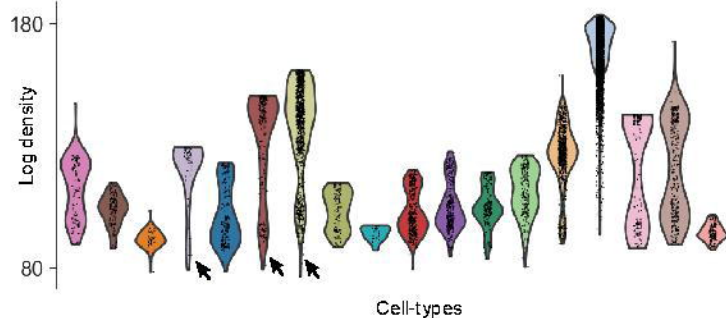
Cell state $x_i \in \mathbb{R}^d, i \in \{1 \dots n\}$

Log-density $f(x_i) \sim \text{GP}(m, \text{Matern52}(l))$

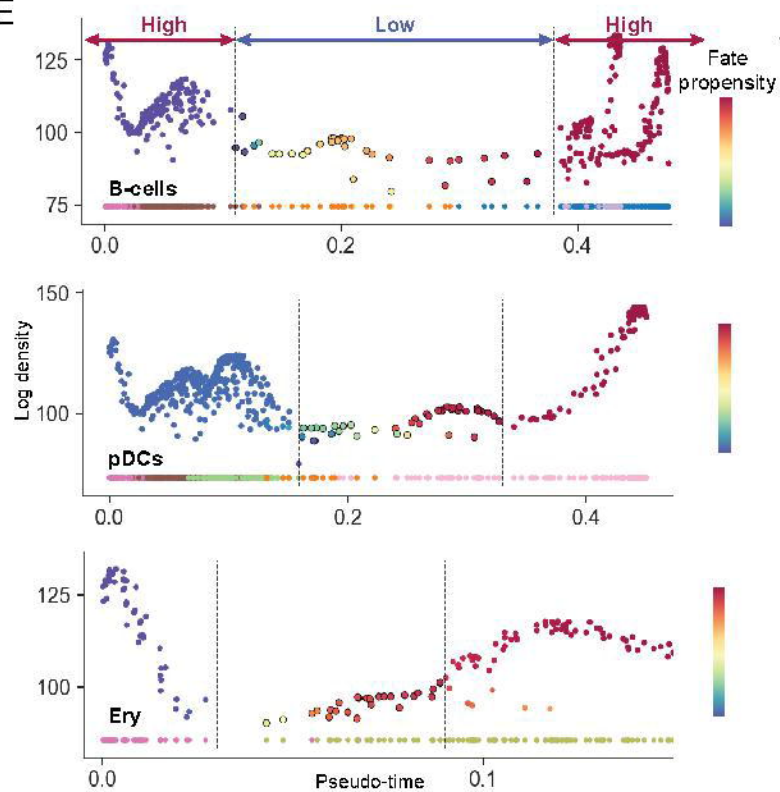
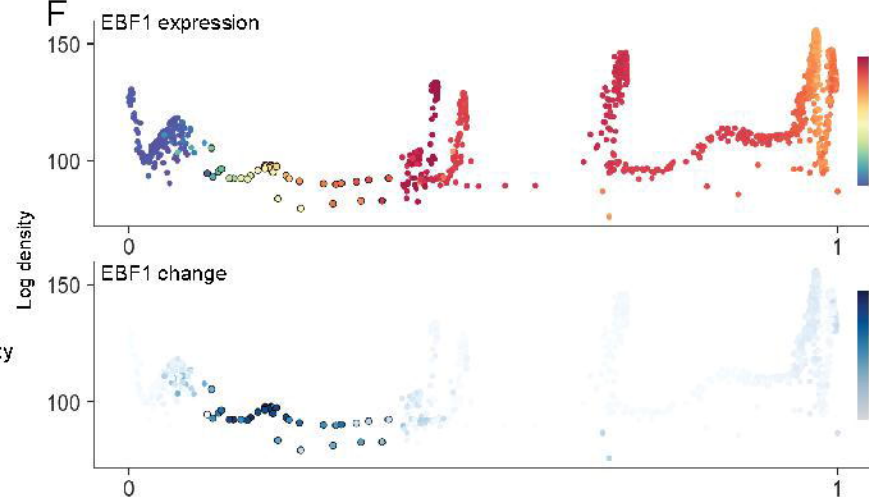
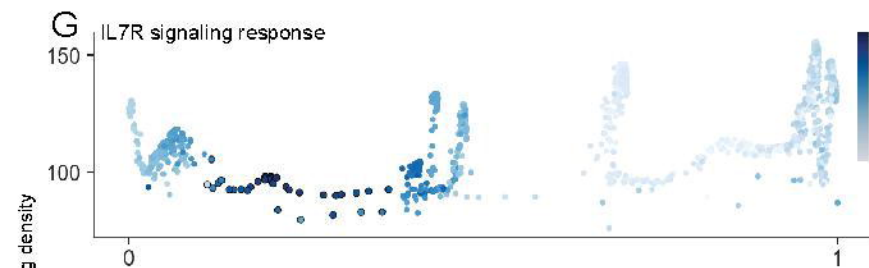
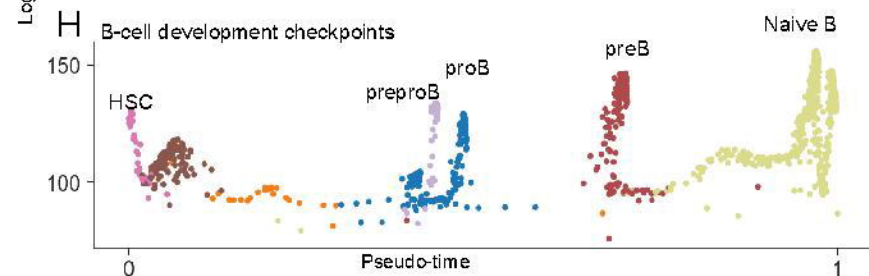
Density $\rho(x_i) = \exp \circ f(x_i)$

NN distance $dn(x_i) \sim \text{NN}(\rho(x_i), d)$

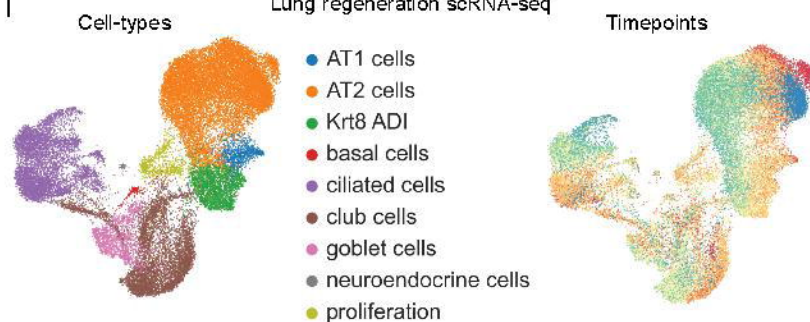


A Human hematopoiesis scRNA-seq**B** Mellon cell-state density**C****D**

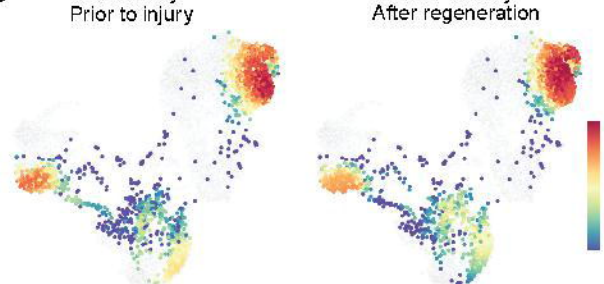
Hematopoietic lineage cells

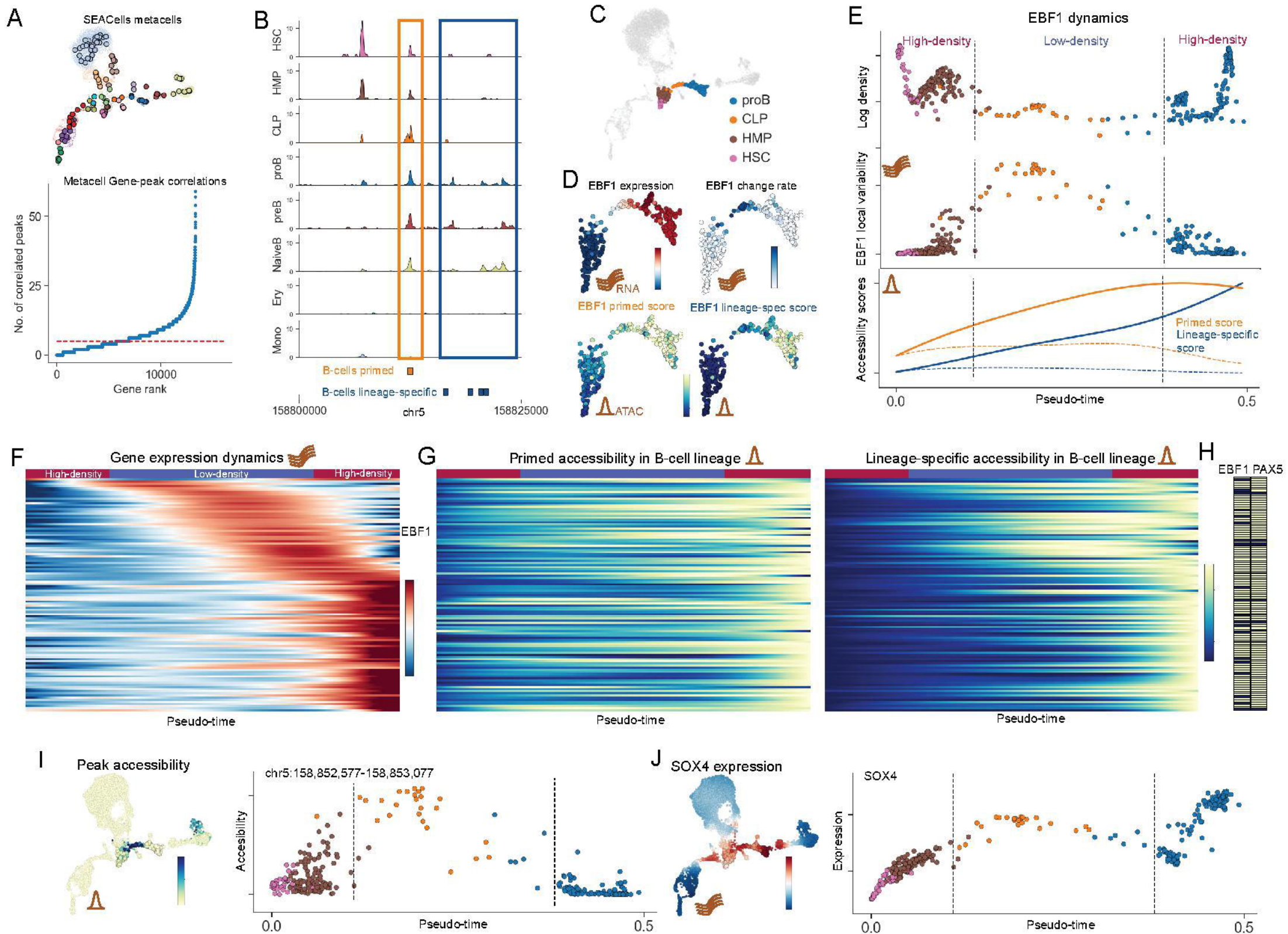
**E****F****G****H****I**

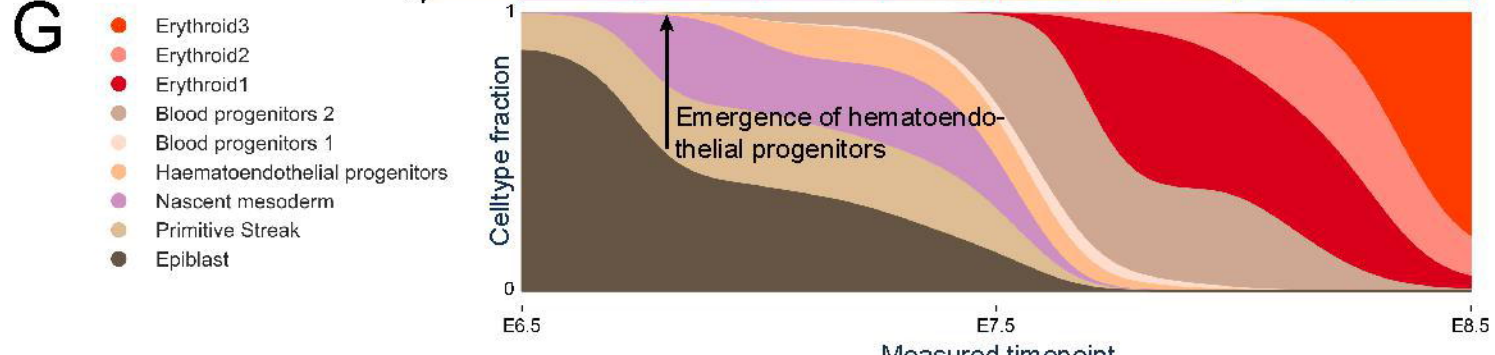
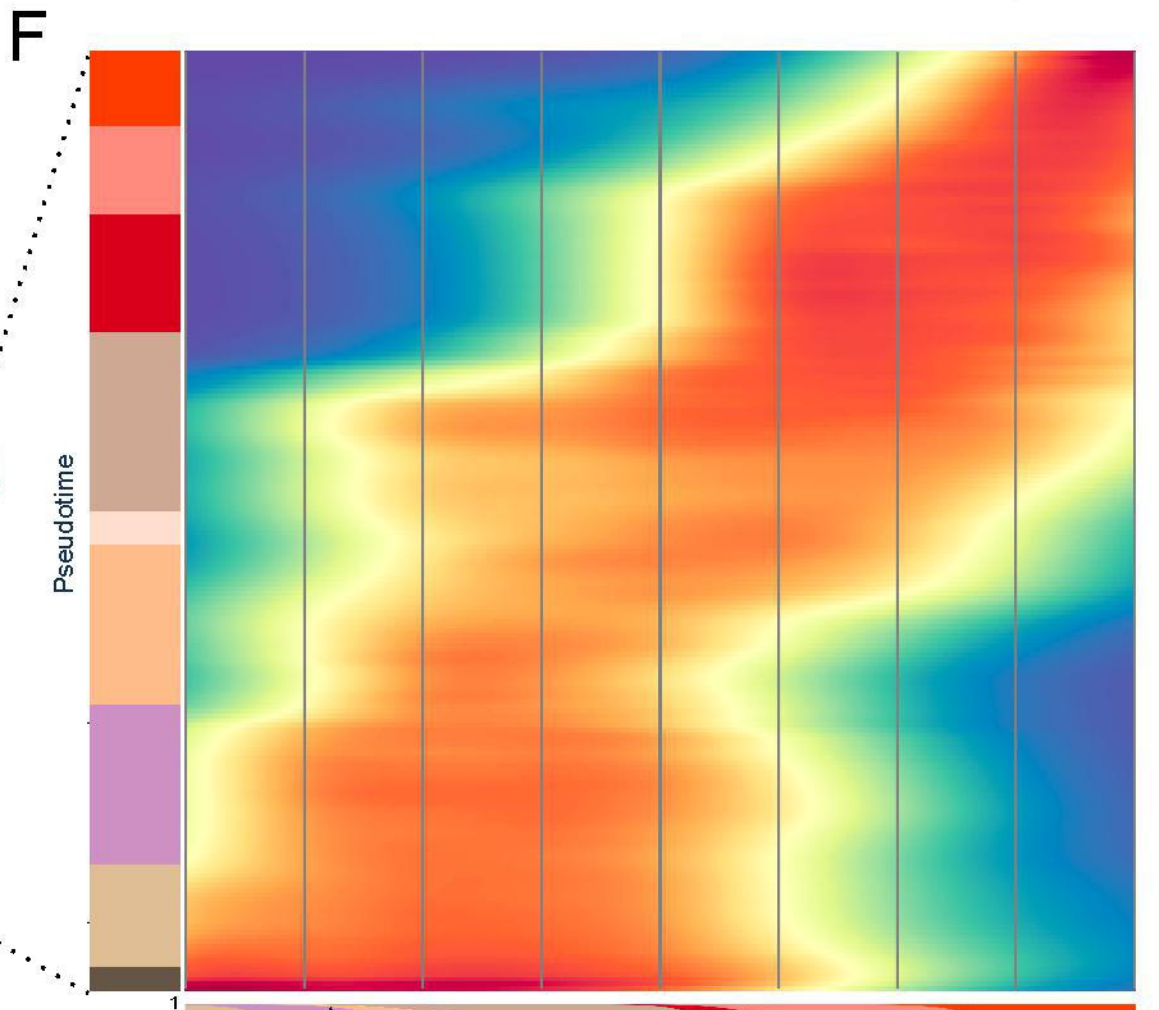
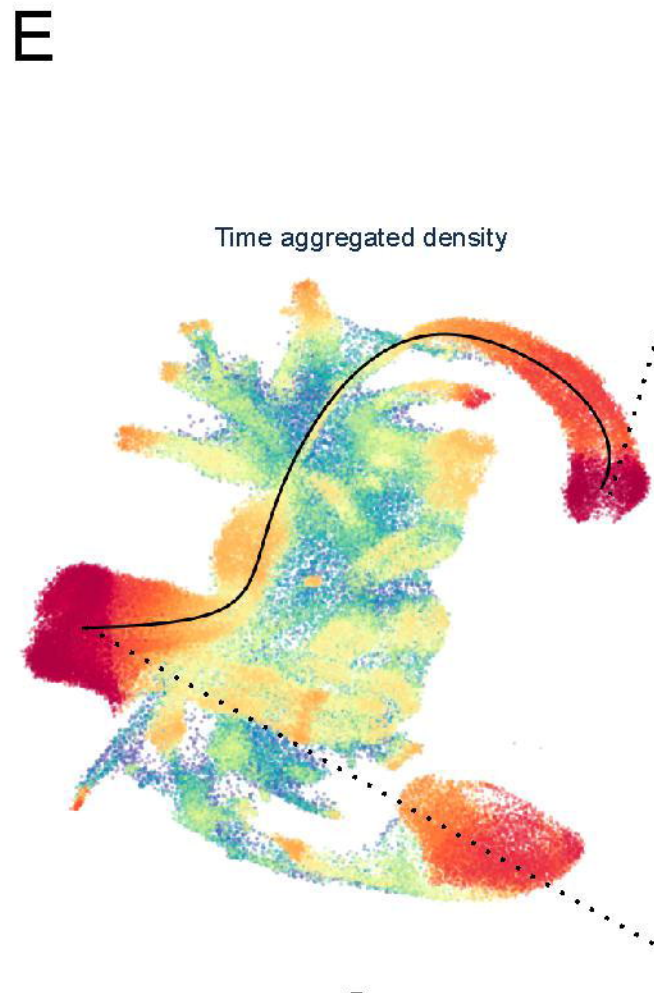
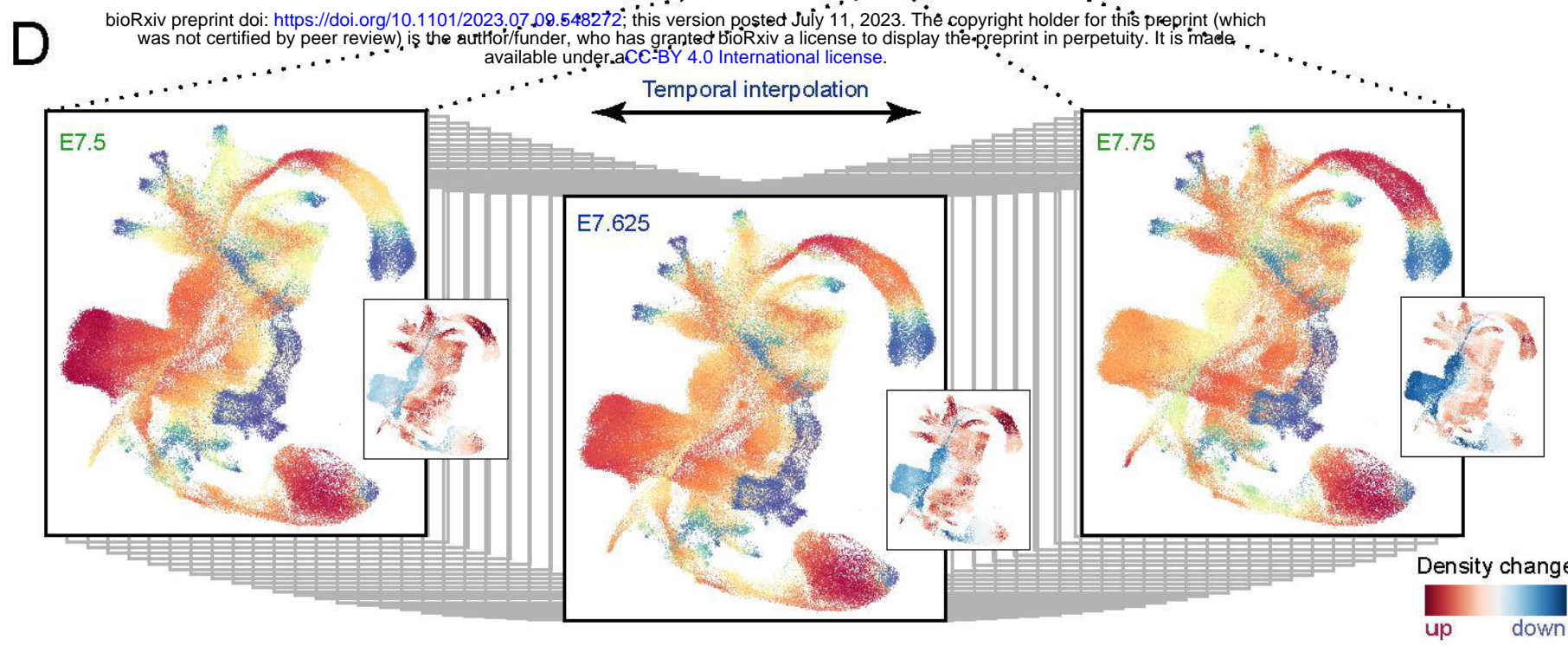
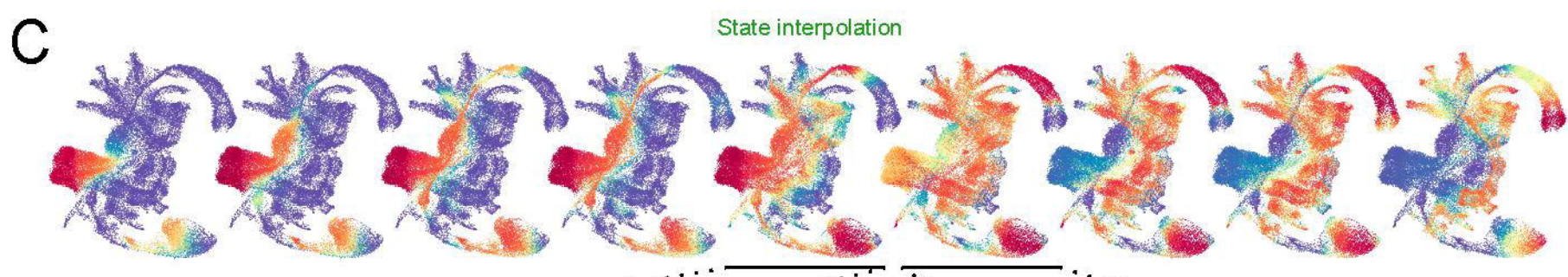
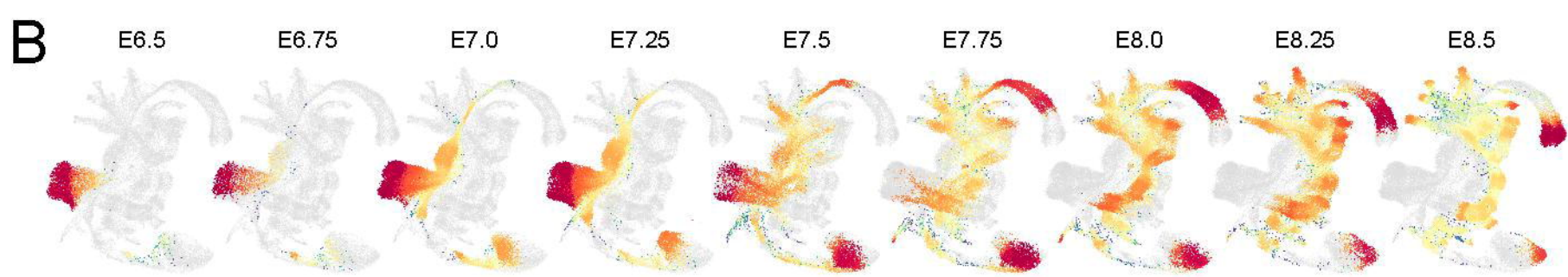
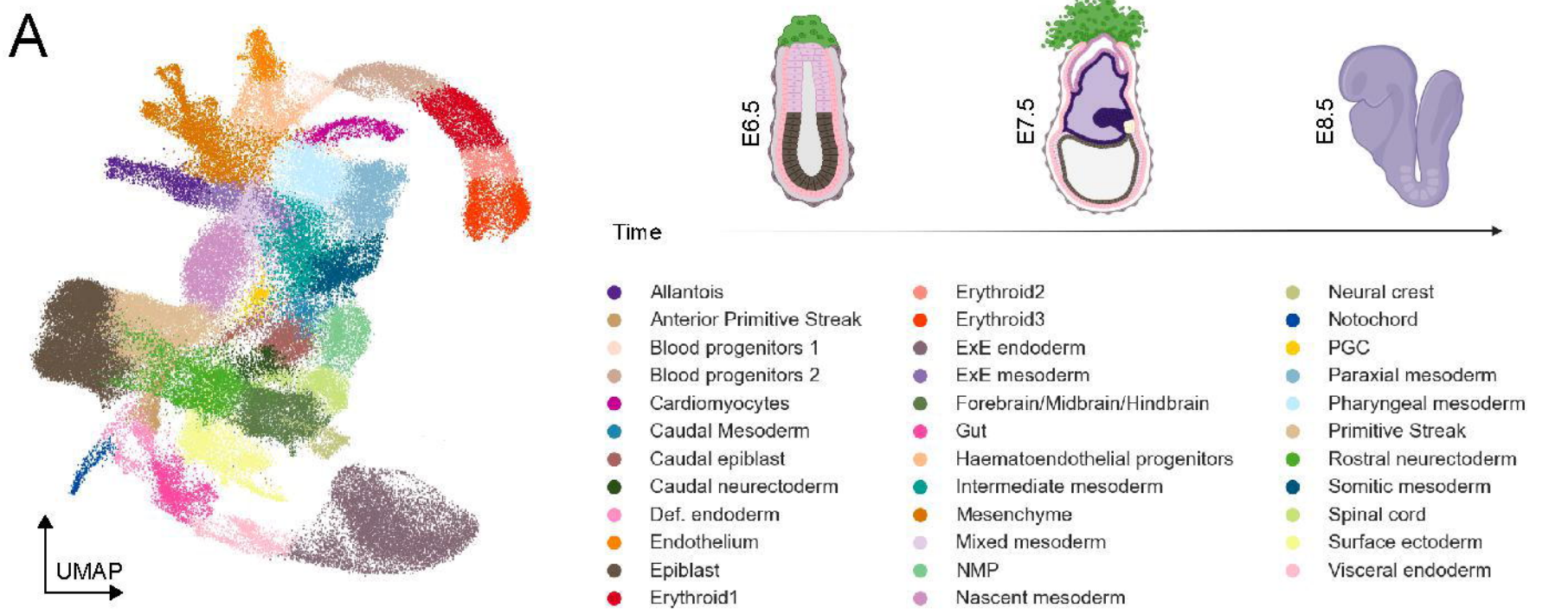
Lung regeneration scRNA-seq

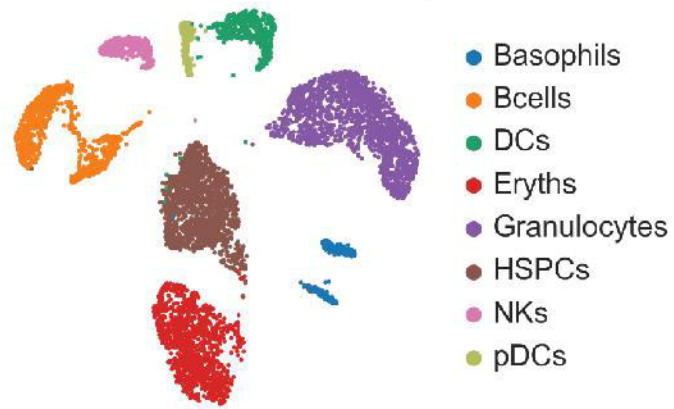
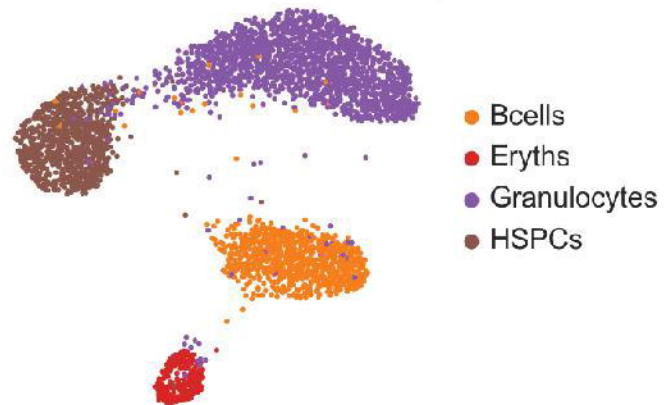
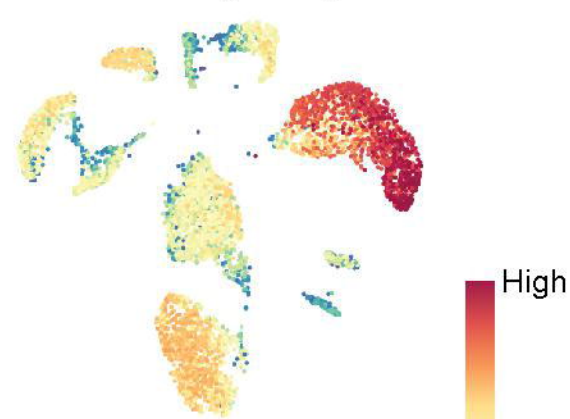
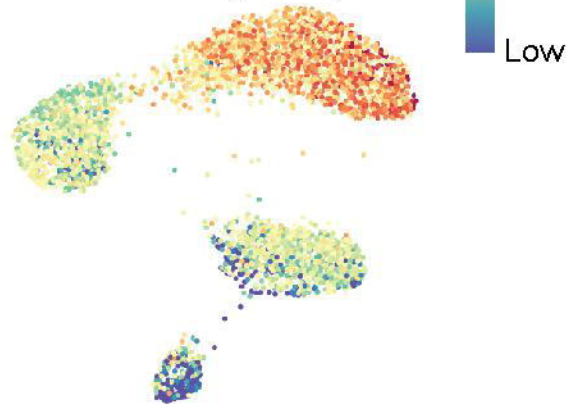
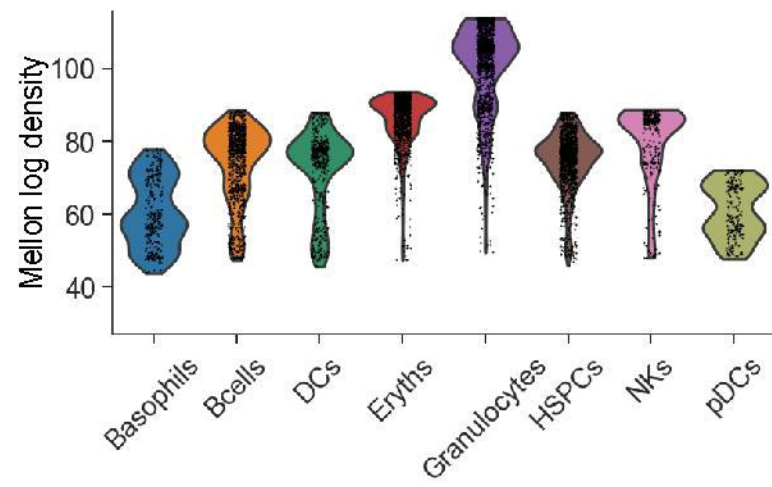
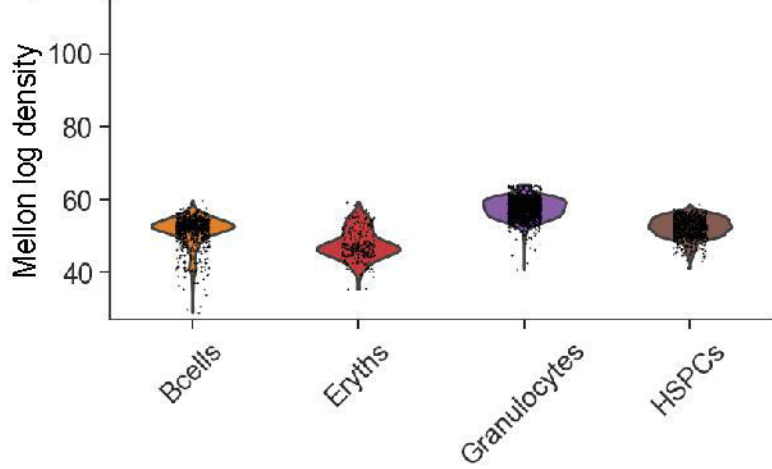
**J**

D0 density







A H3K4me1 - Murine hematopoiesis**B** H3K9me3 - Murine hematopoiesis**C** Mellon log density**D** Mellon log density**E****F****G**