

# ArchVelo: archetypal velocity modeling for single-cell multi-omic trajectories

Received: 30 September 2025

Accepted: 26 May 2026

Cite this article as: Avdeeva, M., Walker, S.K., Veeken, J. *et al.* ArchVelo: archetypal velocity modeling for single-cell multi-omic trajectories. *Nat Commun* (2026). <https://doi.org/10.1038/s41467-026-74000-4>

Maria Avdeeva, Sarah K. Walker, Joris van der Veeken, Alexander Y. Rudensky & Yuri Pritykin

We are providing an unedited version of this manuscript to give early access to its findings. Before final publication, the manuscript will undergo further editing. Please note there may be errors present which affect the content, and all legal disclaimers apply.

If this paper is publishing under a Transparent Peer Review model then Peer Review reports will publish with the final article.

# ArchVelo: Archetypal Velocity Modeling for Single-cell Multi-omic Trajectories

Maria Avdeeva<sup>1,\*</sup>, Sarah K. Walker<sup>2</sup>, Joris van der Veecken<sup>3</sup>,  
Alexander Y. Rudensky<sup>4</sup>, and Yuri Pritykin<sup>2,5,\*</sup>

<sup>1</sup> Center for Computational Biology, Flatiron Institute, Simons Foundation, New York, New York, USA

<sup>2</sup> Lewis-Sigler Institute for Integrative Genomics, Princeton University, Princeton, New Jersey, USA

<sup>3</sup> Research Institute of Molecular Pathology, Vienna, Austria

<sup>4</sup> Howard Hughes Medical Institute, Immunology Program, and Ludwig Center, Memorial Sloan Kettering Cancer Center, New York, NY, USA

<sup>5</sup> Department of Computer Science, Princeton, New Jersey, USA

\* Correspondence: mavdeeva@flatironinstitute.org, pritykin@princeton.edu

## Abstract

Inferring cellular dynamics from static single-cell data remains a central challenge in genomics. We introduce ArchVelo, a computational framework for modeling gene regulation and inferring trajectories from paired single-cell chromatin accessibility (scATAC-seq) and transcriptomic (scRNA-seq) data. ArchVelo represents chromatin accessibility as archetypes—shared regulatory programs—to model their dynamic influence on transcription. It outperforms existing methods in trajectory inference accuracy and gene-level latent time alignment, enables trajectory decomposition into archetypal components, and identifies the underlying transcription factors. After benchmarking on mouse brain and human hematopoiesis datasets, we apply ArchVelo to CD8 T cells in viral infection and reveal distinct trajectories of differentiation and proliferation. Focusing on progenitor exhausted CD8 T cells, critical for sustained immunity and immunotherapy response, we identify differentiation from  $Ccr6^-$  to  $Ccr6^+$  progenitors, shared between acute and chronic infections. ArchVelo provides a principled framework for modeling dynamic gene regulation and trajectory inference in multi-omic single-cell data across biological systems.

## Introduction

Single-cell transcriptomics (scRNA-seq) has been widely used to characterize cellular heterogeneity across diverse biological systems.<sup>1,2</sup> In scRNA-seq analysis, cells are represented as high-dimensional vectors of genome-wide gene expression, providing a quantitative multi-dimensional view of cellular phenotypes.<sup>3,4</sup> These data are commonly structured as a  $k$ -nearest-neighbor (kNN) graph based on transcriptomic similarity between cells, which can be used to identify discrete cell types and states through clustering, or to reconstruct continuous biological processes such as cell differentiation or response to perturbations via trajectory inference.

Trajectory inference methods assume that cells undergo dynamic transcriptomic changes driven by underlying molecular processes, and that scRNA-seq provides snapshots of these processes across many cells. By analyzing the transcriptomic similarity graph, these methods aim to infer the progression of cellular states and assign pseudotime coordinates to cells. However, most trajectory inference methods operate on undirected graphs and rely on additional assumptions or external information to define origins, branch points, or endpoints of trajectories.<sup>3,5</sup>

RNA velocity methods represent a powerful class of trajectory inference approaches that exploit unspliced and spliced RNA reads in scRNA-seq data to model transcriptional dynamics directly.<sup>6-15</sup> The ratio of unspliced to spliced reads provides a proxy for nascent versus mature transcript abundance, which can be modeled using ordinary differential equations (ODEs) to infer gene-specific regulatory kinetics. Aggregating these gene-level models allows estimation of a high-dimensional RNA velocity vector for each cell, enabling directional trajectory inference in gene expression space. While this framework has yielded valuable insights, it also presents challenges in terms of model robustness and interpretability.<sup>8,16</sup>

Although gene expression reflects the output of regulatory processes, upstream regulatory activity, such as transcription factor binding and chromatin accessibility, can provide additional, potentially more direct, insights into cell state dynamics. In particular, scATAC-seq data, which measure chromatin accessibility, offer a window into regulatory element activity and can complement transcriptional data. Recent advances in multi-omic technologies now enable simultaneous profiling of chromatin accessibility and gene expression in the same cells (scATAC+RNA-seq), offering a more comprehensive view of gene regulation in dynamic cellular contexts.<sup>1,17</sup> However, modeling transcriptional dynamics from scATAC+RNA-seq data introduces additional complexity. Compared to scRNA-seq, scATAC-seq data are higher-dimensional, sparser, and lack readily available references of chromatin accessibility regions for preprocessing and quantification (unlike reference gene annotations for scRNA-seq analysis that are well curated).<sup>3,18</sup> Moreover, the relationship between chromatin accessibility and gene expression is not fully understood. While a previously proposed method, MultiVelo,<sup>9</sup> pioneered RNA velocity analysis in scATAC+RNA-seq

data, it relies on aggregating chromatin accessibility across all peaks linked to a gene into a single value. This simplification obscures the complexity of gene regulation, where multiple distinct regulatory elements (enhancers and promoters) may respond to different upstream signals or compete to drive transcription. We argue that trajectory inference can be significantly improved by a more informative representation that captures these distinct regulatory programs rather than averaging them.

Here we introduce ArchVelo, a method for RNA velocity analysis and trajectory inference from scATAC+RNA-seq data. ArchVelo is based on archetypal analysis, which decomposes scATAC-seq profiles into a small number of interpretable archetypes representing “extreme” or characteristic regulatory programs within a dataset.<sup>19,20</sup> This low-rank representation of the chromatin accessibility modality improves RNA prediction performance and provides a biologically meaningful basis for modeling transcriptional dynamics. ArchVelo integrates this archetypal representation into a kinetic model that jointly describes the dynamics of chromatin accessibility and transcription. For each gene, transcriptional dynamics are modeled as a linear combination of contributions from the archetypes, enabling more accurate velocity estimation and trajectory inference. Using publicly available scATAC+RNA-seq data from developing mouse brain and human hematopoiesis, we rigorously benchmark ArchVelo against other state-of-the-art RNA velocity methods: scVelo,<sup>7</sup> MultiVelo,<sup>9</sup> DeepVelo,<sup>13</sup> VeloVI,<sup>12</sup> TFVelo<sup>14</sup> and Cell2Fate.<sup>15</sup> This benchmarking panel was chosen to comprehensively represent the major conceptual advances in the field. Specifically, it spans a widely adopted foundational method (scVelo), a model capturing cell-specific kinetics (DeepVelo), approaches leveraging multi-omic or transcription factor regulatory information (MultiVelo, TFVelo), deep network and deep generative models (DeepVelo, VeloVI), and a fully Bayesian model of complex transcription dynamics (Cell2fate). We show that ArchVelo achieves higher accuracy and improved latent time alignment for individual genes. Furthermore, ground-truth benchmarking confirms that ArchVelo-inferred trajectories more closely recapitulate known biological progressions. Beyond improved accuracy, ArchVelo provides a novel feature: decomposition of the RNA velocity field into archetype-specific components. By coupling this decomposition with transcription factor (TF) motif analysis, ArchVelo enables mechanistic interpretation of cellular trajectories and nominates candidate TFs as drivers of specific regulatory components. To demonstrate the biological utility of ArchVelo, we applied it to a multi-omic dataset profiling cytotoxic CD8 T cell responses to acute and chronic LCMV infection in mice.<sup>21</sup> ArchVelo revealed two major distinct trajectories corresponding to functional differentiation and cell proliferation. Focusing on the Tcf1<sup>+</sup> progenitor CD8 T cells, essential for sustained responses in cancer and infection and for cancer immunotherapies,<sup>22</sup> we uncovered a previously uncharacterized differentiation trajectory from Ccr6<sup>-</sup> to Ccr6<sup>+</sup> progenitors shared between acute and chronic infection.

In summary, ArchVelo is a robust and interpretable framework for RNA velocity analysis and trajectory inference from scATAC+RNA-seq data. By combining multi-omic modeling with archetypal decomposition, ArchVelo enhances the accuracy of dynamic inference and enables mechanistic insights into regulatory programs driving cell state transitions. The open-source ArchVelo package is available at <https://github.com/pritykinlab/ArchVelo>.

## Results

### Archetypal analysis of scATAC-seq improves scRNA-seq modeling

A single-cell multi-omic sequencing experiment includes simultaneous genome-wide measurements of transcriptome (scRNA-seq) and chromatin accessibility (scATAC-seq) in single cells (**Fig. 1A**).

The standard approach is to represent the scRNA-seq component as a matrix  $X$  of unique molecular identifier (UMI) or read counts for every gene in every cell, and to represent the scATAC-seq component as a matrix  $Y$  of read counts at peaks of chromatin accessibility in every cell. This scATAC-seq representation relies on a pre-constructed reference set of peaks, i.e. genomic regions with overall high level of scATAC-seq signal across cells. ATAC-seq peaks presumably correspond to sites of regulatory activity such as promoters, enhancers and repressors. Analysis of such data can be challenging due to high sparsity and dimensionality, and incomplete understanding of how regulatory activity at regions of accessible chromatin drives gene expression. The chromatin accessibility landscape across cells can be quite complex, even within the locus of the same gene, where peaks may exhibit heterogeneous cell type-specific accessibility potentially associated with cell type-specific expression (**Fig. 1B,C**).

To address these challenges, we use archetypal analysis of the scATAC-seq component of the multi-omic data (**Fig. 1D,E**).<sup>19,20</sup> Each archetype  $a_k$  is a profile of chromatin accessibility over cells. For a pre-defined number of archetypes, the signal at every cell is decomposed into archetypal scores (contained in matrix  $A$ ), and chromatin accessibility profile at every peak across cells is approximated by a convex combination of the archetypes (scores denoted by  $z_{pk}$  for peak  $p$ ). To explore the global genome-wide relationship between the scATAC-seq and scRNA-seq modalities in a multi-omic dataset, we formulated a regularized linear regression model to predict scRNA-seq from scATAC-seq, using archetypal featurization of scATAC-seq (**Fig. 1F,G**). We observed that this simple approach had better predictive power than other methods. This includes scATAC-seq gene activity scores from ArchR, a popular package for scATAC-seq analysis,<sup>23</sup> and a baseline regression approach using accessibility values in all peaks corresponding to a gene as features (**Methods**). The archetypal regression also performs significantly better than previously proposed deep learning method BABEL.<sup>24</sup> Thus the low-dimensional archetypal representation of scATAC-

seq data is productive for modeling the regulatory relationship between the scATAC-seq and scRNA-seq modalities in scATAC+RNA-seq data.

### ArchVelo for modeling gene expression regulation dynamics

The standard analysis of scRNA-seq data uses estimates of mRNA abundance, i.e. the transcript product of gene expression. However, scATAC+RNA-seq analysis provides an opportunity to model the process of gene expression regulation rather than only its end product. Therefore, we next turned to modeling dynamics of gene expression regulation, extending the RNA velocity framework. The RNA velocity ordinary differential equations (ODEs)<sup>6,7</sup> model the cascade of transcription, splicing and degradation of individual genes assuming piecewise constant transcription rates, with different constant rates in induction and repression phases. This simplifying assumption was extended for scATAC+RNA-seq data in MultiVelo,<sup>9</sup> where an additional equation modeling chromatin dynamics upstream of transcription was introduced, with the transcription rate proportional to the aggregated level of chromatin accessibility of a gene at any timepoint. However, this approach did not take full advantage of the resolution of the scATAC-seq data which measures chromatin accessibility at individual regulatory elements. Multiple regulatory programs can potentially drive transcription of a gene utilizing different subsets of its regulatory elements. To incorporate this information into the RNA velocity framework, we developed ArchVelo (**Fig. 2, Methods**) that extends MultiVelo equations using the archetypal decomposition of the scATAC-seq modality (**Fig. 1D**). ArchVelo models dynamics of chromatin accessibility archetypes  $a_k(t)$  for every gene  $g$  as a function of the underlying latent time  $t$ , incorporating individual transcription rates  $\alpha_k^g$  for every archetype, archetype accessibility score ( $z_k^g = \sum_{p \in \mathcal{P}(g)} z_{pk}$  where  $\mathcal{P}(g)$  denotes the peaks mapping to  $g$ ) and preserving the splicing and degradation model components  $\beta^g$  and  $\gamma^g$ , respectively. The transcription rates are chosen to depend on transcription state  $\tau_k^g$ . The flexible kinetic model of ArchVelo (**Fig. 2**), with dynamics constrained via joint scATAC-seq archetypal decomposition, could provide more accurate trajectory modeling than MultiVelo and give insight into dynamics of the archetypal regulatory programs driving transcription. To assess this and benchmark ArchVelo against existing methods, we next applied it to multi-omic datasets across biological systems.

### Improved trajectory inference in the mouse embryonic brain

Similar to other RNA velocity frameworks, ArchVelo combines gene-specific dynamic models to infer a global velocity field (**Methods**). To demonstrate its functionality and benchmark performance against other methods, we used a published scATAC+RNA-seq dataset for the embryonic mouse brain (**Fig. 3A**).<sup>25</sup>

We reanalyzed this dataset, identifying clusters corresponding to major developmental populations, annotated using marker gene expression (**Fig. 3B, Supplementary Fig. 1A, Methods**). These included intermediate progenitor cells (IPC), two subpopulations of migrating neurons (MN1, MN2), upper-layer (UL) and deeper-layer (DL) neurons. In the glial compartment, we identified astrocyte progenitors (AP), expressing markers of radial glia, and a subpopulation of astrocytes. In addition, we detected two distinct clusters of cycling cells: one enriched for astrocytic and glial markers (AC) and another (C) comprised of both glial lineage cells including oligodendrocyte progenitors (OPC) and neuronal progenitors such as IPCs. ArchVelo, which was not used for clustering or cell type annotations, revealed a separation of trajectories in the glial and neuronal compartment. Within the neuronal compartment, ArchVelo correctly inferred the separation of the upper and deeper layer lineages.<sup>26,27</sup> In particular, UL neurons were predicted to arise from IPC via MN1, a migrating intermediate state expressing upper layer markers such as *Satb2*, *Cux2* and *Plxna4*, while most DL neurons were inferred to originate from MN2, the migrating neuronal subpopulation expressing deeper-layer markers such as *Bcl11b* and *Fezf2*.

We next compared the performance of ArchVelo with other state-of-the-art RNA velocity methods on this dataset (**Fig. 3C–H, Supplementary Fig. 1B–D**). As a result of its kinetic flexibility, ArchVelo produced more accurate fits of gene expression dynamics in the  $u$ - $s$  phase space than MultiVelo (**Fig. 3C,D**). Both methods infer gene-specific latent times, but in ArchVelo, latent time estimation is constrained by shared chromatin archetypes, which improves consistency across genes. Consequently, ArchVelo produced better-aligned latent time profiles (**Fig. 3E,F, Methods**). Overall, owing to this shared structure and richer kinetic modeling, ArchVelo simultaneously achieved higher model likelihood and greater latent time robustness than MultiVelo across a wide range of the  $w_c$  parameters controlling ATAC-seq versus RNA-seq contributions, and also outperformed other methods with respect to these metrics (**Fig. 3G, Supplementary Fig. 1B**). Cross-boundary direction correctness (CBDir) analysis<sup>28</sup> using the curated reference cell type transitions revealed that ArchVelo overall achieved higher scores than other methods (**Fig. 3H, Supplementary Fig. 1C,D, Methods**). Thus, ArchVelo more accurately recovered the differentiation dynamics in the developing brain than previous methods.

## Human hematopoietic stem cell differentiation trajectories

We next applied ArchVelo to a multi-omic dataset of differentiating human hematopoietic stem cells (HSCs) previously generated in the MultiVelo study.<sup>9</sup> In this experiment, purified CD34<sup>+</sup> HSCs were cultured for seven days before sequencing. Our reanalysis revealed cell populations consistent with known hematopoietic differentiation hierarchies and the expression of lineage-specific markers (**Fig. 4A,B, Methods**).<sup>29–35</sup>

In particular, we identified a subpopulation of HSCs closely resembling the multipotent progenitors (MPP), which seeded all major hematopoietic differentiation trajectories. Downstream differentiation branches included lymphoid-primed multipotent progenitors (LMPP), granulocyte–macrophage progenitors (GMP), progenitor neutrophils (N) and progenitor dendritic cells (DC). We also identified a subpopulation expressing intermediate levels of markers associated with erythroid (*TFRC*, *TFR2*), basophil-eosinophil-mast cell (*CPA3*), or both (*RHEX*) lineages. This mixed transcriptional profile, together with expression of *GATA2*, is consistent with erythroid-myeloid progenitors (EMP), a previously characterized early hematopoietic subpopulation.<sup>32,34,35</sup> Additional clusters resembled megakaryocyte-erythroid progenitors (MEP), basophil-eosinophil-mast cell progenitors (BEM), progenitor megakaryocytes (MK) and platelets, as well as early erythrocytes (Ery).

We next compared ArchVelo with other methods using this dataset (**Fig. 4C–F, Supplementary Fig. 2**). ArchVelo achieved higher per-gene log-likelihood and greater consistency of per-gene latent time profiles than other methods (**Fig. 4C–E, Supplementary Fig. 2A**). Then we constructed a curated reference graph of hematopoietic differentiation for the identified subpopulations (**Fig. 4F, Methods**).<sup>31–35</sup> Notably, BEMs appear transcriptomically closer to erythroid and megakaryocyte progenitors than to other hematopoietic lineages.<sup>30,36</sup> Furthermore, early bifurcation between the monocyte–neutrophil–lymphocyte compartment and the megakaryocyte–erythroid–basophil–eosinophil–mast compartment has been proposed.<sup>30,32–34</sup> In the latter branch, since BEMs have been reported to arise from EMPs<sup>32,34</sup> and from MEPs,<sup>30</sup> we incorporated both possibilities into the ground truth graph. Evaluating trajectories inferred by ArchVelo against this curated reference, we found that ArchVelo successfully reconstructed most expected lineage relationships. Furthermore, on average ArchVelo achieved higher CDBir scores than other methods across cell type transition edges of the reference differentiation graph (**Fig. 4F, Supplementary Fig. 2B–D**). Thus ArchVelo provides a more accurate and robust representation of hematopoietic differentiation dynamics than other RNA velocity methods.

Together these comprehensive benchmarking results demonstrate that ArchVelo outperforms previous methods in trajectory inference in single-cell multi-omic scATAC+RNA-seq data.

## Velocity decomposition into archetypal components

Beyond improving trajectory inference, ArchVelo provides a unique capability to decompose transcriptional dynamics into interpretable regulatory components (**Fig. 5A–C, Methods**). This is enabled by the linear structure of the underlying ODE model, where transcriptional variables  $u^g(t)$  and  $s^g(t)$  for each gene  $g$  can be expressed as linear combinations of archetype-specific contribu-

M. Avdeeva et al.

tions:

$$u^g(t) = \sum_k u_k^g(t), \quad s^g(t) = \sum_k s_k^g(t),$$

where  $u_k^g(t)$  and  $s_k^g(t)$  denote the components driven by archetype  $a_k(t)$ . Consequently, the gene-specific velocity

$$v^g(t) = \dot{s}^g(t)$$

is similarly decomposable into archetypal components:

$$v^g(t) = \sum_k v_k^g(t), \quad v_k^g(t) = \dot{s}_k^g(t).$$

Each component represents the portion of a gene's expression and velocity attributable to regulatory elements associated with a particular chromatin accessibility archetype.

Aggregating these contributions across all genes enables ArchVelo to extract distinct components of the global velocity field, each potentially corresponding to a specific regulatory program (**Methods**). This decomposition is computed at the single-cell level, and the resulting components can be directly projected onto the UMAP embedding, revealing archetype-specific trajectories that can be independently visualized and interpreted.

For example, in the embryonic mouse brain dataset, two archetypes A5 and A8 produced velocity fields concentrated in the glial compartment (**Fig. 5D–J, Supplementary Fig. 3**). Although they overlapped in cell coverage, the inferred dynamics were clearly distinct. Archetype A5 was enriched in cycling cells and associated with cell cycle-related genes such as *E2f8* and *Cdc25c* (**Supplementary Fig. 3A,B**). Consistent with this, the latent time inferred from A5 aligned with progression through cell cycle phases (**Fig. 5E–G**). In contrast, archetype A8 encompassed both cycling and non-cycling cells in the glial compartment and was enriched for astrocyte lineage markers such as *Aldh1l1* (**Fig. 5H–J, Supplementary Fig. 3C,D**). The corresponding velocity field suggested differentiation trajectories toward astrocytes, supported by the progressive upregulation of *Aldh1l1*, *Slc6a11*, and *Sparcl1* during astrocyte lineage specification. Thus, ArchVelo disentangled proliferation and differentiation components within the same precursor population, clarifying trajectories of astrocyte-committed glioblasts.

These findings highlight how ArchVelo archetypal decomposition disentangles the superposition of multiple regulatory influences that simultaneously shape the trajectory of a cell. In conventional trajectory analyses, the observed dynamics represent an aggregate view of all underlying forces, which can obscure distinct processes such as cell cycle progression and lineage commitment. By resolving the velocity field into archetype-specific components, ArchVelo isolates these regulatory contributions, enabling separate interpretation of concurrent programs and clarifying how different forces drive cells in potentially divergent directions.

In the HSC dataset, decomposition of the velocity field highlighted archetypes A5, A4 and A8, which revealed diverging trajectories from the progenitor MEP subpopulation toward megakaryocyte (Prog MK), BEM and erythroid (Ery) progenitors, respectively (**Fig. 6A,B**). TF motif analysis for these archetypes uncovered top factors potentially driving these differentiation trajectories (**Fig. 6C–E, Methods**). For example, the top motif for archetype A5 corresponded to factors ERG, ETV3, and FLI1, with FLI1 being a well-established driver of megakaryocyte identity.<sup>37</sup> Consistent with this role, FLI1 motif scores, expression, and the expression of its downstream target ITGA2B showed expected patterns (**Fig. 6C**).<sup>38</sup> Analogously, GATA2 emerged as the top regulator of BEM cell differentiation (archetype A4, **Fig. 6D**), while the GATA family and KLF1 scored highest for the erythroid trajectory (archetype A8, **Fig. 6E**), consistent with their well-established roles in these lineage specifications.<sup>39,40</sup> Together, these results validate the ability of ArchVelo decomposition to link distinct velocity components to candidate TFs driving lineage-specific regulatory programs. Moreover, as illustrated in the hematopoiesis analysis, individual archetypal components can in some cases correspond to *bona fide* lineage bifurcations, rather than merely overlapping regulatory influences, thereby directly uncovering the branching structure of differentiation trajectories.

In sum, decomposition of ArchVelo trajectories into multiple, potentially divergent trajectories driven by different regulatory programs within the same cells, represented by scATAC-seq archetypes, provides a conceptually new view of regulatory dynamics beyond cell trajectories.

### Trajectories in CD8 T cells responding to viral infection

We next applied ArchVelo to a distinct biological system, cytotoxic CD8 T cells responding to viral infection. Specifically, we analyzed a multi-omic scATAC+RNA-seq dataset collected at day 7 post-infection with either the Armstrong (acute) or clone 13 (chronic) strains of lymphocytic choriomeningitis virus (LCMV). This dataset was generated as part of our larger recent study characterizing splenic T cell populations across infection conditions, building on our prior work.<sup>21,41,42</sup>

Using scATAC-seq accessibility signatures and gene activity profiles, we first identified antigen-experienced CD8 T cells, the population most relevant for studying infection-induced dynamics (**Methods**).<sup>21</sup> Within this compartment, we refined cell state annotations, resolving multiple subpopulations consistently observed in both LCMV Armstrong and clone 13 infections (**Fig. 7A–C**). Marker gene expression confirmed activated effector (e.g. *Klrg1*, *Ccl5*, *Gzmb*, *Gzmk*, *Gzma*), memory-like effector (e.g. *Il7r*), and early exhausted (e.g. *Pdcd1*, *Tox*) states, as well as intermediate subpopulations, alongside distinct cycling subsets characterized by *Mki67*, *Top2a*, *Pola1* and *Pcna* (**Fig. 7A,B**). Cycling cells were largely devoid of terminal differentiation markers, suggesting independent regulatory programs. We also identified a progenitor population defined by expression of *Tcf7* (encoding protein Tcf1), *Cxcr5*, *Slamf6*, *Id3*, *Sell* and *Il7r* but lacking activation-associated

genes *Gzmb*, *Ccl5* and *Prf1*. This well-studied subset, originally called “progenitor exhausted” and observed in contexts of chronic infection and cancer, has been implicated in long-term immune responses and linked to immunotherapy outcomes.<sup>22</sup> Consistent with recent work by us and others, we found the progenitors present in both acute and chronic infection (**Fig. 7A–C**).<sup>42–44</sup> Separately, we resolved an effector subpopulation with high expression of activation genes (*Gzmb*, *Ccl5*, *Prf1*) but also enriched for the memory marker *Il7r*, approaching levels seen in progenitors. Importantly, all identified subsets were shared between acute and chronic responses (**Fig. 7A–C**). This refined annotation provided the foundation for applying ArchVelo to dissect dynamic trajectories in infection-induced CD8 T cell responses.

ArchVelo revealed two major largely orthogonal sets of trajectories: one corresponding to cell cycle progression and the other to functional T cell differentiation (**Fig. 7A,D,E**). Transitions between cell states were consistent across both infection conditions and could be further illustrated by phase plots of ArchVelo fits for individual genes (**Fig. 7D,E, Supplementary Fig. 4, Methods**). For example, *Mki67* dynamics aligned with transitions through cell cycle phases, while *Sell* expression captured differentiation from  $Il7r^+$  effectors into progenitors. In sum, ArchVelo identified distinct trajectories of CD8 T cell differentiation and proliferation that were largely shared between responses to LCMV Armstrong and clone 13.

To further leverage the resolution of ArchVelo, we focused on the CD8 T cell progenitor compartment, a population of major translational interest. Within this compartment, we identified two subpopulations of  $Ccr6^-$  and  $Ccr6^+$  progenitor CD8 T cells, present in both Armstrong and clone 13 infections (**Fig. 8A,B, Supplementary Fig. 5A**). Both subsets expressed canonical progenitor markers, but also exhibited distinct gene expression profiles, indicating functional heterogeneity. This observation is consistent with previous reports of analogous  $Ccr6^-$  vs.  $Ccr6^+$  progenitor subsets in cancer, as well as  $CD62L^+$  vs.  $CD62L^-$  progenitors in chronic infection, whose gene signatures resembled those of our  $Ccr6^-$  and  $Ccr6^+$  subsets, respectively (**Supplementary Fig. 5B**).<sup>45–47</sup> In cancer, the  $Ccr6^+$  subset was reported to respond poorly to checkpoint blockade, with functional responsiveness instead attributed to the less differentiated  $Ccr6^-$  subset.<sup>45</sup> Consistent with this, in LCMV infection we observe that  $Ccr6^-$  progenitors differentiate to the  $Ccr6^+$  progenitors (**Fig. 8A,B, Supplementary Fig. 5C**).

Analysis of the archetypal loadings reinforced the distinction between the  $Ccr6^-$  and  $Ccr6^+$  progenitor subsets and highlighted their associated velocity components A8 and A5, respectively (**Fig. 8C,D**). Archetype A8 was associated with the transition from  $Ccr6^-$  to  $Ccr6^+$  progenitors, whereas archetype A5 appeared to capture further differentiation of  $Ccr6^+$  progenitors and was weaker in the  $Ccr6^-$  progenitors. Using archetypal latent times, we reconstructed a shared differentiation trajectory across both infection contexts (**Fig. 8E**). This analysis revealed smooth and coherent gene expression dynamics, encompassing many of the differentially expressed genes

(**Fig. 8B**), and supported a continuous progression between progenitor subsets. TF motif analysis of these archetypal components identified candidate factors driving this progression (**Fig. 8F**). In particular, archetype A8, associated with the less differentiated  $Ccr6^-$  progenitors, was enriched for Smad and E2f family motifs, consistent with regulatory programs maintaining proliferative potential.<sup>48–51</sup> In contrast, archetype A5, linked to further differentiation of  $Ccr6^+$  progenitors, showed enrichment of Nfat and Nr4a family motifs, factors with well-established roles in T cell activation and effector differentiation.<sup>52–55</sup> Thus using ArchVelo we identified a previously undescribed differentiation from  $Ccr6^-$  to  $Ccr6^+$  progenitor CD8 T cells in both acute and chronic LCMV infection.

Together, these results demonstrate that ArchVelo enables fine-grained dissection of overlapping dynamic programs such as proliferation and differentiation, both hallmarks of active T cell response *in vivo* to antigen stimulation. By uncovering shared trajectories across infection states, ArchVelo offers a powerful framework for understanding cellular decision-making.

## Discussion

We presented a new computational method ArchVelo for improved modeling of gene expression regulation dynamics and trajectory inference using multi-omic scATAC+RNA-seq data. Methodologically, two elements were critical. First, representing scATAC-seq as archetypes mitigates sparsity and shares information across regulatory elements with similar cell-wise activity, yielding stronger RNA prediction and a more stable regulatory basis for kinetic modeling. Second, embedding this representation into a gene-wise ODE system lets each archetype contribute linearly to transcription, providing flexible kinetics without sacrificing interpretability. Together, these choices improved quantitative performance (higher  $u-s$  likelihoods; more coherent latent times) and produced velocity fields that better captured expected developmental progressions. We applied ArchVelo across three biological systems, embryonic mouse brain, human hematopoiesis, and virus-responding CD8 T cells, recovering expected lineage progressions in the first two and, in T cells, delineating orthogonal proliferation and differentiation programs and a previously uncharacterized trajectory from  $Ccr6^-$  to  $Ccr6^+$  progenitors across acute and chronic infection.

While ArchVelo currently relies on simultaneous multi-omic profiling to learn gene-specific regulatory weights, the underlying archetypal framework is effective for summarizing chromatin accessibility landscapes even without using scRNA-seq data, as exemplified in our recent study.<sup>21</sup> Future extensions could adapt ArchVelo to unpaired datasets by leveraging computational integration of separately profiled scATAC-seq and scRNA-seq data, or matching bulk ATAC-seq data, thereby expanding its applicability to the wealth of single-cell genomics studies.

M. Avdeeva et al.

We also combined ArchVelo archetypal trajectory decomposition with TF motif analysis, enabling nomination of candidate TFs associated with distinct regulatory components. While these associations are correlative, they provide testable hypotheses about TF programs that couple chromatin dynamics to transcriptional change and may guide perturbation studies. These results highlight the potential of integrating regulatory analysis with velocity modeling. In future work, ArchVelo could be extended with a more seamless framework to systematically connect TF activity, chromatin accessibility dynamics, and transcriptional trajectories in multi-omic data.

A current limitation of the ArchVelo framework is the assumption of a linear mapping between chromatin accessibility archetypes and transcriptional activation rates. While this linear approximation ensures mathematical tractability and protects against kinetic overfitting, *in vivo* gene regulation frequently involves complex, non-linear dynamics driven by cooperative TF binding and intricate enhancer-promoter logic. Furthermore, epigenetic and transcriptional decoupling can occur in disease contexts.<sup>56</sup> Future iterations of the model could address this by incorporating non-linear mapping and jointly analyzing samples from across conditions, to more accurately capture the full complexity of these regulatory relationships.

ArchVelo is a method in the broader family of trajectory inference methods built on the RNA velocity framework. Many such approaches proceed in two stages: first, fitting a dynamic model for each gene independently by leveraging variation across cells; second, aggregating these gene-level fits into high-dimensional velocity fields that are then combined into trajectories. In ArchVelo, archetypal representation of scATAC-seq data improves both stages by stabilizing gene-wise modeling and producing more coherent velocity fields. Looking forward, this concept could be extended further by merging the two stages into a joint model of gene regulation across many genes simultaneously, moving beyond purely gene-by-gene modeling. The archetypal decomposition of scATAC-seq already provides a natural foundation for this direction.

Thus, ArchVelo illustrates a productive new direction for RNA velocity modeling, combining improved accuracy with interpretability. Beyond its immediate applications, the framework provides a foundation that can be extended in multiple ways, opening opportunities for future methodological advances in single-cell multi-omics.

## Methods

### Archetypal analysis on ATAC modality

We reduce dimensionality of the ATAC modality of each multi-omic dataset by applying archetypal analysis (AA).<sup>19,20</sup> For an experiment with  $N$  cells and  $S$  peak summits, let us denote the  $S \times N$  Pearson residual normalized matrix of peak or peak summit accessibility profiles as  $Y$ . For a fixed dimensionality parameter  $K$ , AA solves the following optimization problem:

$$\min_{Z,B} \|Y - ZA\|_F^2, A = BY$$

where  $B$  and  $Z$  satisfy additional constraints  $B \geq 0, Z \geq 0, B1 = 1, Z1 = 1$ . The matrix  $A \in \mathbb{R}^{K \times N}$  contains the archetypal chromatin accessibility profiles in the rows. The condition  $Z1 = 1$  on the  $S \times K$  matrix of non-negative peak summit loadings ensures that the accessibility profile of every summit, i.e., every row of the ATAC matrix  $Y$ , is approximated by a convex combination of archetypes. Moreover, the conditions on the  $K \times S$  matrix  $B$  ensure that every archetype lies on the boundary of the convex hull of the accessibility profiles at peak summits. AA can be viewed as a soft clustering technique without the orthogonality constraint. We solve the optimization problem using the efficient implementation of the AA algorithm<sup>57</sup> available via `py_pcha` Python package. The implementation was developed for large scale AA, and is based on projected gradient as well as a robust FurthestSum initialization algorithm.<sup>57</sup> The implementation also includes a relaxation of the AA problem where the archetypes can reside outside of the convex hull of the observations which the authors call the AA- $\delta$  problem. In our notation, the problem modification relaxes the  $B1 = 1$  condition into  $\max|B1 - 1| < \delta$ . By default, we use  $\delta = 0.1$ . For more details on archetypal analysis and the AA- $\delta$  algorithm we refer the reader to available literature.<sup>19,20,57</sup> We used elbow method for choosing the optimal number of archetypes for each dataset used in this study (**Supplementary Fig. 6**).

### Linear model to predict RNA from ATAC modality

We train a model to predict the log-transformed library-size normalized RNA count matrix  $X$  from the ATAC modality  $Y$ . We choose to use the archetypes as features in a ridge-regularized linear regression model:

$$\min_{w^g} \sum_{c=1}^N \left| x_c^g - \sum_{k=1}^K w_k^g \sum_{s \in P(g)} z_{sk} a_{kc} \right|^2 + \lambda^g \sum_{k=1}^K |w_k^g|^2.$$

Here  $c$  and  $k$  are the cell and archetype indices,  $g$  denotes a gene and its accessibility score for archetype  $k$  is given by the sum of the loadings  $z$  over the set of all peak summits  $P(g)$  mapping

to  $g$ . We train the model for every gene independently. We assume a linear relationship between the number of gene transcripts and archetypal chromatin accessibility of the gene, with the desired coefficients  $w_k^g$  describing the overall input of every archetype into the pool of transcripts.

For training and evaluation, we split every dataset into two subsets, using  $N_{\text{train}}$  cells to form the matrix  $Y_{\text{train}}$  for training and  $N_{\text{test}}$  cells to form the matrix  $Y_{\text{test}}$  for testing. By default, we use 66.7% and 33.3% of the cells for these purposes, respectively. To train and test our linear model, we first apply AA- $\delta$  on  $Y_{\text{train}}$  to extract the archetypal chromatin accessibility profiles on the training set:

$$Y_{\text{train}} \approx Z_{\text{train}} A_{\text{train}} = Z_{\text{train}} B_{\text{train}} Y_{\text{train}}$$

where  $Z_{\text{train}} \in \mathbb{R}_{\geq 0}^{S \times K}$ ,  $B_{\text{train}} \in \mathbb{R}_{\geq 0}^{K \times S}$ . Then, to validate the model on the held-out test set, we use the training archetypes to perform prediction on the test set. More precisely, we use the  $A_{\text{test}} = B_{\text{train}} Y_{\text{test}}$  as the archetypal features on the test set and keep  $Z_{\text{train}}$  as the corresponding loading matrix.

To optimize the regularization hyperparameter  $\lambda^g$  in the ridge regression, we apply a 3-fold cross-validation (CV) procedure by splitting  $Y_{\text{train}}$  further into 3 independent subsets and using two of them for training and one for validation for every fold. We use Pearson correlation as the scoring metric, and pick the best regularization parameter from 50 values evenly spaced between  $-7$  and  $5$  on a  $\log_{10}$  scale. After selecting the hyperparameter with best average performance, we use it to retrain the model on  $Y_{\text{train}}$  and test the final performance of the model on  $Y_{\text{test}}$  as described above.

We observed that the regression results are robust with respect to the choice of the number of archetypes (**Supplementary Fig. 7**).

## ArchVelo model

We build an ordinary differential equation (ODE) model of chromatin dynamics, transcription, splicing and degradation cascade for every gene. We propose that gene transcription can be decomposed into a sum of components driven by chromatin dynamics of accessibility archetypes  $a_k(t)$ , and assume individual kinetics at each archetype. Following MultiVelo,<sup>9</sup> we model dynamics of chromatin opening and closing using a simple two-state model. Normalizing each chromatin accessibility archetype between 0 and 1, we assume that it asymptotically approaches 1 in the open state and asymptotically decays to 0 in the closed state:  $\dot{a}_k(t) = \alpha_k^{(co)} - \alpha_k^{(co)} a_k(t)$  for some rate  $\alpha_k^{(co)}$  if chromatin is opening and  $\dot{a}_k(t) = -\alpha_k^{(cc)} a_k(t)$  for some rate  $\alpha_k^{(cc)}$  if it is closing. This modeling choice was motivated in MultiVelo by empirical observation and by an assumption that initially rapid chromatin remodeling gets slowed down by biochemical constraints such as the structures of histone complexes and their inter-molecular interactions.<sup>9</sup> With the simplifying assumption that

the chromatin opening and closing kinetics are mirror images of each other, i.e. opening and closing have the same rates, we can join these equations into one:  $\dot{a}_k(t) = \tau_k^{(c)}(t)\alpha_k^{(c)} - \alpha_k^{(c)}a_k(t)$  where

$$\tau_k^{(c)}(t) = \begin{cases} 1 & \text{if } t \leq t_k^{(c)}, a_k(t) \text{ is opening,} \\ 0 & \text{if } t > t_k^{(c)}, a_k(t) \text{ is closing,} \end{cases} \quad (1)$$

and we assume individual chromatin switch times  $t_k^{(c)}$  for every archetype. In what follows we will assume  $\alpha_k^{(co)} = \alpha_k^{(cc)}$  just for simplicity of notation while the ArchVelo model fits these rates separately.

We aim to solve the following system of ODEs:

$$\dot{a}_k(t) = \tau_k^{(c)}(t)\alpha_k^{(c)} - \alpha_k^{(c)}a_k(t), \quad k = 1 \dots K; \quad (2)$$

$$\dot{u}^g(t) = \sum_{k=1}^K \tau_k^g(t)\alpha_k^g\delta_k^g a_k(t) - \beta^g u^g(t); \quad (3)$$

$$\dot{s}^g(t) = \beta^g u^g(t) - \gamma^g s^g(t). \quad (4)$$

Here eq. (2) describes chromatin dynamics while eqs. (3) and (4) describe dynamics of unspliced and spliced counts, respectively. We assume that the transcription can be active and repressed at every archetype, with  $\tau_k^g(t)$  a piecewise constant function analogous to eq. (1) switching from 0 to 1 at transcription initialization time  $t_k^{g,(i)}$  and back to 0 at transcription repression time  $t_k^{g,(r)}$ . We assume that every archetype has its own transcription kinetic parameters, with every archetype  $a_k$ , when active, contributing to transcription with a gene-dependent rate  $\alpha_k^g$ . Simultaneously, for every gene  $g$ , we set transcription rate at an archetype to be proportional to its accessibility score for  $g$ , i.e. its total loading for the peak summits mapping to  $g$ ,  $z_k^g = \sum_{s \in P(g)} z_{sk}$ . Finally, we correct for prior normalization of the chromatin archetypes here by renormalizing them back by  $n_k = \max(a_k) - \min(a_k)$ . This results in  $\delta_k^g = z_k^g n_k$  term in eq. (3). Variables  $\beta^g$  and  $\gamma^g$  denote splicing and degradation rates, respectively.

Due to the linear nature of the above model,  $u^g$  and  $s^g$  can be decomposed into a linear combination of unspliced and spliced profiles resulting from transcription at every archetype:  $u^g = \sum_k u_k^g(t)$ ,  $s^g = \sum_k s_k^g(t)$  where  $(a_k, u_k^g, s_k^g)$  solve a system of ODEs

$$\begin{aligned} \dot{a}_k(t) &= \tau_k^{(c)}(t)\alpha_k^{(c)} - \alpha_k^{(c)}a_k(t); \\ \dot{u}_k^g(t) &= \tau_k^g(t)\alpha_k^g\delta_k^g a_k(t) - \beta^g u_k^g(t); \\ \dot{s}_k^g(t) &= \beta^g u_k^g(t) - \gamma^g s_k^g(t). \end{aligned} \quad (5)$$

For every  $k$ , this system mimics a MultiVelo model.<sup>9</sup> Each system can be solved analytically, with smooth solutions for every continuous interval of time characterized by the same chromatin

and transcriptional state. However, we emphasize that the solution to the ArchVelo problem cannot be obtained by the repeated application of MultiVelo at every archetype. The optimization is performed jointly for all  $k$  due to the presence of shared parameters  $\beta^g$ ,  $\gamma^g$  and jointly defined likelihood. We iteratively update the model parameters and latent time assignments in an expectation-maximization procedure (EM) which is adapted to the ArchVelo model likelihood (eq. (6) below).

We observed that ArchVelo results are robust with respect to the choice of the number of archetypes (**Supplementary Fig. 8**).

### Parameter initialization for ArchVelo

Due to an increase in the number of parameters in ArchVelo compared to MultiVelo, we sought to find a suitable initialization for ArchVelo parameters. To this end, on the same dataset, we first apply a modification of MultiVelo that we call MultiVelo-AA. MultiVelo-AA differs from MultiVelo in how the chromatin accessibility profiles for individual genes are defined. In MultiVelo, a single chromatin accessibility profile  $c^g$  for every gene is defined using the sum of accessibility profiles at its promoter and other linked peaks.<sup>9</sup> In MultiVelo-AA, the chromatin profiles for every gene are instead predicted using the results of AA on the Pearson normalized ATAC matrix (see **Archetypal analysis on ATAC modality**). More precisely, for every gene  $g$ , in MultiVelo-AA, we use  $\bar{c}^g = \sum_k z_k^g a_k$  in lieu of the aggregated chromatin accessibility profile  $c^g$  and MultiVelo is run on the  $(\bar{c}^g, u^g, s^g)$  triple. We use the output of MultiVelo-AA to initialize latent time assignments, and kinetic parameters for transcription, splicing and degradation for every gene. More precisely, while the chromatin dynamics parameters are fit separately for every archetype, for other kinetic parameters, the same initialization is used for all components. Note that, analogously to MultiVelo, in both ArchVelo and MultiVelo-AA, weighted nearest neighbor (WNN) smoothing is first applied to the archetypes  $a_k$ , yielding a set of smoothed archetypes. These smoothed archetypes are used in ArchVelo to fit eq. (2) and in MultiVelo-AA to define  $\bar{c}^g$ ; we will nevertheless not include any smoothing notation and will keep denoting the archetypes by  $a_k$  in the description of ArchVelo.

### Archetypal velocity components

Velocity profiles  $v^g = \dot{s}^g$  for individual genes are obtained from the solution of eq. (4). The velocity graph (transition matrix for the cells), common latent time assignment and velocity embeddings are produced using the corresponding methods in scVelo package.<sup>7</sup> By default, as in MultiVelo, to bring the velocity profiles for individual genes to the same scale, we apply  $l_1$  normalization to each of them. By default, likelihood cutoff 0.05 is used in ArchVelo.

Archetypal decomposition of chromatin dynamics automatically provides downstream linear decomposition for the velocity profiles  $v^g = \dot{s}^g = \sum_{k=1}^K \dot{s}_k^g = \sum_{k=1}^K v_k^g$  (eq. (5)). As a result,

cell trajectories can be decomposed into components potentially driven by different regulatory programs. Moreover, the individual latent time estimates for every component can provide a more accurate ordering of cells in case of conflicting signal or presence of cell cycle component. The normalized velocity  $\tilde{v}^g$  can then be decomposed into a sum of normalized velocity components  $\tilde{v}_k^g$ :  $\tilde{v}^g = \frac{v^g}{\|v^g\|_{l_1}} = \sum_{k=1}^K \frac{v_k^g}{\|v^g\|_{l_1}} = \sum_{k=1}^K \tilde{v}_k^g$ . Note that the universal normalization with the  $l_1$  norm of the full velocity is applied to each component, rather than individual  $l_1$  normalization. The velocity graph, latent time assignments and velocity embedding can then be calculated for each  $k$  to infer archetypal velocity fields. In MultiVelo and ArchVelo, velocity profiles are smoothed among nearest neighbors using the connectivity matrix of the RNA modality; however, the unsmoothed velocities are stored in ArchVelo and can be used.

Furthermore, archetypal velocity components and spliced count decomposition can aid in the downstream cell and gene selection. Specifically, we define cell-specific archetypal velocity loadings as  $\|v_k(c)\|_2 = \sqrt{\sum_g |v_k^g(c)|^2}$  (see **Supplementary Fig. 3A,C** for an example). Such loadings can be defined for both velocity, unspliced and spliced count components and potentially highlight cells sharing a regulatory program. By default, in ArchVelo,  $k$ -means clustering on the spliced count component ( $s_k^g$ ) loadings is used to identify cells with significant archetypal transcriptional input. Analogously, we leverage the archetypal decomposition to cluster genes using their spliced count components. Specifically, genes are grouped according to the spliced count component with the highest  $l_1$ -norm (see **Supplementary Fig. 3B,D** for an example). For each archetype, the latent time can then be further refined on the cells selected for this archetype by subsetting the matrix to the selected genes and rerunning the velocity analysis using the corresponding component of velocity.

For archetypal velocity components, we also developed a variation of the velocity embedding visualization on a grid. In scVelo and MultiVelo, the velocity graph is built using the cosine similarity between the velocity vector and the vector of change in gene expression between cells. This approach is agnostic to the cell-specific norm of the velocity vector and is not effective for visualizing the velocity components which are usually concentrated on a subset of the manifold. To alleviate this problem, we developed a modification of the scVelo grid embedding technique which scales the velocity components back by  $\|v_k(c)\|_2$  estimates on the grid. This method was applied to visualize all archetypal velocity components in this paper.

### ArchVelo model likelihood

We assume Gaussian noise in the observations which come in the form of spliced, unspliced counts and archetypal chromatin accessibility profiles from AA on the ATAC modality  $Y$ . The archetypes are weighed with their corresponding weights  $\delta_k^g$  (see **ArchVelo model**). Additionally, we introduce a parameter  $w_c$  scaling the influence of  $Y$  on the velocity inference. Denoting

M. Avdeeva et al.

$\tilde{a}_g^k = w_c \delta_k^g a_g^k$ , for cell  $i$ , the observation becomes  $\mathbf{x}_i = (u_i, s_i, \tilde{a}_{1,i} \dots \tilde{a}_{K,i})$ . For the prediction of the form  $\mathbf{f}(t_i, \theta) = (\hat{u}_i, \hat{s}_i, \hat{a}_{1,i} \dots \hat{a}_{K,i})$ , we optimize the negative log-likelihood:

$$\begin{aligned} -\log \mathcal{L}(\theta) &= \frac{K+2}{2} \log(2\pi\sigma^2) + \frac{1}{2n\sigma^2} \sum_{i=1} \|\mathbf{x}_i - \mathbf{f}(t_i, \theta)\|^2 \\ &= \frac{K+2}{2} \log(2\pi\sigma^2) + \frac{1}{2n\sigma^2} \sum_{i=1} \left[ (u_i - \hat{u}_i)^2 + (s_i - \hat{s}_i)^2 + w_c^2 \sum_{k=1}^K (\delta_k^g (a_{k,i} - \hat{a}_{k,i}))^2 \right] \end{aligned} \quad (6)$$

where  $\sigma^2$  parametrizes the variance. Weight  $w_c$  is also available as a parameter in MultiVelo, with  $w_c = 0.6$  as default. For benchmarking experiments,  $w_c$  was varied between 0 and 1, with 11 different parameters at increments of 0.1. As the result of benchmarking, for ArchVelo,  $w_c = 0.3$  was chosen as the default parameter.

### ArchVelo expectation-maximization details

We iteratively update the model parameters and latent time assignments in an expectation-maximization procedure (EM) adapted to the model likelihood. More precisely, during the maximization step, the model parameters are updated via maximum likelihood. For each  $k = 1 \dots K$ , the model includes  $2 + 6K$  parameters:  $3K$  parameters describing chromatin dynamics  $(\alpha_k^{(co)}, \alpha_k^{(cc)}, t_k^{(c)})$ ,  $3K$  parameters for the transcription component  $(\alpha_k^g, t_k^{(i)}, t_k^{(r)})$ , and 2 parameters for splicing and degradation. We first update the parameters fitting eq. (2) to the archetypes using current latent time estimates (via the dual annealing method implemented in scipy Python package). We then optimize the component of the likelihood corresponding to eqs. (3, 4) to update the parameters corresponding to transcription, splicing and degradation. We use the Nelder-Mead method implementation from scipy Python package to optimize the parameters.

The E-step estimates the latent time assignment of every cell given the current updated ODE parameters. We also optionally adopt the strategy described in MultiVelo for the E-step. There the uniformly spaced anchor points are maintained, and cells are assigned with the latent time of the nearest anchor point. 500 anchor points are similarly maintained by default. The key difference here is the nearest neighbor graph which utilizes the distances in the  $(K+2)$ -dimensional phase space of unspliced, spliced counts and  $K$  archetypes rather than the 3-dimensional space of  $(c, u, s)$  in MultiVelo (see **ArchVelo model likelihood**). Nearest neighbor assignment is obtained via the kd-tree quick nearest neighbor lookup implemented in scipy Python package.

### Benchmarking metrics and details

To benchmark ArchVelo against other methods, we used two metrics: latent time silhouette and mean log-likelihood in the  $u$ - $s$  phase space. With the underlying model assuming a unique common latent time profile, an appropriate metric of robustness should measure the degree of alignment of

latent time estimates. To achieve this, we correlate gene-specific inferred latent time profiles using Spearman correlation. We have found that, typically, individual latent time assignments are correlated between genes within several gene modules and anticorrelated between those modules. To assess latent time robustness, we apply hierarchical clustering to the Spearman correlation profiles (average method, Euclidean metric) and cut the tree to form  $n = 2, \dots, 20$  clusters. For every  $n$ , we calculate silhouette score (`silhouette_score` method implemented in `sklearn` Python package) of the resulting clustering. Mean silhouette score over  $n$  is reported as latent time silhouette. We use the common  $(u, s)$  component of the log-likelihood to assess the quality of the fit. To obtain a fair likelihood comparison, the same subset of cells was used for all methods for every gene. For this, we used the cell subset chosen by the outlier removal procedure in MultiVelo. For this benchmark, we included ArchVelo, MultiVelo, scVelo and VeloVI, as they report both individual latent times and  $u$ - $s$  fits for genes. Finally, we note that in these benchmarking experiments, for every method, we used default parameters (except  $w_c$  which was varied for multi-omic methods) and the same set of cells and genes and, as a result, the same kNN graph.

To evaluate cross-boundary direction correctness, we used `cross_boundary_correctness` method from VeloAE Python package, version 0.2.0. For the mouse embryonic brain, we constructed the ground truth graph of cell type transitions for the neuronal and glial compartments. The graph contains an edge from astrocyte progenitors to astrocytes, and two branches originating from IPCs. IPCs have been shown to be capable to produce both upper-layer (UL) and deeper-layer (DL) neurons.<sup>26,27</sup> In the UL branch, we placed MN1 between IPC and UL due to its expression of upper-layer markers such as *Cux1*, *Cux2*, *Plxna4* and present but weaker expression of *Satb2*. For the deeper layer, we used *Bcl11b* (*Ctip2*) and *Fezf2*, well-known markers of DL, for establishing cluster identities. Expression of these markers placed MN2 cluster into the DL branch. Strong expression of markers associated with migration and intermediate levels of these and other well-known markers of DL identity led us to place MN2 between IPC and DL on this branch. Our tree is also consistent with previous literature.<sup>58</sup> For the HSC dataset, we constructed the ground truth graph of cell type transitions for all clusters. Here we refer the reader to several reviews on hematopoiesis.<sup>33,59,60</sup> Previous publications suggested an early split between the neutrophil-monocyte-lymphoid lineage and the BEM-erythroid-megakaryocyte progenitors.<sup>30,32-34</sup> The common progenitor to the latter lineage was termed erythro-myeloid progenitor (EMP),<sup>32</sup> the notation that we adopt. Notably, our ground truth model of hematopoiesis is also similar to the model reconstructed by Zheng et al.<sup>35</sup> in another single-cell study. Finally, we note that while MEPs are capable of giving rise to bona fide basophils in basophilic conditions,<sup>30</sup> we included an additional MEP-to-BEM edge to refine the previously available model.

Published methods were applied for benchmarking purposes according to their documentation. In particular, due to similar RNA preprocessing pipelines, scVelo, MultiVelo, DeepVelo and VeloVI

were applied on our preprocessed datasets, with default parameters. ATAC modality for MultiVelo was preprocessed in accordance with the original pipeline. Cell2Fate was applied on raw count data on the same set of cells as ArchVelo for both datasets. Parameters `min_shared_counts = 10`, `n_var_genes = 1000` were used for the mouse embryonic brain and `min_shared_counts = 100`, `n_var_genes = 1500` on HSC, similar to our preprocessing pipeline. Similarly, TFVelo was applied on the raw counts for the same set of cells as ArchVelo, with preprocessing steps and parameters adapted to most closely correspond to our preprocessing pipeline.

### ArchVelo transition scores

To calculate ArchVelo transition probabilities, cell by cell matrix of transition probabilities was extracted applying `transition_matrix` method from `scvelo.tools` on the ArchVelo velocity field. For this analysis, to fix the number of neighbors used for every cell for normalization purposes, the velocity field was recalculated with parameter `mode_neighbors = 'connectivities'`. By default, we use 10 nearest neighbors for this analysis. Transition probabilities were then averaged over every pair of clusters (to normalize for cluster sizes) and normalized to condition on the source cluster.

### Single-cell ATAC-seq data preprocessing

All single-cell analysis was done using `scanpy`.<sup>61</sup> By default, for all ArchVelo analyses, the peak summit counts were filtered to the ones present in at least 1% of cells. For every peak summit count matrix  $Y$ , we applied Poisson correction  $\lceil Y/2 \rceil$  on the raw counts.<sup>62</sup> Our default normalization step before archetypal analysis was to apply Pearson residual approach to the resulting matrix, with  $\theta = 1$  as the default overdispersion parameter. The residuals were clipped at  $\sqrt{N}$  where  $N$  is the number of cells. We used an implementation of this normalization technique from Python `scanpy` package.

### Single-cell RNA-seq data preprocessing

`Scanpy` was used for all analysis of scRNA-seq data. Typical QC steps were followed to remove cells based on low or extremely high signal as well as high mitochondrial counts. For initial analysis, normalization was performed (Pearson residuals or log-transformed library size normalized counts, depending on dataset, see below) and then standard scRNA-seq processing steps were followed: genes with low counts were filtered, highly variable genes were selected, PCA was run, kNN graph was built, and clustering performed. By default, standard imputation step was run for unspliced and spliced counts on each dataset with the same parameters that were used for kNN graph construction. Cell types were determined based on marker gene expression. For differential expression analysis, `rank_genes_groups` from `scanpy` with `method = 'wilcoxon'` was used. Full details of RNA-seq analysis run on each dataset can be found below.

## **BABEL gene expression predictions**

BABEL code was downloaded from github and then edited to enable the assignment of train-test splits to compare directly using the same training and test sets as our archetypal regression. Each dataset was then run using the `train_model.py` function with mouse v25 gencode for mouse datasets or the HG38 version supplied in the BABEL codebase. Anndata files were read in and modified to be in the expected format for BABEL (i.e. concatenating the ATAC and RNA modalities into one anndata object). The RNA predictions from ATAC were extracted from the resulting `atac_rna_adata.h5ad` anndata. Pearson correlations were calculated with the true gene expression values using `np.log1p` on the predicted  $X$  matrix and log library size normalization on the real gene expression values.

## **ArchR gene scores**

Gene scores were generated using the `addGeneScoreMatrix` with all default parameters. The gene score matrix was extracted using the `'getMatrixFromProject'` function, with `useMatrix = "GeneScoreMatrix"`.

## **Peak and peak summit calling**

MACS2 was run on single base pair Tn5 insertion sites for all datasets to generate the set of peaks. Additional 150bp-wide summit regions of maximum signal within these peaks were then generated, using bigwig files of signal in the peak regions, using the `biosignals` package with a custom script. Counts over these 150 bp regions (referred to as peak summits throughout the text) were generated (typically using ArchR, dataset-specific details can be found below) and used for all subsequent analyses.

## **Nearest genes mapping**

Gene annotations were used to assign peak summit regions to nearest genes, using the `"bedtools closest"` function. For mouse datasets, the GTF annotations file from Gencode V25 was used for mouse mm10 genome, downloaded from [https://ftp.ebi.ac.uk/pub/databases/gencode/Gencode\\_mouse/release\\_M25/](https://ftp.ebi.ac.uk/pub/databases/gencode/Gencode_mouse/release_M25/). For human, HG38 v46 annotation was downloaded from Gencode at [https://ftp.ebi.ac.uk/pub/databases/gencode/Gencode\\_human/release\\_46](https://ftp.ebi.ac.uk/pub/databases/gencode/Gencode_human/release_46).

Peak summits were sequentially labeled as promoter (if within 2kb), exon (if directly within an exon), intronic (if distance to gene was 0), or intergenic (if within 50kb). The rest of the summits were left unlabeled.

## TF motif analysis

TF motif binding weights for peak summits were generated using motifmatchr<sup>63</sup> with a p-value threshold of  $10^{-3}$ , over the 693 clustered motif PWMs from Vierstra et al.<sup>64</sup> For the CD8 T cell data analysis, motif binding weights were generated using PWMs selected for T cells from Pritykin et al.<sup>42</sup> Thus, peak summit by motif matrix of motif binding weights  $M$  was generated for each analysis.

Cell-specific motif scores were defined as the columns of  $Y^T M$ , where  $Y$  is the peak summit by cell matrix of ATAC counts normalized using Pearson residual normalization. These scores are shown on UMAPs in **Fig. 6C-E**. Then motif correlation scores  $c_{m,a}$  for every motif  $m$  and archetype  $a$  were defined as the Pearson correlation between  $a$  and the vector of cell-specific scores corresponding to  $m$ . In **Fig. 8F**, correlation scores were calculated for the progenitor compartment only. We further argued that for full datasets the archetypal analysis results can be also used to calculate the motif scores more directly as  $Z^T M$ , where  $Z$  is the peak summit by archetype matrix of archetypal loadings. The resulting matrix  $Z^T M$  was column-normalized. To ensure functional relevance of such associations between motives and archetypes, in **Fig. 6C-E**, the motifs were first filtered by their archetypal correlation score ( $c_{m,a} > 0.1$ ) and then the archetypal motif scores were shown using barplots.

## PBMC data analysis

Data was downloaded from the 10x website.<sup>65</sup>

**ATAC modality.** The data preprocessing and analysis was done as previously described.<sup>21</sup> An initial analysis was performed to determine a final set of cells to use, with ArchR. ArrowFiles were generated using the ‘createArrowFiles’ function, setting minTSS=4, minFrag = 1000, addTileMat=TRUE, and addGeneScoreMat = True. Doublets were inferred using the ‘addDoubletScores’ function, using k=10, knnMethod = “UMAP”, and LSIMethod = 1. An ArchR project was then created using these ArrowFiles, and 1,249 doublets were filtered using the ‘filterDoublets’ function, with a filterRatio set to 1. These cells were then used as the final cells for subsequent analysis. A bed file of single base pair insertion sites was generated from the fragment file, only incorporating this final list of cells. Peaks were called using this bed file with MACS2 callpeak, using shift set to -75 and extsize set to 150, qvalue to 0.05, with nomodel, bdg to True, SPMR to True, and keep-dup to all. 150bp-wide summit regions of maximum signal within the peaks were generated using the biosignals package with a custom script. To generate the peak summit read count matrix, a new set of ArrowFiles was generated, using minTSS=0. minFrag=1, maxFrag=1e20, addTileMat = TRUE, and addGeneScoreMat=TRUE. An ArchRProject was then generated on these ArrowFiles, and then filtered to the cell list generated previously, using the subsetCells function. Next, addIterativeLSI, addClusters, addUMAP, and addImputeWeights

were run using default parameters. The peak summit read count matrix was generated using `addPeakSet` and `addPeakMatrix` (run with a ceiling of 20). The peak summit read count matrix was retrieved using the `getMatrixFromProject` function with `useMatrix = "PeakMatrix"`. This was then read into `scanpy`, and filtered to peaks present in at least 1% of cells.

**RNA modality.** The filtered feature matrix was downloaded from the 10x website and loaded into `scanpy`. Genes present in fewer than 1% of cells were filtered out (23,118 genes), and then cells with fewer than 1% of this set of genes was filtered out (23 cells). 309 cells with more than 30% mitochondrial reads were removed, 14 cells with total counts greater than 20k were removed, and 26 cells with more than 5000 genes present were removed. Pearson residual normalization was run with `theta` of 100. PCA was run with 100 PCs. The scRNA-seq matrix was then filtered to include only the cells present in the scATAC-seq modality, resulting in 9,952 cells. The scRNA-seq kNN graph was built using cosine metric, 30 nearest neighbors, and 30 PCs, and UMAP was run with all default parameters. Leiden clustering was performed with resolution of 1.2, and cell type annotations were determined based on marker gene expression. Specifically, *MS4A1*, *CD19* and *CD79A* were used to determine B cells; *CD3D*, *CD3E* and *CD4* were used to classify CD4 T cells, some of which were annotated as Treg cells based on expression of *FOXP3* and *IKZF2*; *CD3D*, *CD3E*, *CD8A* and *CD8B1* were used to identify CD8 T cells, of which some were defined as Activated CD8 T cells if also expressing *GZMK*, *PRF1*, and *KLRG1*, the rest expressed high levels of *CCR7*, *TCF7* and *SELL* and were assigned as CD8 memory/naive; a few cells expressed low levels of *CD3G*, *CD4*, *CD8A* and *CD8B1* and were assigned as T cells; NK cells were determined using *NKG7* and *NCAM1*; pDC were identified using *TCF4*, *MZB1*, *IL3RA*, *IRF4*, *SERPINF1*, *PLAC8*, *LILRB4*, *GZMB*, *JCHAIN*, *IRF7*, *SLC15A4*, *SPIB*, *ITM2C* and *IRF8*; CD14+ monocytes were identified using *LYZ*, *CD14*, *VCAN*, *CD36* and *SERPINA1*; and CD16+ monocytes were identified using *TCF7L2*, *MTSS1*, *SERPINA1*, *RHOC*, *MS4A7*, *IFITM2*, *FCGR3A*, *SIGLEC10*, *CX3CR1* and *LILRB1*.

### Mouse brain data analysis

Data was downloaded from the 10x website.<sup>25</sup>

**ATAC modality.** The data preprocessing and analysis was done as previously described.<sup>21</sup> Single base pair insertion sites were generated from the fragment file. Peaks were called for each cell type separately based on the previous cell type annotations from the paper,<sup>9</sup> using `MACS2`, with `shift` set to `-75` and `extsize` set to `150`, `qvalue` to `0.05`, `bdg` to `True`, `SPMR` to `True`, and `keep-dup` to `all`. These peaks were then merged using the R package `GenomicRanges` `reduce` function. A bigwig file was generated for each cell type using the `bedGraphToBigWig` tool. 150 bp summit regions of maximum signal within the peaks were generated using the `biosignals` package by adding the signal from these bigwig files together in an R script. The peak summit read count

matrix was generated using ArchR. ArrowFiles were first generated from the fragment file, using `minTSS=0`, `minFragments=1`, `maxFragments=1e20`, `addTileMat = TRUE`, and `addGeneScoreMat=TRUE`. An ArchRProject was then generated on these ArrowFiles. The peak summit matrix was generated using the `addPeakMatrix` function with a ceiling of 1000. It was then retrieved using the `'getMatrixFromProject'` function with `useMatrix = "PeakMatrix"` and `binarize` set to `FALSE`. This matrix was then filtered to the cells that were used in MultiVelo.<sup>9</sup> After standard preprocessing steps described above, archetypal analysis was applied with  $k = 10$  on this dataset.

**RNA modality.** The cells were first filtered according to the MultiVelo tutorial resulting in 3653 cells. Initially, a liberal cutoff of `min_shared_counts=10` via `filter_and_normalize` method from `scVelo` package was applied. Log-transformed library normalized counts were used for this dataset. 30 PCs and 50 nearest neighbors were used to construct the kNN graph for this dataset and Leiden clustering with `resolution=4` was applied to identify subpopulations relevant to our analysis. Two small subpopulations were filtered out: one based on expression of vascular-associated ECM and adhesion markers including *Col4a1/2*, *Lama4*, *Itgb1*, and one based on expression of GABAergic neuron markers such as *Gad2*, *Nrxn3*, and *Dlx1/2*. The resulting matrix was filtered again with parameters `min_shared_counts = 10` and 1000 top highly variable genes were selected. The final matrix after cells and genes were filtered for high-quality ATAC signal contained 3252 cells and 930 genes. The kNN graph for this matrix was recalculated with the same parameters. Cluster annotation was performed according to the available literature. Specifically, for both neuronal and glial compartments, we used the marker gene sets from a mouse embryonic brain single-cell transcriptomic atlas,<sup>58</sup> Extended Data Figures 2a, 4d and 5a,c. The main markers used for annotation were *Satb2*, *Cux1*, *Cux2*, *Plxna4* for UL, *Fezf2*, *Bcl11b*, *Tle4*, *Sox5*, *Ldb2* for DL, *Nrp1* as a marker of migration, *Pax6*, *Hes5*, *Fabp7*, *Dbi*, *Slc1a3* for RG, *Eomes*, *Neurog2*, *Btg2* for IPC, *Aldh1l1*, *Sparcl1* for astrocytes, *Olig1*, *Olig2*, *Pdgfra* for OPCs, and *Neurod2*, *Neurod6*, *Tubb3* as general neuronal markers. We also supplemented these gene sets with additional astrocytic markers *Slc6a11*, *Tnc* and *Slco1c1*,<sup>66,67</sup> and migration/axon guidance markers *Dcc*, *Sema3c* and *Unc5d*.<sup>68</sup>

## HSC data analysis

Data for day 7 (MV2) was downloaded from GEO using accession number GSE209878.

**ATAC modality.** Single-base pair insertion sites were generated from the fragment file over cells from the paper. The same pipeline as for the mouse brain data was run, with the same parameters. The final peak summit matrix from ArchR was filtered so that all cells contained  $> 3000$  total peak summit count. Peak summits were filtered according to the standard preprocessing pipeline described above. Archetypal analysis was applied with  $k = 9$  to this dataset.

**RNA modality.** 14895 cells in the original dataset were first filtered according to standard QC pipelines and filtered for doublets using `scrublet` algorithm<sup>69</sup> available in Python through `scanpy` re-

sulting in 12289 cells. Genes were filtered with parameters  $\text{min\_cells} = 20$ ,  $\text{min\_shared\_counts} = 100$ . To annotate this dataset, due to a strong cell cycle signal, cell cycle regression was performed. We applied the standard pipeline for cell cycle regression, with a slight modification allowing to add the intercept back after regression. The cell cycle regression was performed on the log-transformed library size normalized counts. 1500 highly variable genes were then selected, and 50 PCs and 50 nearest neighbors were used for this dataset. Data was clustered with  $\text{resolution} = 3$  and re-annotated based on available literature.<sup>29–35</sup> Specifically, we used the following subpopulation markers: HSC/MPP: *HLF*, *AVP*, *HOXA9*, LMPP: *FLT3*, *TCF4*, Prog DC: *LYZ*, *IRF8*, *JAML*, GMP: *MPO*, Prog N: *ELANE*, *AZU1*, *PRTN3*, EMP: *GATA2*, *TFRC*, *CPA3*, *RHEX*, MEP: *TFRC*, *TFR2*, Ery: *KLF1*, *HBB*, Prog MK: *ITGA2B*, Platelet: *PLEK*, *PF4*. The final matrix after cells and genes were filtered for high-quality ATAC signal contained 10935 cells and 1325 genes. The kNN graph was recalculated for this matrix using 50 PCs and 50 nearest neighbors.

### CD8 T cell data analysis

The CD8 T cell dataset analyzed here was generated as part of a larger study characterizing splenic T cells in LCMV infection.<sup>21</sup> We focused on the subset of cells previously annotated as effector CD8 T cell, progenitor CD8 T cell and exhausted CD8 T cell subpopulations, i.e. all CD8 T cells activated by the LCMV antigen, from both the LCMV Armstrong (acute infection) and LCMV clone 13 (chronic infection) conditions.

**ATAC modality.** Archetypal analysis was applied on ATAC data from both infection conditions simultaneously. Our standard preprocessing pipeline was applied to the ATAC matrix, followed by archetypal analysis with  $k = 9$ .

**RNA modality.** The same pipeline was applied to both clones. Genes were first filtered with `filter_and_normalize()` method from `scVelo` package with  $\text{min\_shared\_counts} = 30$ ,  $\text{n\_top\_genes} = 1500$  parameters. Pearson residual normalized counts from the full dataset were then used.<sup>21</sup> 50 PCs, 50 nearest neighbors and cosine metric were used for kNN graph generation (and unspliced and spliced count imputation). For MultiVelo, standard preprocessing was applied, with cells filtered to having  $> 2000$  of aggregated ATAC counts. The final matrix for velocity analyses (after cells and genes were filtered for high-quality ATAC signal) contained 3707 cells and 1410 genes for the Armstrong clone dataset and 2520 cells and 1439 genes for clone 13. To annotate this dataset, cells from both clones (6409 effector, progenitor and exhausted CD8 T cells prior to cell filtering) were jointly processed. We performed this analysis on the union of genes for both datasets. For finer clustering, we used 50 PCs, 30 nearest neighbors for the kNN graph and performed Leiden clustering with resolution 2. This yielded 23 clusters that were initially grouped into 10 functional and differentiation cell states, reproducible between the two infection conditions, as described in the main text and **Fig. 7B**. To further split the progenitor subpopulation into  $\text{Ccr6}^-$  and  $\text{Ccr6}^+$

M. Avdeeva et al.

subsets, we applied Leiden clustering to this subpopulation with resolution 0.4. For cell ordering over latent time in the progenitor compartment, we used the difference between archetype 5 and archetype 8. To calculate scores for genes that were previously reported to be differentially expressed in the CD62L<sup>+</sup> and CD62L<sup>-</sup> precursor exhausted T cell clusters,<sup>46</sup> we used `score_genes()` method from `scanpy`.

ARTICLE IN PRESS

## Data availability

The scATAC+RNA-seq data for PBMC, mouse brain and human hematopoiesis used in this study is publicly available<sup>9,25,65</sup> at 10x Genomics website and NCBI GEO using accession code GSE209878 at <https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE209878>.

The scATAC+RNA-seq data for CD8 T cells used in the analysis presented here was generated as part of a larger study characterizing splenic T cells in LCMV infection.<sup>21</sup> The data have been deposited to NCBI GEO under the accession code GSE300984, available at

<https://www.ncbi.nlm.nih.gov/geo/query/acc.cgi?acc=GSE300984>. Raw and preprocessed data required for full ArchVelo analysis and benchmarking have been deposited to Figshare and are available at <https://doi.org/10.6084/m9.figshare.27931914>,

<https://doi.org/10.6084/m9.figshare.30114808>,

<https://doi.org/10.6084/m9.figshare.31843234>, and

<https://doi.org/10.6084/m9.figshare.31909822>.

## Code availability

The open-source ArchVelo package is publicly available at <https://github.com/pritykinlab/ArchVelo> under BSD-3-Clause license. The computational code used in the analysis presented in this manuscript is publicly available at

[https://github.com/pritykinlab/ArchVelo\\_notebooks](https://github.com/pritykinlab/ArchVelo_notebooks) under BSD-3-Clause license. The specific version of the code associated with this publication is archived in Zenodo and is accessible via <https://doi.org/10.5281/zenodo.20085695>.<sup>70</sup>

## References

1. A. Baysoy, Z. Bai, R. Satija, and R. Fan. The technological landscape and applications of single-cell multi-omics. *Nature Reviews Molecular Cell Biology*, 24(10):695–713, 2023.
2. A. Regev, S. A. Teichmann, E. S. Lander, I. Amit, C. Benoist, E. Birney, B. Bodenmiller, P. Campbell, P. Carninci, M. Clatworthy, et al. The human cell atlas. *elife*, 6:e27041, 2017.
3. L. Heumos, A. C. Schaar, C. Lance, A. Litinetskaya, F. Drost, L. Zappia, M. D. Lücken, D. C. Strobl, J. Henao, F. Curion, et al. Best practices for single-cell analysis across modalities. *Nature Reviews Genetics*, 24(8):550–572, 2023.
4. P. V. Kharchenko. The triumphs and limitations of computational methods for scRNA-seq. *Nature methods*, 18(7):723–732, 2021.
5. D. E. Wagner and A. M. Klein. Lineage tracing meets single-cell omics: opportunities and challenges. *Nature Reviews Genetics*, 21(7):410–427, 2020.
6. G. La Manno, R. Soldatov, A. Zeisel, E. Braun, H. Hochgerner, V. Petukhov, K. Lidschreiber, M. E. Kastriiti, P. Lönerberg, A. Furlan, et al. RNA velocity of single cells. *Nature*, 560(7719):494–498, 2018.
7. V. Bergen, M. Lange, S. Peidli, F. A. Wolf, and F. J. Theis. Generalizing RNA velocity to transient cell states through dynamical modeling. *Nature biotechnology*, 38(12):1408–1414, 2020.
8. V. Bergen, R. A. Soldatov, P. V. Kharchenko, and F. J. Theis. RNA velocity—current challenges and future perspectives. *Molecular systems biology*, 17(8):e10282, 2021.
9. C. Li, M. C. Virgilio, K. L. Collins, and J. D. Welch. Multi-omic single-cell velocity models epigenome–transcriptome interactions and improves cell fate prediction. *Nature biotechnology*, 41(3):387–398, 2023.
10. A. Riba, A. Oravecz, M. Durik, S. Jiménez, V. Alunni, M. Cerciati, M. Jung, C. Keime, W. M. Keyes, and N. Molina. Cell cycle gene regulation dynamics revealed by RNA velocity and deep-learning. *Nature communications*, 13(1):2865, 2022.
11. M. Lange, V. Bergen, M. Klein, M. Setty, B. Reuter, M. Bakhti, H. Lickert, M. Ansari, J. Schniering, H. B. Schiller, et al. CellRank for directed single-cell fate mapping. *Nature methods*, 19(2):159–170, 2022.
12. A. Gayoso, P. Weiler, M. Lotfollahi, D. Klein, J. Hong, A. Streets, F. J. Theis, and N. Yosef. Deep generative modeling of transcriptional dynamics for RNA velocity analysis in single cells. *Nature methods*, 21(1):50–59, 2024.
13. H. Cui, H. Maan, M. C. Vladiou, J. Zhang, M. D. Taylor, and B. Wang. DeepVelo: deep learning extends RNA velocity to multi-lineage systems with cell-specific kinetics. *Genome biology*, 25(1):27, 2024.
14. J. Li, X. Pan, Y. Yuan, and H.-B. Shen. TFvelo: gene regulation inspired RNA velocity estimation. *Nature Communications*, 15(1):1387, 2024.
15. A. Aivazidis, F. Memi, V. Kleshchevnikov, S. Er, B. Clarke, O. Stegle, and O. A. Bayraktar. Cell2fate infers RNA velocity modules to improve cell fate prediction. *Nature methods*, 22(4):698–707, 2025.
16. G. Gorin, M. Fang, T. Chari, and L. Pachter. RNA velocity unraveled. *PLOS Computational Biology*, 18(9):e1010492, 2022.
17. J. D. Buenrostro, B. Wu, U. M. Litzenburger, D. Ruff, M. L. Gonzales, M. P. Snyder, H. Y. Chang, and W. J. Greenleaf. Single-cell chromatin accessibility reveals principles of regulatory variation. *Nature*, 523(7561):486–490, 2015.
18. S. Preissl, K. J. Gaulton, and B. Ren. Characterizing cis-regulatory elements using single-cell epigenomics. *Nature Reviews Genetics*, 24(1):21–43, 2023.
19. A. Cutler and L. Breiman. Archetypal Analysis. *Technometrics*, 36(4):338–347, 1994.
20. T. Hastie, R. Tibshirani, J. H. Friedman, and J. H. Friedman. *The elements of statistical learning: data mining, inference, and prediction*, volume 2. Springer, 2009.
21. S. K. Walker, J. van der Veecken, A. Rudensky, and Y. Pritykin. Single-cell multiomics reveals archetypal regulatory programs shared across CD4 and CD8 T cell subsets in viral infection. *bioRxiv*, 2025.09.08.675014, 2025.
22. A. Kallies, D. Zehn, and D. T. Utzschneider. Precursor exhausted T cells: key to successful immunotherapy? *Nature Reviews Immunology*, 2019. PMID: 31591533.
23. J. M. Granja, M. R. Corces, S. E. Pierce, S. T. Bagdatli, H. Choudhry, H. Y. Chang, and W. J. Greenleaf. ArchR is a scalable software package for integrative single-cell chromatin accessibility analysis. *Nature genetics*, 53(3):403–411, 2021.

24. K. E. Wu, K. E. Yost, H. Y. Chang, and J. Zou. BABEL enables cross-modality translation between multiomic profiles at single-cell resolution. *Proceedings of the National Academy of Sciences*, 118(15):e2023070118, 2021.
25. Single Cell Multiome ATAC + Gene Expression dataset analyzed using Cell Ranger ARC 1.0.0. Fresh Embryonic E18 Mouse Brain (5k). <https://www.10xgenomics.com/datasets/fresh-embryonic-e-18-mouse-brain-5-k-1-standard-1-0-0>. Accessed: 2024-10-20.
26. N. A. Vasistha, F. García-Moreno, S. Arora, A. F. P. Cheung, S. J. Arnold, E. J. Robertson, and Z. Molnár. Cortical and Clonal Contribution of Tbr2 Expressing Progenitors in the Developing Mouse Brain. *Cerebral cortex (New York, N.Y. : 1991)*, 25(10):3290–3302, oct 2015. PMID: 24927931.
27. T. Kowalczyk, A. Pontious, C. Englund, R. A. M. Daza, F. Bedogni, R. Hodge, A. Attardo, C. Bell, W. B. Huttner, and R. F. Hevner. Intermediate neuronal progenitors (basal progenitors) produce pyramidal-projection neurons for all layers of cerebral cortex. *Cerebral cortex (New York, N.Y. : 1991)*, 19(10):2439–2450, oct 2009. PMID: 19168665.
28. C. Qiao and Y. Huang. Representation learning of RNA velocity reveals robust cell transitions. *Proceedings of the National Academy of Sciences of the United States of America*, 118(49), dec 2021. PMID: 34873054.
29. A. M. Ranzoni, A. Tangherloni, I. Berest, S. G. Riva, B. Myers, P. M. Strzelecka, J. Xu, E. Panada, I. Mohorianu, J. B. Zaugg, and A. Cvejic. Integrative Single-Cell RNA-Seq and ATAC-Seq Analysis of Human Developmental Hematopoiesis. *Cell stem cell*, 28(3):472–487.e7, mar 2021. PMID: 33352111.
30. D. Pellin, M. Loperfido, C. Baricordi, S. L. Wolock, A. Montepeloso, O. K. Weinberg, A. Biffi, A. M. Klein, and L. Biasco. A comprehensive single cell transcriptional landscape of human hematopoietic progenitors. *Nature Communications*, 10(1):2395, 2019.
31. J. D. Buenrostro, M. R. Corces, C. A. Lareau, B. Wu, A. N. Schep, M. J. Aryee, R. Majeti, H. Y. Chang, and W. J. Greenleaf. Integrated Single-Cell Analysis Maps the Continuous Regulatory Landscape of Human Hematopoietic Differentiation. *Cell*, 173(6):1535–1548.e16, may 2018. PMID: 29706549.
32. A. Görgens, S. Radtke, M. Möllmann, M. Cross, J. Dürig, P. A. Horn, and B. Giebel. Revision of the Human Hematopoietic Tree: Granulocyte Subtypes Derive from Distinct Hematopoietic Lineages. *Cell Reports*, 3(5):1539–1552, 2013.
33. N. Schippel and S. Sharma. Dynamics of human hematopoietic stem and progenitor cell differentiation to the erythroid lineage. *Experimental Hematology*, 123:1–17, jul 2023.
34. R. Drissen, S. Thongjuea, K. Theilgaard-Mönch, and C. Nerlov. Identification of two distinct pathways of human myelopoiesis. *Science immunology*, 4(35), may 2019. PMID: 31126997.
35. S. Zheng, E. Papalexi, A. Butler, W. Stephenson, and R. Satija. Molecular transitions in early progenitors during human cord blood hematopoiesis. *Molecular Systems Biology*, 14(3):e8041, mar 2018.
36. H. Li, P. Côté, M. Kuoch, J. Ezike, K. Frenis, A. Afanassiev, L. Greenstreet, M. Tanaka-Yano, G. Tarantino, S. Zhang, J. Whangbo, V. L. Butty, E. Moiso, M. Falchetti, K. Lu, G. G. Connelly, V. Morris, D. Wang, A. F. Chen, G. Bianchi, G. Q. Daley, S. Garg, D. Liu, S. T. Chou, A. Regev, E. Lummertz da Rocha, G. Schiebinger, and R. G. Rowe. The dynamics of hematopoiesis over the human lifespan. *Nature methods*, 22(2):422–434, feb 2025. PMID: 39639169.
37. Y. Ben-David, B. Gajendran, K. M. Sample, and E. Zacksenhaus. Current insights into the role of Fli-1 in hematopoiesis and malignant transformation. *Cellular and molecular life sciences : CMLS*, 79(3):163, feb 2022. PMID: 35412146.
38. L. S. Bastian, B. A. Kwiatkowski, J. Breininger, S. Danner, and G. Roth. Regulation of the megakaryocytic glycoprotein IX promoter by the oncogenic Ets transcription factor Fli-1. *Blood*, 93(8):2637–2644, apr 1999. PMID: 10194443.
39. Y. Li, X. Qi, B. Liu, and H. Huang. The STAT5–GATA2 pathway is critical in basophil and mast cell differentiation and maintenance. *The Journal of Immunology*, 194(9):4328–4338, 2015.
40. A. B. Cantor and S. H. Orkin. Transcriptional regulation of erythropoiesis: an affair involving multiple partners. *Oncogene*, 21(21):3368–3376, may 2002. PMID: 12032775.
41. Y. Zhong, S. K. Walker, Y. Pritykin, C. S. Leslie, A. Y. Rudensky, and J. van der Veecken. Hierarchical regulation of the resting and activated T cell epigenome by major transcription factor families. *Nature immunology*, 23:122–134, 2022.
42. Y. Pritykin, J. van der Veecken, A. R. Pine, Y. Zhong, M. Sahin, L. Mazutis, D. Pe’er, A. Y. Rudensky, and C. S. Leslie. A unified atlas of CD8 T cell dysfunctional states in cancer and infection. *Molecular Cell*, 81(11):2477–2493, 2021.
43. T. Chu, M. Wu, B. Hoellbacher, G. P. de Almeida, C. Wurmser, J. Berner, L. V. Donhauser, A.-K. Gerullis, S. Lin, J. D. Cepeda-Mayorga, I. I. Kilb, L. Bongers, F. Toppeta, P. Strobl, B. Youngblood, A. M. Schulz, A. Zippelius, P. A.

M. Avdeeva et al.

- Knolle, M. Heinig, C.-P. Hackstein, and D. Zehn. Precursors of exhausted T cells are pre-emptively formed in acute infection. *Nature*, 640(8059):782–792, April 2025.
44. C. Gago da Graça, A. A. Sheikh, D. M. Newman, L. Wen, S. Li, J. Shen, Y. Zhang, S. S. Gabriel, D. Chisanga, J. Seow, A. Poch, L. Rausch, M.-H. T. Nguyen, J. Singh, C.-H. Su, L. A. Cluse, C. Tsui, T. N. Burn, S. L. Park, B. Von Scheidt, L. K. Mackay, A. Vasanthakumar, D. Bending, W. Shi, W. Cui, J. Schröder, R. W. Johnstone, A. Kallies, and D. T. Utzschneider. Stem-like memory and precursors of exhausted T cells share a common progenitor defined by ID3 expression. *Science Immunology*, 10(103):eadn1945, January 2025.
  45. M. L. Burger, A. M. Cruz, G. E. Crossland, G. Gaglia, C. C. Ritch, S. E. Blatt, A. Bhutkar, D. Canner, T. Kienka, S. Z. Tavana, A. L. Barandiaran, A. Garmilla, J. M. Schenkel, M. Hillman, I. de Los Rios Kobara, A. Li, A. M. Jaeger, W. L. Hwang, P. M. K. Westcott, M. P. Manos, M. M. Holovatska, F. S. Hodi, A. Regev, S. Santagata, and T. Jacks. Antigen dominance hierarchies shape TCF1(+) progenitor CD8 T cell phenotypes in tumors. *Cell*, 184(19):4996–5014.e26, sep 2021. PMID: 34534464.
  46. C. Tsui, L. Kretschmer, S. Rapelius, S. S. Gabriel, D. Chisanga, K. Knöpper, D. T. Utzschneider, S. Nüssing, Y. Liao, T. Mason, S. V. Torres, S. A. Wilcox, K. Kanev, S. Jarosch, J. Leube, S. L. Nutt, D. Zehn, I. A. Parish, W. Kastenmüller, W. Shi, V. R. Buchholz, and A. Kallies. MYB orchestrates T cell exhaustion and response to checkpoint inhibition. *Nature*, 609(7926):354–360, 2022.
  47. T. J. Stewart and S. I. Abrams. Altered immune function during long-term host-tumor interactions can be modulated to retard autochthonous neoplastic growth. *The Journal of Immunology*, 179(5):2851–2859, 2007.
  48. J. W. Zhu, S. J. Field, L. Gore, M. Thompson, H. Yang, Y. Fujiwara, R. D. Cardiff, M. Greenberg, S. H. Orkin, and J. DeGregori. E2F1 and E2F2 determine thresholds for antigen-induced T-cell proliferation and suppress tumorigenesis. *Molecular and cellular biology*, 21(24):8547–8564, 2001.
  49. L. Wu, C. Timmers, B. Maiti, H. I. Saavedra, L. Sang, G. T. Chong, F. Nuckolls, P. Giangrande, F. A. Wright, S. J. Field, et al. The E2F1–3 transcription factors are essential for cellular proliferation. *Nature*, 414(6862):457–462, 2001.
  50. D. DeRyckere and J. DeGregori. E2F1 and E2F2 are differentially required for homeostasis-driven and antigen-induced T cell proliferation in vivo. *The Journal of Immunology*, 175(2):647–655, 2005.
  51. M. A. Travis and D. Sheppard. TGF- $\beta$  activation and function in immunity. *Annual review of immunology*, 32(1):51–82, 2014.
  52. F. Macian. NFAT proteins: key regulators of T-cell development and function. *Nature Reviews Immunology*, 5(6):472–484, 2005.
  53. M. R. Müller and A. Rao. NFAT, immunity and cancer: a transcription factor comes of age. *Nature Reviews Immunology*, 10(9):645–656, 2010.
  54. A. E. Moran, K. L. Holzappel, Y. Xing, N. R. Cunningham, J. S. Maltzman, J. Punt, and K. A. Hogquist. T cell receptor signal strength in Treg and iNKT cell development demonstrated by a novel fluorescent reporter mouse. *Journal of Experimental Medicine*, 208(6):1279–1289, 2011.
  55. J. Chen, I. F. López-Moyado, H. Seo, C.-W. J. Lio, L. J. Hempleman, T. Sekiya, A. Yoshimura, J. P. Scott-Browne, and A. Rao. NR4A transcription factors limit CAR T cell function in solid tumours. *Nature*, 567(7749):530–534, 2019.
  56. H. Mao, M. Jia, M. Di, E. Valenzi, X. T. Cai, R. Lafyatis, K. Zhang, and P. V. Benos. HALO: hierarchical causal modeling for single cell multi-omics data. *Nature Communications*, 16(1):8892, 2025.
  57. M. Mørup and L. K. Hansen. Archetypal analysis for machine learning and data mining. *Neurocomputing*, 80:54–63, 2012.
  58. D. J. Di Bella, E. Habibi, R. R. Stickels, G. Scalia, J. Brown, P. Yadollahpour, S. M. Yang, C. Abbate, T. Biancalani, E. Z. Macosko, F. Chen, A. Regev, and P. Arlotta. Molecular logic of cellular diversification in the mouse cerebral cortex. *Nature*, 595(7868):554–559, 2021.
  59. E. Laurenti and B. Göttgens. From haematopoietic stem cells to complex differentiation landscapes. *Nature*, 553(7689):418–426, 2018.
  60. L. A. Liggett and V. G. Sankaran. Unraveling Hematopoiesis through the Lens of Genomics. *Cell*, 182(6):1384–1400, sep 2020.

61. F. A. Wolf, P. Angerer, and F. J. Theis. SCANPY: large-scale single-cell gene expression data analysis. *Genome biology*, 19:1–5, 2018.
62. L. D. Martens, D. S. Fischer, V. A. Yépez, F. J. Theis, and J. Gagneur. Modeling fragment counts improves single-cell ATAC-seq analysis. *Nature Methods*, 21(1):28–31, 2024.
63. A. Schep and S. University. Motifmatchr: Fast Motif Matching in R, 2022.
64. J. Vierstra, J. Lazar, R. Sandstrom, J. Halow, K. Lee, D. Bates, M. Diegel, D. Dunn, F. Neri, E. Haugen, E. Rynes, A. Reynolds, J. Nelson, A. Johnson, M. Frerker, M. Buckley, R. Kaul, W. Meuleman, and J. A. Stamatoyannopoulos. Global reference mapping of human transcription factor footprints. *Nature*, 583(7818):729–736, July 2020.
65. Single Cell Multiome ATAC + Gene Expression dataset analyzed using Cell Ranger ARC 2.0.0. PBMC from a Healthy Donor - Granulocytes Removed Through Cell Sorting (10k). <https://www.10xgenomics.com/datasets/pbmc-from-a-healthy-donor-granulocytes-removed-through-cell-sorting-10-k-1-standard-2-0-0>. Accessed: 2024-10-20.
66. Y. Gao, C. T. J. van Velthoven, C. Lee, E. D. Thomas, R. Mathieu, A. P. Ayala, S. Barta, D. Bertagnolli, J. Campos, T. Cardenas, D. Carey, T. Casper, A. B. Chakka, R. Chakrabarty, M. Chiang, L. Ching, M. Clark, M. J. Desierto, R. Ferrer, J. Gloe, J. Goldy, N. Guilford, J. Guzman, C. R. Halterman, S. D. Hastings, D. Hirschstein, W. Ho, K. James, Z. Juneau, N. Martin, R. McCue, E. Meyerdierks, A. C. Mitchell, J. S. Nagra, B. Nguy, T. N. Nguyen, P. Olsen, A. A. Oyama, N. Pena, J. Quon, Q. Ren, A. Ruiz, N. V. Shapovalova, J. Sulc, A. Torkelson, A. Tran, H. Tung, N. Valera Cuevas, J. Wang, J. Ariza, D. A. M. McMillen, J. Waters, M. Kunst, K. Ronellenfitch, B. Levi, M. J. Hawrylycz, C. Pagan, N. Dee, K. A. Smith, B. Tasic, Z. Yao, and H. Zeng. Continuous cell-type diversification in mouse visual cortex development. *Nature*, 647(8088):127–142, 2025.
67. M. E. Schroeder, D. M. McCormack, L. R. Metzner, J. Kang, K. X. Li, E. Yu, L. Melamed, K. M. Levandowski, H. Zaniwski, Q. Zhang, E. S. Boyden, F. M. Krienen, and G. Feng. A transcriptomic atlas of astrocyte heterogeneity across space and time in mouse and marmoset. *Neuron*, 113(23):3942–3965.e19, dec 2025.
68. E. Y. van Battum, M. H. van den Munkhof, and R. J. Pasterkamp. Novel insights into the regulation of neuron migration by axon guidance proteins. *Current Opinion in Neurobiology*, 92:103012, 2025.
69. S. L. Wolock, R. Lopez, and A. M. Klein. Scrublet: Computational Identification of Cell Doublets in Single-Cell Transcriptomic Data. *Cell Systems*, 8(4):281–291.e9, April 2019.
70. M. Avdeeva. ArchVelo: Archetypal Velocity Modeling for Single-cell Multi-omic Trajectories. *Zenodo*, <https://doi.org/10.5281/zenodo.20085695>, 2026.

## Acknowledgments

We thank all members of the Pritykin lab and the Developmental Dynamics group at the Flatiron Institute for helpful discussions. We thank Lucy Reading-Ikkanda (Flatiron Institute) for assistance with the graphic design of figures. The Flatiron Institute is a division of the Simons Foundation.

## Funding

This work was supported by the NIH/NIAID grant DP2AI171161, AACR-Bristol-Myers Squibb Immuno-oncology Research Fellowship (19-40-15-PRIT), Ludwig Institute for Cancer Research and Weill Cancer Hub East (Y.P.); Boehringer Ingelheim, the Austrian Science Fund (FWF, 10.55776/PAT4163824), and the ERC grant ERC-2023-STG 101116251 (J.v.d.V.); NIH grants P30 CA008748 and R01 AI034206 (A.Y.R.). A.Y.R. is an investigator with the Howard Hughes Medical Institute.

## Author contributions

Conceptualization, M.A. and Y.P.; methodology, M.A. and Y.P.; investigation, M.A., S.K.W. and J.v.d.V.; software, M.A.; writing – original draft, M.A. and Y.P.; writing – review & editing, M.A., S.K.W., J.v.d.V., A.Y.R. and Y.P.; funding acquisition, J.v.d.V., A.Y.R. and Y.P.; supervision, A.Y.R. and Y.P.

## Competing interests

A.Y.R. is an SAB member and has equity in Sonoma Biotherapeutics, RAPT Therapeutics, Coherus BioSciences, Santa Ana Bio, Odyssey Therapeutics, Nilo Therapeutics, and Vedanta Biosciences; he is also an SAB member of BioInvent and Amgen and a co-inventor of a CCR8+ Treg cell depletion IP licensed to Takeda, which is unrelated to the content of this publication. The remaining authors declare no competing interests.

## Figure Legends

### Fig. 1. Single-cell ATAC-seq archetypes improve modeling of single-cell RNA-seq data.

A multi-omic dataset of 12K human peripheral blood mononuclear cells (PBMCs) from 10x Genomics,<sup>65</sup> subsampled to 3,000 cells, is used to illustrate archetypal analysis.

(A) Left: A scATAC+RNA-seq experiment yields a gene-by-cell scRNA-seq matrix  $X$  and a peak-by-cell scATAC-seq matrix  $Y$ . Right: Uniform Manifold Approximation and Projection (UMAP) plot of the scRNA-seq data modality, with cell type annotations based on canonical marker genes.

(B) UMAP showing normalized scRNA-seq expression of gene BCL11A, which is active in multiple cell types.

(C) Pseudo-bulk scATAC-seq signal at the BCL11A locus.

(D) Schematic of archetypal analysis of scATAC-seq data. After preprocessing and normalization, the chromatin accessibility profile of each cell (matrix  $Y$ ) is approximated as a convex combination of  $K$  archetypal states. This decomposition yields  $K$  characteristic, archetypal chromatin accessibility profiles for cells (matrix  $A$ ), alongside  $K$  corresponding loadings for peaks (matrix  $Z$ ) that allow to approximate the accessibility at each individual peak with a convex combination of the  $K$  archetypes.

(E) Archetypal analysis with  $K = 8$  applied to the PBMC dataset. Left: archetypal scores of peak summits (shorter, equally-sized regions of maximal accessibility within the peaks) in the BCL11A locus. Right: average cell type archetypal scores, min-max normalized.

(F) Schematic of gene expression (scRNA-seq) prediction using chromatin accessibility (scATAC-seq) archetypes in paired multi-omic (scATAC+RNA-seq) data. To evaluate model performance, cells are partitioned into independent training and test sets. A cross-validated ridge regression model is trained individually for each gene, using the archetypal accessibility features as predictors and the scRNA-seq gene expression as the target. To prevent data leakage, the archetypal analysis is performed exclusively on the training set, and these archetypes are then projected onto the held-out test set to generate unbiased gene expression predictions for evaluation.

(G) Comparison of our regression model (panel F) with alternative approaches. Pearson correlation ( $r$ ) between predicted and observed gene expression profiles in held-out test cells is plotted. Left: baseline regression using raw scATAC-seq-derived peak accessibility linked to each gene. Center: ArchR gene activity scores.<sup>23</sup> Right: predictions from BABEL.<sup>24</sup> Two-sided Wilcoxon signed-rank test. P-values are below numerical precision limit ( $10^{-308}$ ).

**Fig. 2. Schematic of the ArchVelo model.**

(A) ArchVelo models the dynamics of chromatin accessibility, transcription, splicing, and degradation. Chromatin accessibility is represented by time-dependent profiles  $a_k(t)$  which open and close over latent time  $t$ . Each archetype contributes linearly to the transcription rate of each gene, combining gene-level loadings  $z_k^g = \sum_{p \in \mathcal{P}(g)} z_{pk}$  with archetype-specific transcription rates  $\alpha_k^g$ . Figure adapted from BioRender, <https://BioRender.com/kpe5dzb>.

(B) Example model fits for a representative gene. Left: chromatin accessibility profile for one archetype  $a_k(t)$  (orange line); middle: unspliced normalized imputed counts; right: spliced normalized imputed counts. Data points from single cells are shown in black.

**Fig. 3. ArchVelo enables improved cell trajectory inference in the mouse embryonic brain.**

Comparison of ArchVelo against state-of-the-art methods on a multi-omic scATAC+RNA-seq dataset from mouse embryonic brain.<sup>25</sup>

(A) UMAP of the scRNA-seq component with streamplot of the velocity field inferred by ArchVelo.

(B) Dotplot of marker gene expression across cell types. Dot size, fraction of cells expressing the gene; color intensity, min-max normalized expression. Abbreviations as in (A); also: M: migrating neurons, RG: radial glia, OPC: oligodendrocyte progenitor cells, N: neurons.

(C) Kinetic fits for the gene *Satb2* in the unspliced-spliced ( $u-s$ ) phase space. ArchVelo is compared with MultiVelo, the only other evaluated method incorporating scATAC-seq data.

(D) Left: log-likelihood (pseudocount = 1) of  $u-s$  phase space fits for *Satb2* using ArchVelo, MultiVelo, scVelo and VeloVI (see also (G)). Right: scatterplot comparing per-gene log-likelihoods for MultiVelo ( $x$ -axis) and ArchVelo ( $y$ -axis). Two-sided Wilcoxon signed-rank test.

(E) Clustermaps of the Spearman correlations between gene-specific latent time profiles from ArchVelo and MultiVelo.

(F) Scatterplot of silhouette scores across varying numbers of clusters from (E).  $x$ -axis: MultiVelo;  $y$ -axis: ArchVelo. Two-sided Wilcoxon signed-rank test.

(G) Benchmarking across all genes using two metrics: mean log-likelihood (pseudocount = 1) of  $u-s$  fits ( $x$ -axis) and silhouette score of latent time consistency ( $y$ -axis). Colors: as in panel D. A parameter sweep across eleven values of the  $w_c$  parameter (the relative ATAC-seq vs. RNA-seq contribution) was performed for ArchVelo and MultiVelo. Cross, default parameters.

(H) Benchmarking the reconstruction of established cell transitions. Left: reference graph of known cell transitions in the mouse embryonic brain differentiation. Annotations as in (A). Right: quantitative evaluation of inferred edge directionality using Cross-Boundary Direction Correctness (CBDir) scores<sup>28</sup> for the reference graph edges. The average score across all reference edges is shown for each method. (See **Methods** and **Supplementary Fig. 1** for details.)

**Fig. 4. ArchVelo recapitulates human hematopoietic stem cell differentiation trajectories and outperforms existing methods.**

Comparison of ArchVelo against other methods on a multi-omic dataset for human hematopoietic stem cell differentiation.<sup>9</sup>

(A) UMAP of the HSC differentiation dataset with cell type annotations and velocity fields inferred by ArchVelo. HSC: human hematopoietic stem cells; MPP: multipotent progenitors; LMPP: lymphoid-primed multipotent progenitors; EMP: erythroid-myeloid progenitors; GMP: granulocyte–macrophage progenitors; Prog DC: progenitor dendritic cells; Prog N: progenitor neutrophils; BEM: basophil-eosinophil-mast progenitors; MEP: megakaryocyte–erythrocyte progenitors; Prog MK: progenitor megakaryocytes; Ery: early erythrocytes.

(B) Dotplot of marker gene expression across annotated cell types. Dot size, fraction of cells expressing the gene; color intensity, min–max normalized expression.

(C) Scatterplot comparing per-gene  $u$ - $s$  log-likelihoods (pseudocount = 1) for MultiVelo ( $x$ -axis) and ArchVelo ( $y$ -axis). Two-sided Wilcoxon signed-rank test.

(D) Scatterplot of silhouette scores for latent time correlation clusters across varying numbers of clusters.  $x$ -axis: MultiVelo;  $y$ -axis: ArchVelo. Two-sided Wilcoxon signed-rank test.

(E) Benchmarking across all genes using two metrics: mean log-likelihood (pseudocount = 1) of  $u$ - $s$  fits ( $x$ -axis) and silhouette score of latent time consistency ( $y$ -axis). Colors correspond to methods (as in (F) and Fig. 3D,G). A parameter sweep across eleven values of the  $w_c$  parameter (the relative ATAC-seq vs. RNA-seq contribution) was performed for ArchVelo and MultiVelo. Cross, default parameters.

(F) Benchmarking the reconstruction of established cell transitions. Left: reference graph of known cell transitions in hematopoietic differentiation. Nodes represent annotated cell populations (as in (A) and (B)). Right: quantitative evaluation of inferred edge directionality using CBDir scores for the reference graph edges. The average score across all reference edges is shown for each method. (See Methods and Supplementary Fig. 2 for details.)

**Fig. 5. ArchVelo enables decomposition of the velocity field into archetypal regulatory components.**

Analysis was performed on the same multi-omic (scATAC+RNA-seq) embryonic mouse brain dataset as in **Fig. 3**.

**(A)** Heatmap of the average archetypal loadings across annotated cell clusters.

**(B-C)** Linear decomposition of ordinary differential equation variables into components driven by individual chromatin accessibility archetypes. **(B)** Decomposition of the spliced counts  $s$  for the gene *Satb2* into distinct components  $s_k$  associated with archetypes  $a_k$ . Colored dashed lines show individual components; the solid black line indicates the total predicted  $s$ . The scatterplot shows imputed spliced counts. **(C)** Same decomposition for the velocity  $v = \dot{s}$ , using the same color scheme as in **(B)**.

**(D)** Grid visualization of the ArchVelo velocity field projected onto the UMAP embedding, focusing on the glial compartment. Colors represent cell type annotations, as in **Fig. 3A**.

**(E-J)** Example velocity components for archetypes A5 and A8, shown on the same grid as in **(D)**, and the associated analysis. Although these archetypal components are concentrated on similar subsets of cells, they are associated with distinct local velocity patterns: A5 reflects cell cycle progression, whereas A8 reflects astrocyte lineage differentiation.

**(F)** The velocity field for archetype A5 overlaid with the imputed normalized expression of the proliferation marker *Mki67*.

**(G)** Cell cycle phase scores (G1, S and G2M), inferred from gene signatures and represented as the fraction of cells most enriched for each signature ( $y$ -axis), plotted against the latent time estimated from A5 (smoothed over 100 nearest neighbors). The temporal alignment demonstrates a strong correspondence with cell cycle dynamics.

**(I)** The velocity field for archetype A8 overlaid with imputed normalized expression of *Aldh1l1*, a canonical astrocyte marker that initiates in radial glia and is progressively upregulated during astrocyte differentiation.

**(J)** Normalized expression of the astrocyte marker genes *Aldh1l1*, *Slc6a11* and *Sparcl1* along the A8 latent time. The patterns indicate correspondence with astrocyte differentiation.

**Fig. 6. ArchVelo identifies candidate transcription factor drivers of distinct regulatory components of RNA velocity through motif analysis.**

Analysis was performed on the same multi-omic (scATAC+RNA-seq) human HSC dataset as in **Fig. 4**.

(A) Heatmap of the average archetypal loadings across annotated cell clusters.

(B) Velocity field components corresponding to archetypes A5, A4 and A8, visualized on a common embedding grid (same UMAP as in **Fig. 4A**). Each archetypal component captures a distinct differentiation trajectory. Top left: a subgraph of the reference cell type transitions (from **Fig. 4F**) for MEP differentiation to prog. MK and platelet cells (captured by archetype A5), to BEM cells (captured by A4), and to Ery cells (captured by A8).

(C) Transcription factor (TF) motif analysis for archetype A5. Left: Barplot of the top-scoring motifs for A5 (**Methods**). Right: UMAP plots showing per-cell TF motif scores, normalized scRNA-seq expression of a representative gene encoding the corresponding TF, and a marker gene that is a downstream target of the TF.

(D) Same as (C), for archetype A4.

(E) Same as (C), for archetype A8.

**Fig. 7. ArchVelo uncovers shared trajectories of cell differentiation and proliferation in CD8 T cells responding to acute and chronic viral infection.**

ArchVelo analysis of a multi-omic scATAC+RNA-seq dataset for cytotoxic CD8 T cells at day 7 post-infection with either acute (Armstrong) or chronic (clone 13) lymphocytic choriomeningitis virus (LCMV).<sup>21</sup>

(A) UMAP visualization of the scRNA-seq modality for Armstrong (top) and clone 13 (bottom), overlaid with ArchVelo-inferred trajectory stream plots. Datasets were clustered jointly. Cluster annotations: S: S phase of the cell cycle; G2M.1: early G2/M phase; G2M.2: late G2/M phase; Inter eff: intermediate effector cells; Klr $g1^{++}$  eff: effector cells expressing the highest levels of *Klr $g1$* ; Gzm $^{+}$  eff: effector cells enriched with expression of *Gzma* and *Gzmk*; Il7r $^{+}$  eff: effector cells with enriched expression of *Il7r*; Inter exh: intermediate exhausted cells; Exh: exhausted cells; Prog: progenitor cells.

(B) Dotplot of marker gene expression across annotated clusters, shown separately for each of the two LCMV clones. Dot size, fraction of cells expressing the gene; color intensity, min-max normalized expression.

(C) Composition of non-cycling clusters across the two infection conditions, LCMV Armstrong (acute infection) and LCMV clone 13 (chronic infection). Dashed vertical line: expected proportion based on condition-specific total non-cycling cell numbers.

(D) Heatmaps of ArchVelo transition probabilities between clusters (Methods), for the two infection conditions.

(E) Phase plots for genes *Mki67* and *Sell* for the two infection conditions. Cluster colors are the same as in (A). Black curve: ArchVelo fit.

**Fig. 8. ArchVelo uncovers a differentiation trajectory from  $Ccr6^{-}$  to  $Ccr6^{+}$  progenitor CD8 T cells.**

The analysis focused on the progenitor compartment of cytotoxic CD8 T cells at day 7 following acute (Armstrong) or chronic (clone 13) LCMV infection (same data as in **Fig. 7**).

(A) Zoomed-in view of ArchVelo-inferred trajectories within the progenitor CD8 T cell compartment and adjacent clusters in LCMV Armstrong (top) and LCMV clone 13 (bottom). The  $Ccr6^{-}$  and  $Ccr6^{+}$  progenitor subpopulations are highlighted.

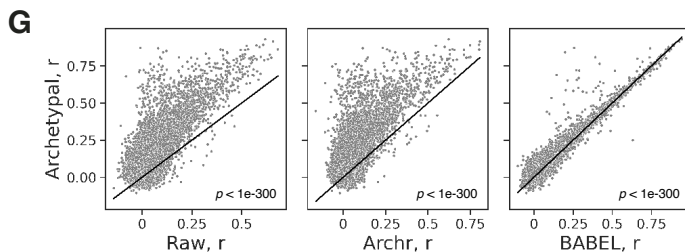
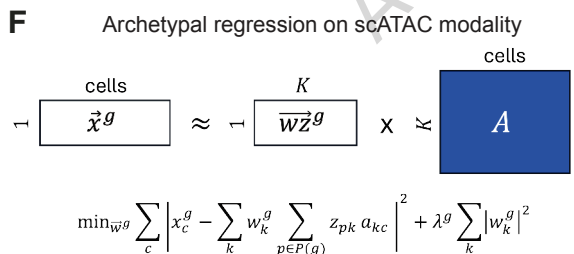
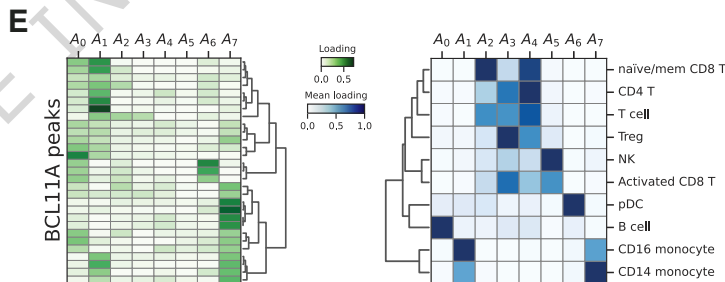
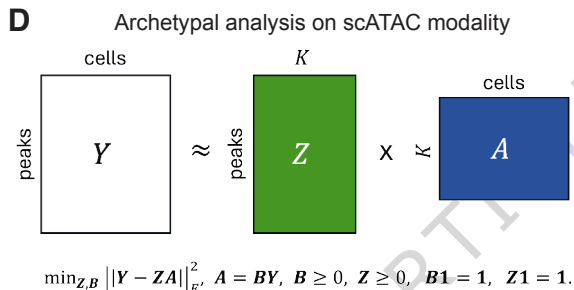
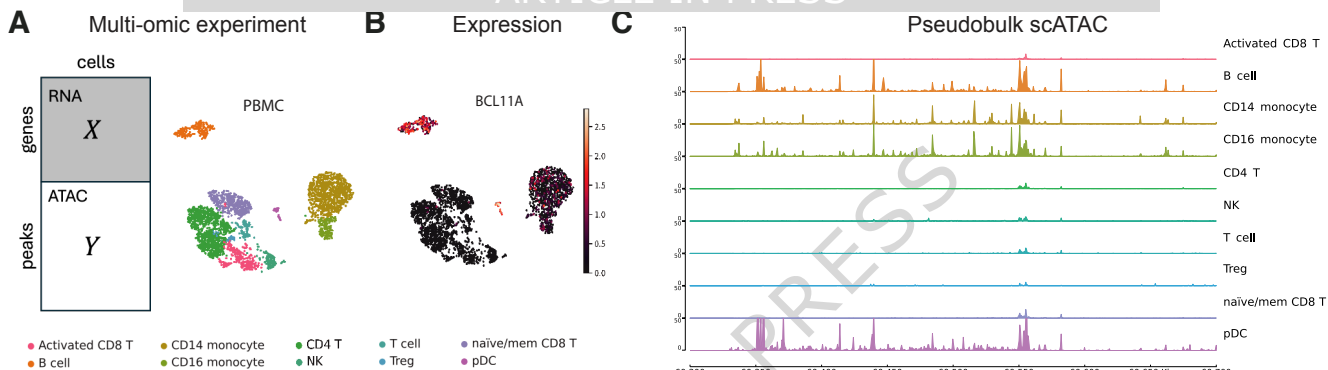
(B) Dotplots showing expression of known progenitor marker genes across  $Ccr6^{-}$  progenitors,  $Ccr6^{+}$  progenitors, and all other cells, as well as genes differentially expressed between  $Ccr6^{-}$  and  $Ccr6^{+}$  progenitors. Dot size, fraction of cells expressing the gene; color intensity, min–max normalized expression levels; bar size, cluster size.

(C) Heatmap of the average archetypal loadings across annotated cell clusters, min-max normalized, incorporating the refined annotation of progenitor subpopulations.

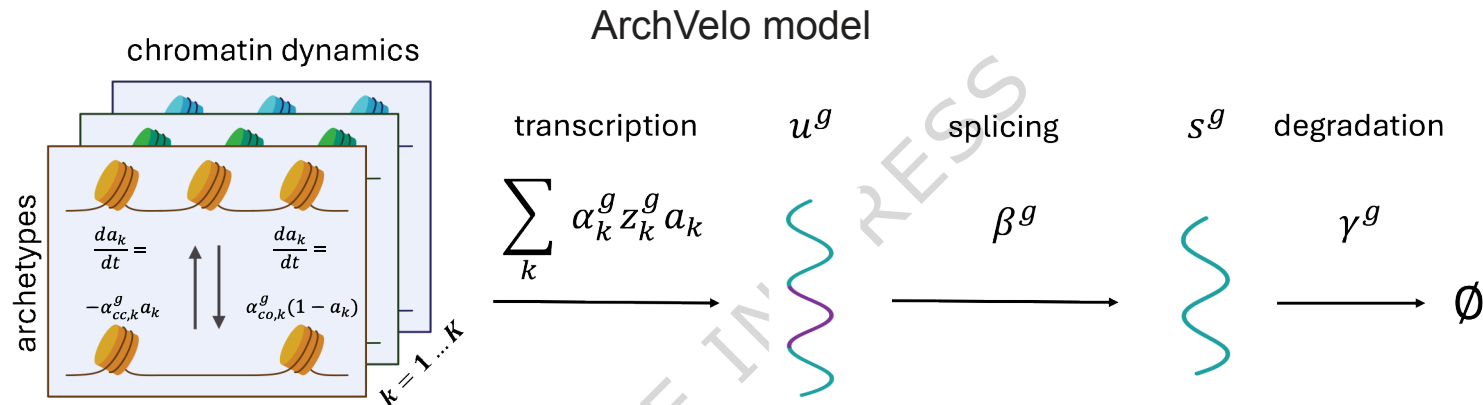
(D) Velocity field components corresponding to archetypes A8 (left) and A5 (right), visualized on the same UMAP embedding as in **Fig. 7A**. Top: LCMV Armstrong; bottom: LCMV clone 13.

(E) Archetypal decomposition reveals a shared differentiation trajectory across infection conditions within the progenitor subpopulation. Heatmaps show min–max normalized imputed scRNA-seq expression of selected genes (averaged over a sliding window of 20 cells), aligned by latent time. Color bars: latent time; subcluster  $Ccr6^{-}$  and  $Ccr6^{+}$  progenitor annotations; min–max normalized archetypal loadings for archetypes A8 and A5 (averaged over a sliding window of 30 cells).

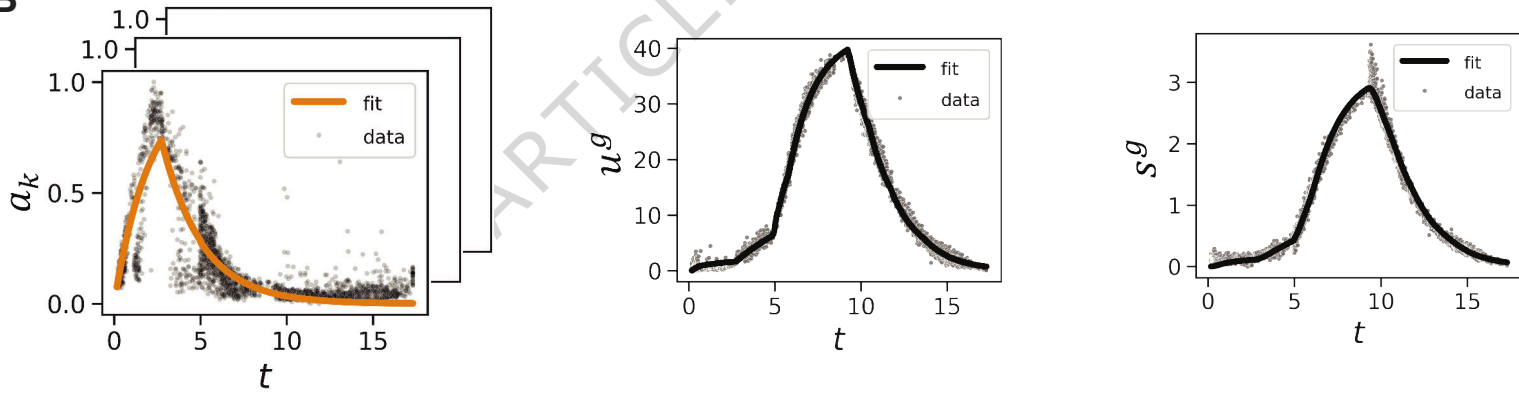
(F) Heatmap of TF motif correlation scores (for top and bottom 10 motives, **Methods**) for archetypes A5 and A8 within the progenitor subpopulation.



A

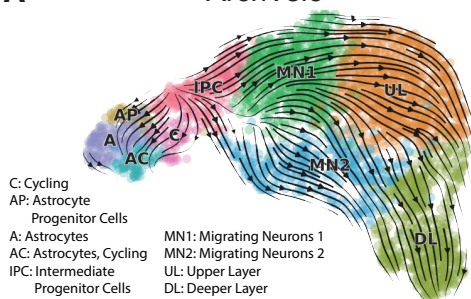


B

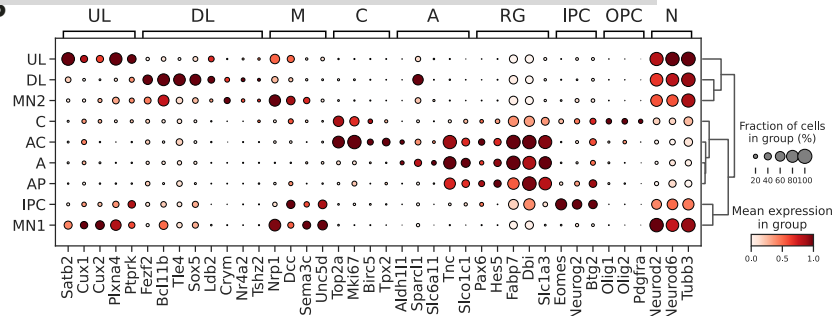


A

ArchVelo



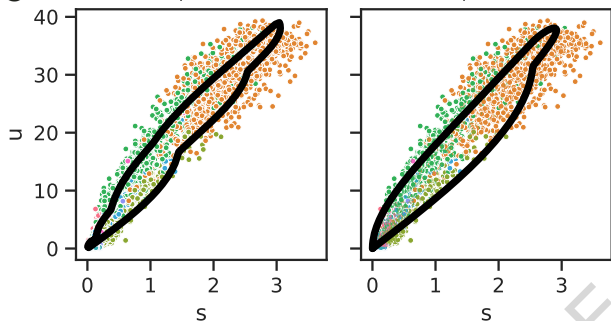
D



C

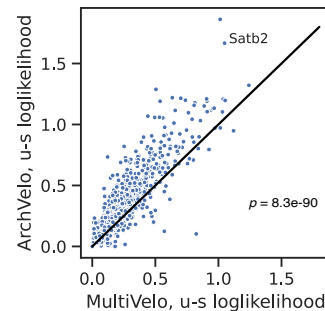
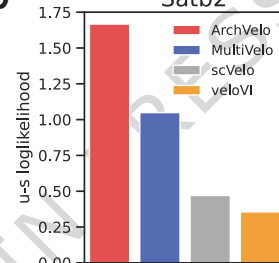
Satb2, ArchVelo

Satb2, MultiVelo

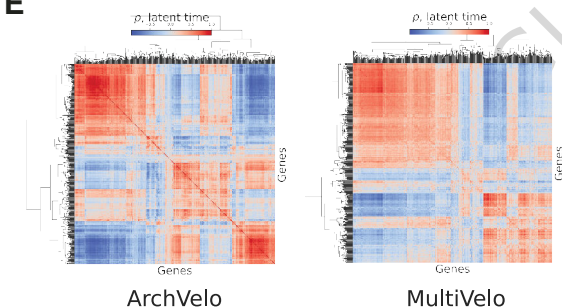


D

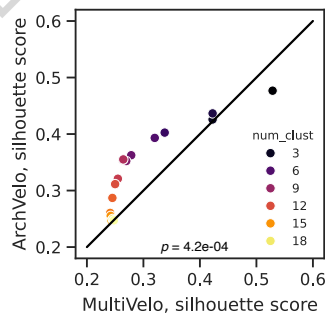
Satb2



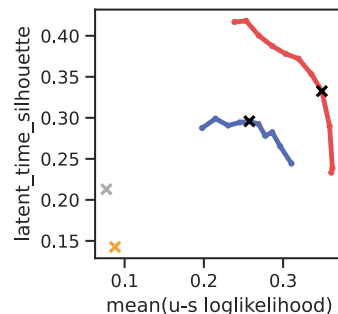
E



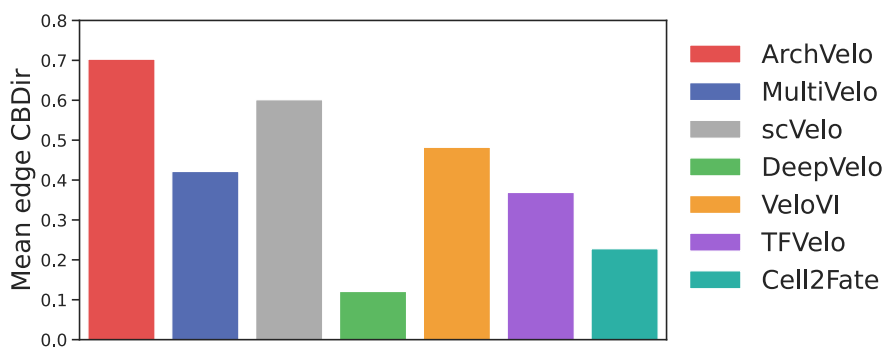
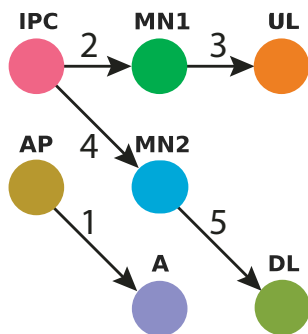
F

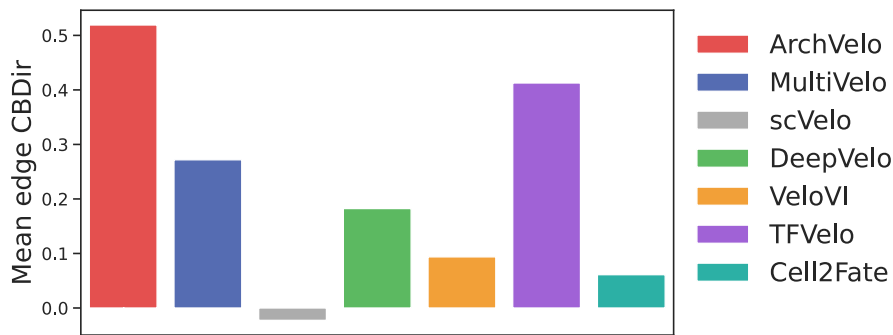
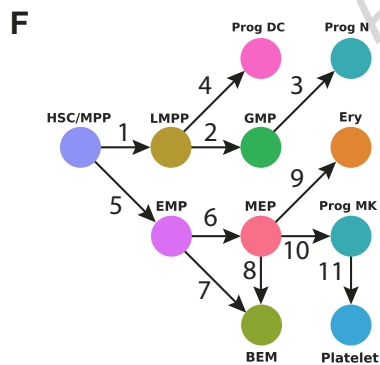
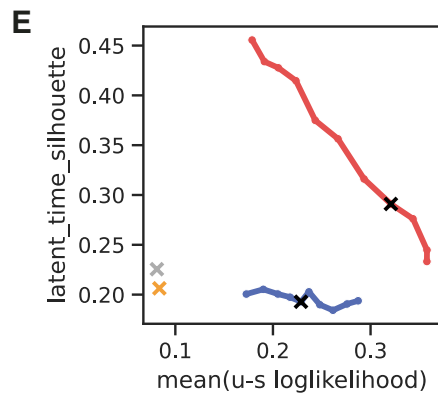
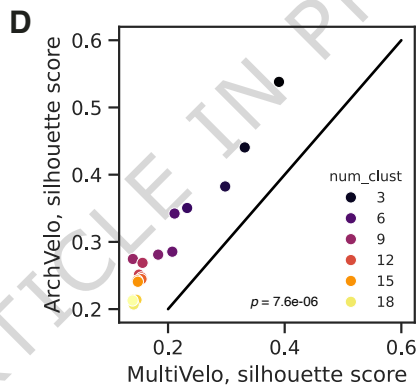
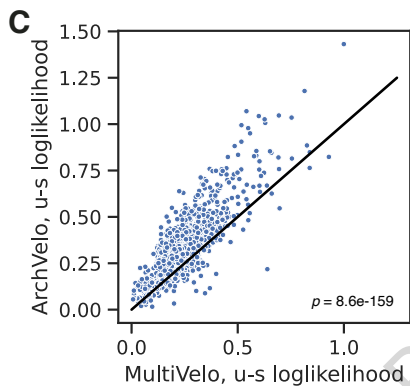
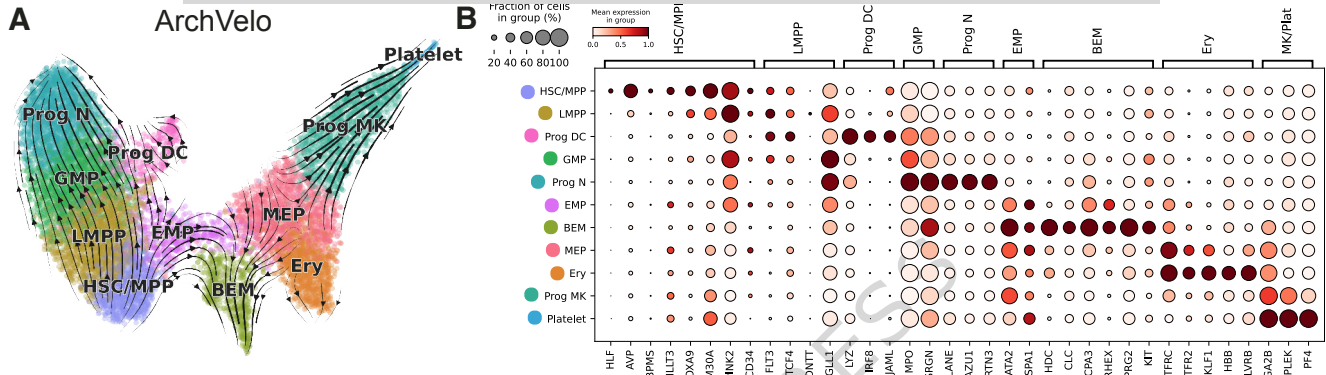


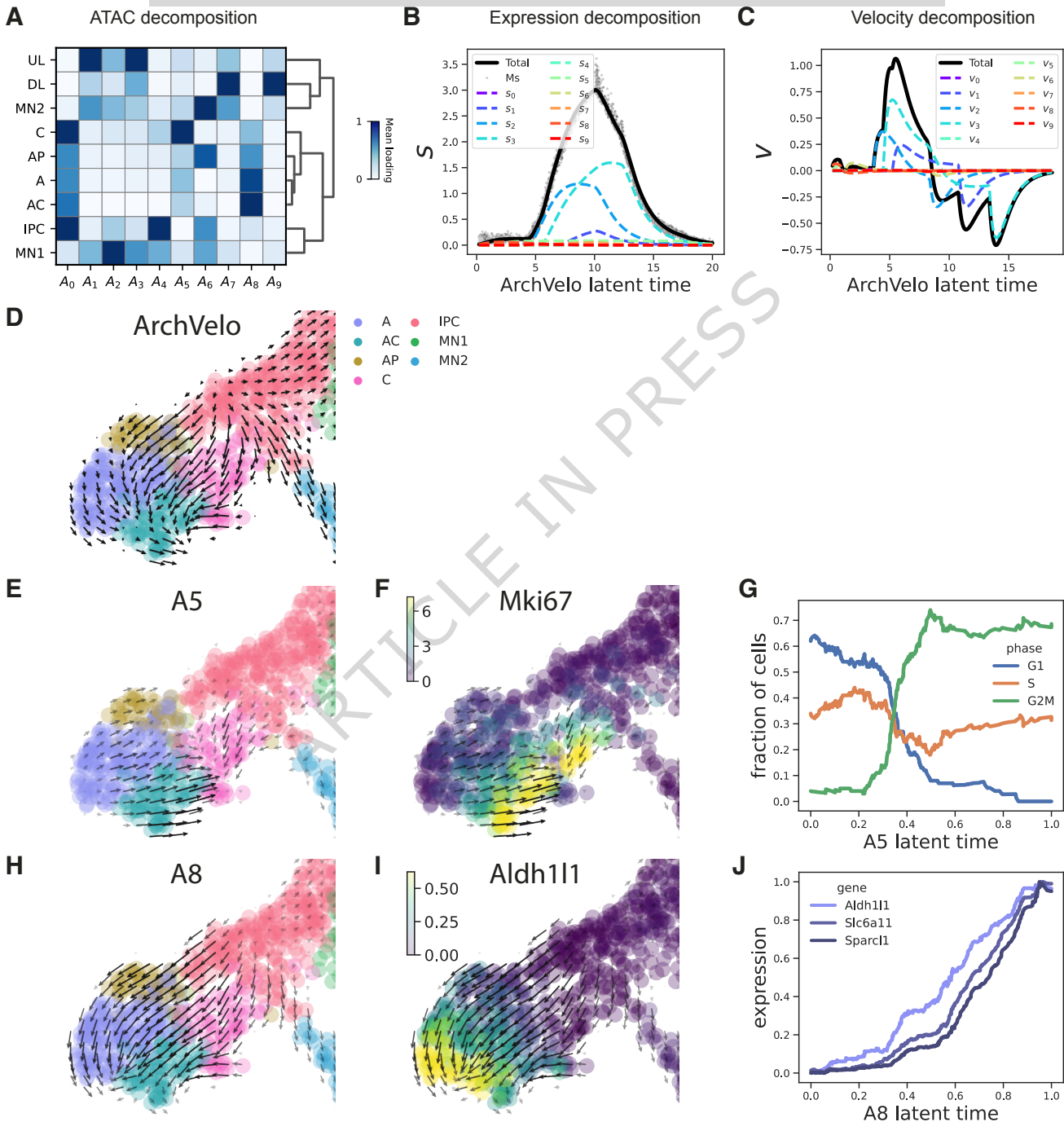
G

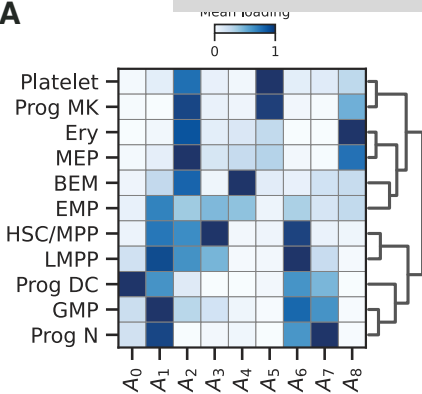
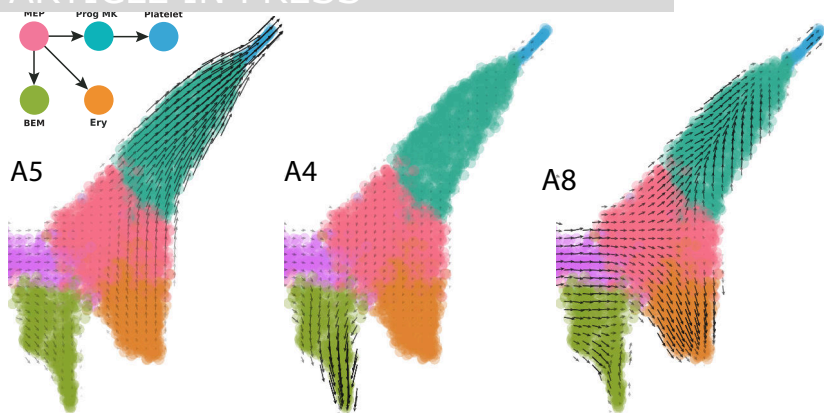
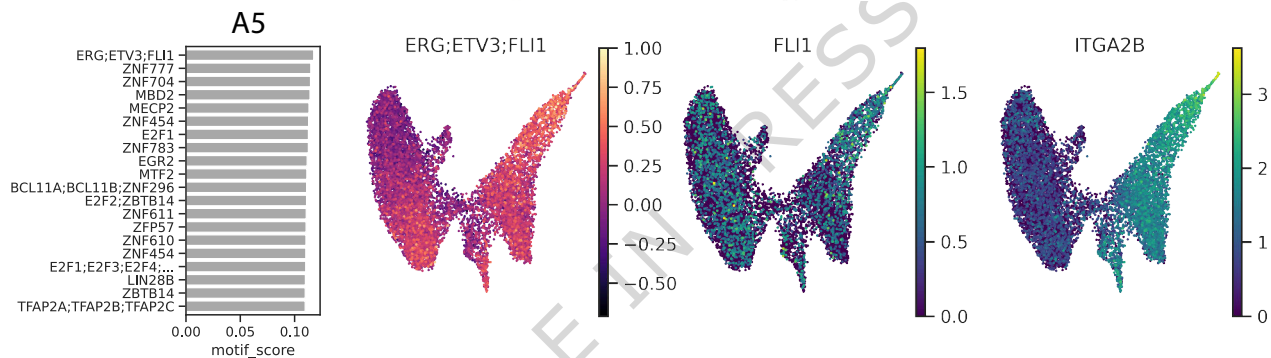
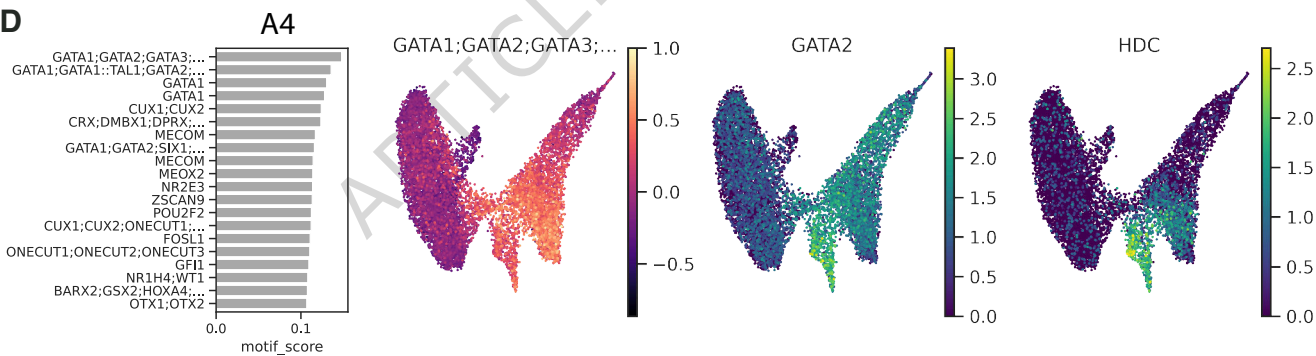
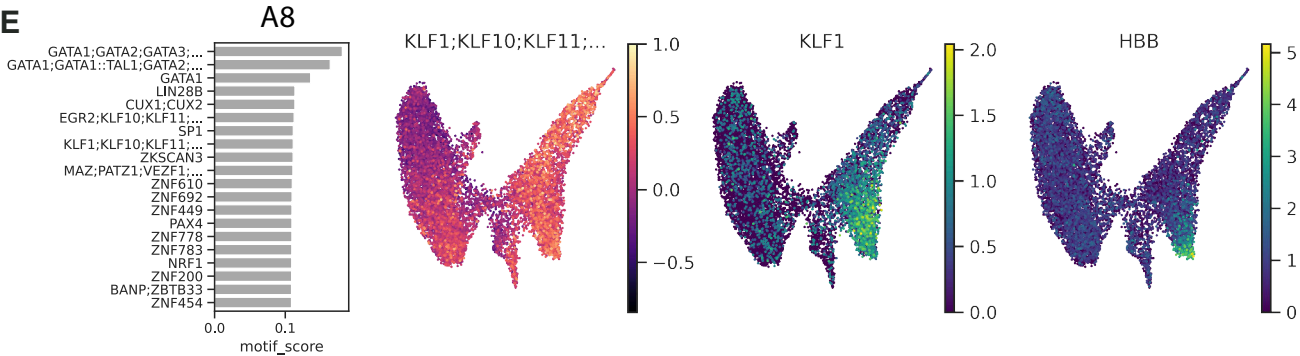


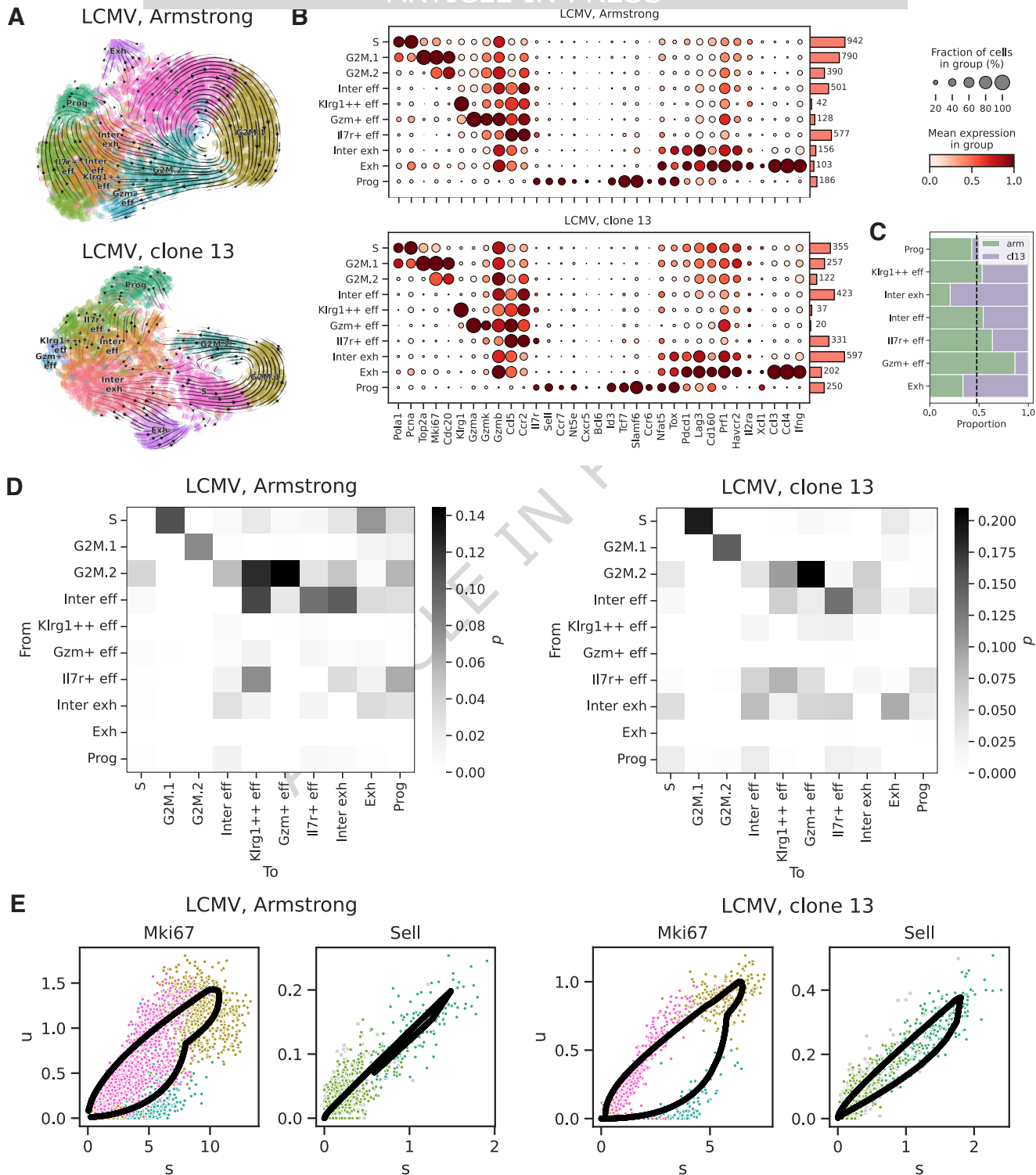
H

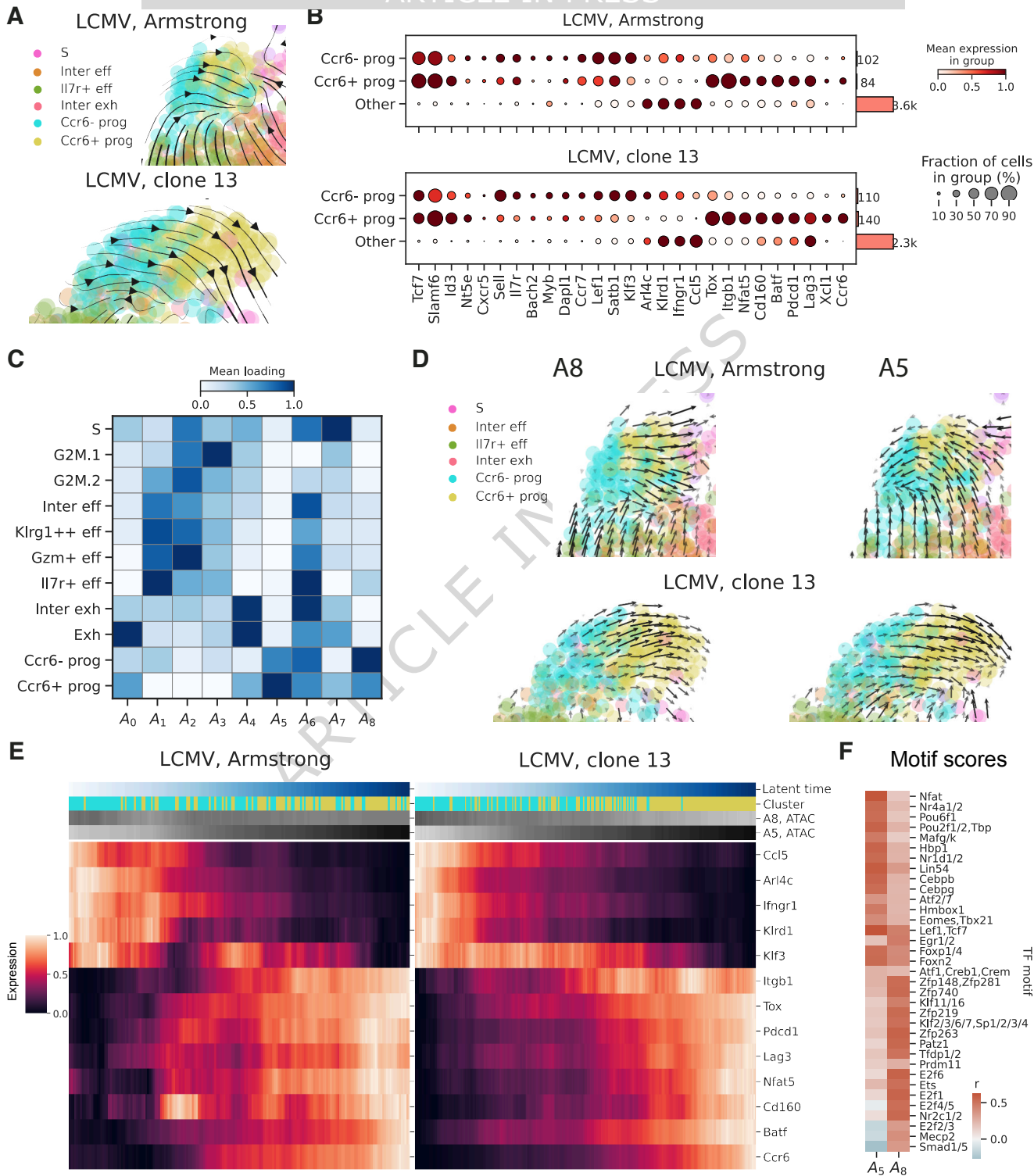






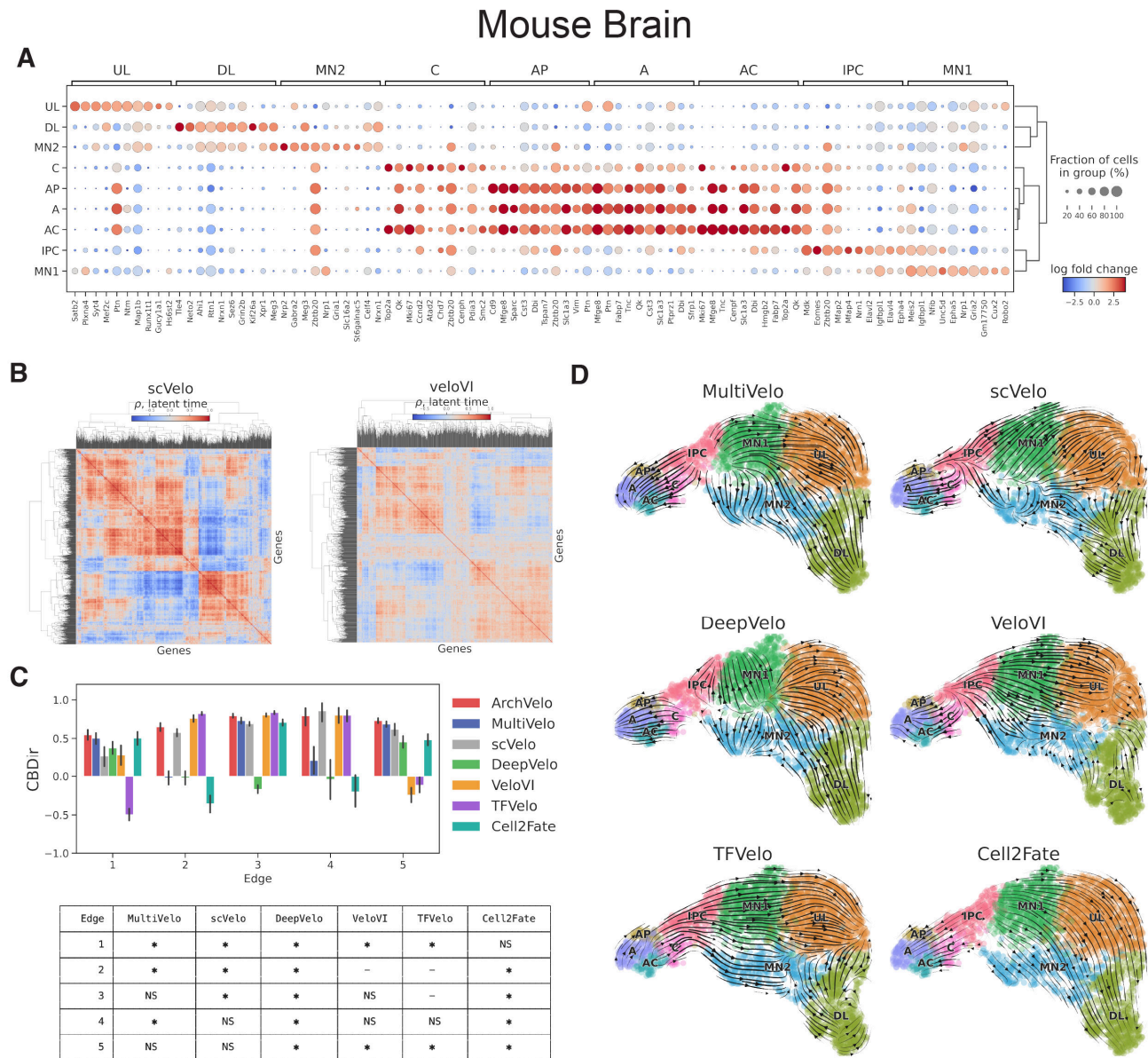
**A****B****C****D****E**





# Supplementary Information for “ArchVelo: Archetypal Velocity Modeling for Single-cell Multi-omic Trajectories”

## Supplementary Figures



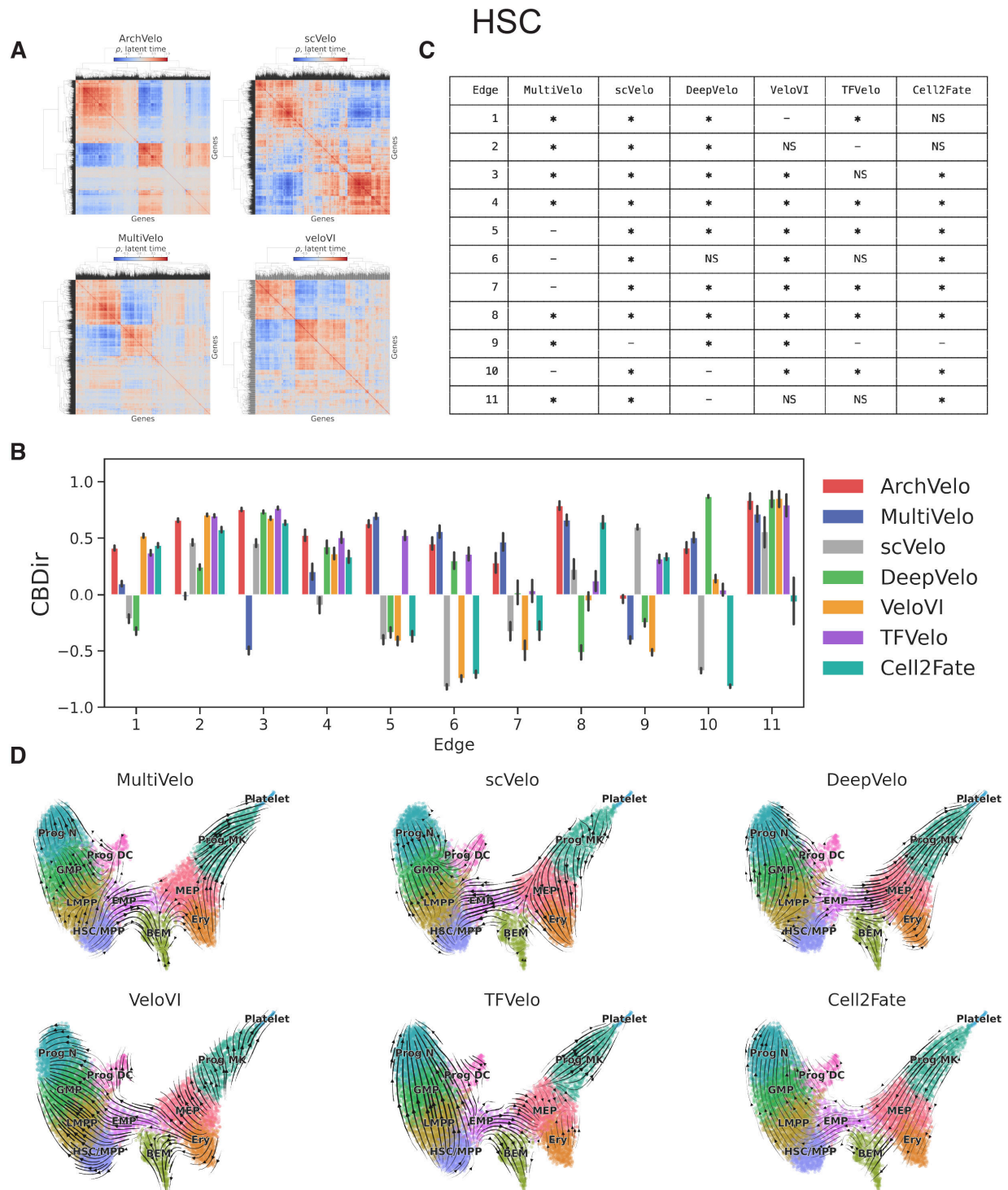
**Supplementary Figure 1. Extended benchmarking results for the mouse embryonic brain dataset.**

(A) Differential expression analysis between the identified clusters. For each cluster, the top 10 most significantly differentially overexpressed genes are shown. Color indicates the log-fold change of mean expression in the cluster compared to all other cells.

(B) Clustermaps of the Spearman correlations between gene-specific latent time profiles inferred by scVelo and VeloVI.

(C) Top: quantitative evaluation of inferred edge directionality using CDir scores for the reference graph edges (shown in **Fig. 3H**). Barplot shows the average score across boundary cell pairs from the cell types connected by an edge; error bar shows the 95% confidence interval. Bottom: table summarizing the statistical significance when comparing ArchVelo against each of the other methods (columns) for every reference edge (rows). Significance was determined using a two-sided paired Wilcoxon signed-rank test, comparison of medians was used to determine directionality of difference. To correct for multiple hypothesis testing, p-values were adjusted independently for each edge using the Benjamini-Hochberg procedure with a false discovery rate (FDR) threshold of 0.01. Asterisks (\*) indicate that the CDir scores for ArchVelo are significantly higher, minus signs (-) indicate that the CDir scores for the competing method are significantly higher than for ArchVelo; "NS" indicates no significant difference.

(D) Velocity stream plots inferred by other methods, shown on the same UMAP embedding as in **Fig. 3A**.

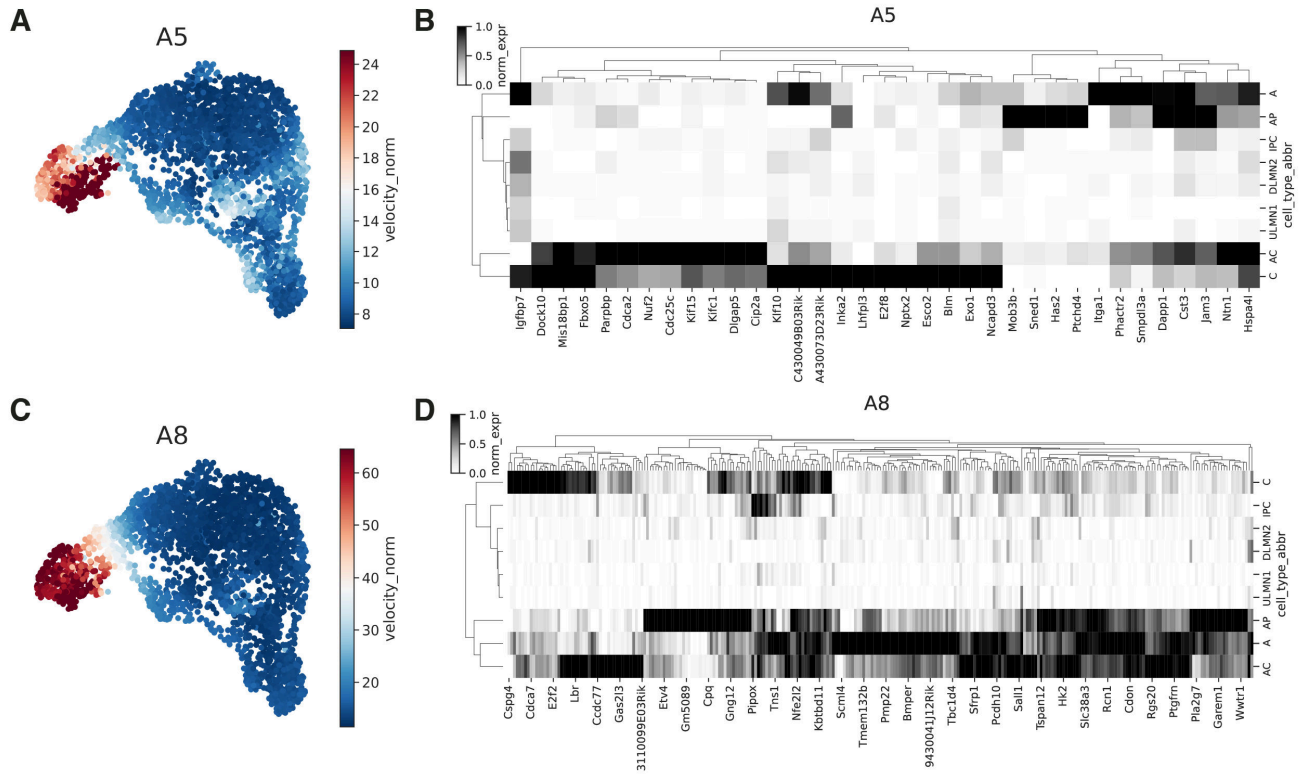


**Supplementary Figure 2. Extended benchmarking results for the human hematopoietic stem cell differentiation dataset.**

(A) Clustermaps of the Spearman correlations between gene-specific latent time profiles inferred by ArchVelo, scVelo, MultiVelo and VeloVI.

(B-C) Same as **Supplementary Figure 1C**, but for the HSC dataset and the reference graph in **Fig. 4F**.

(D) Velocity stream plots inferred by other methods, shown on the same UMAP embedding as in **Fig. 4A**.



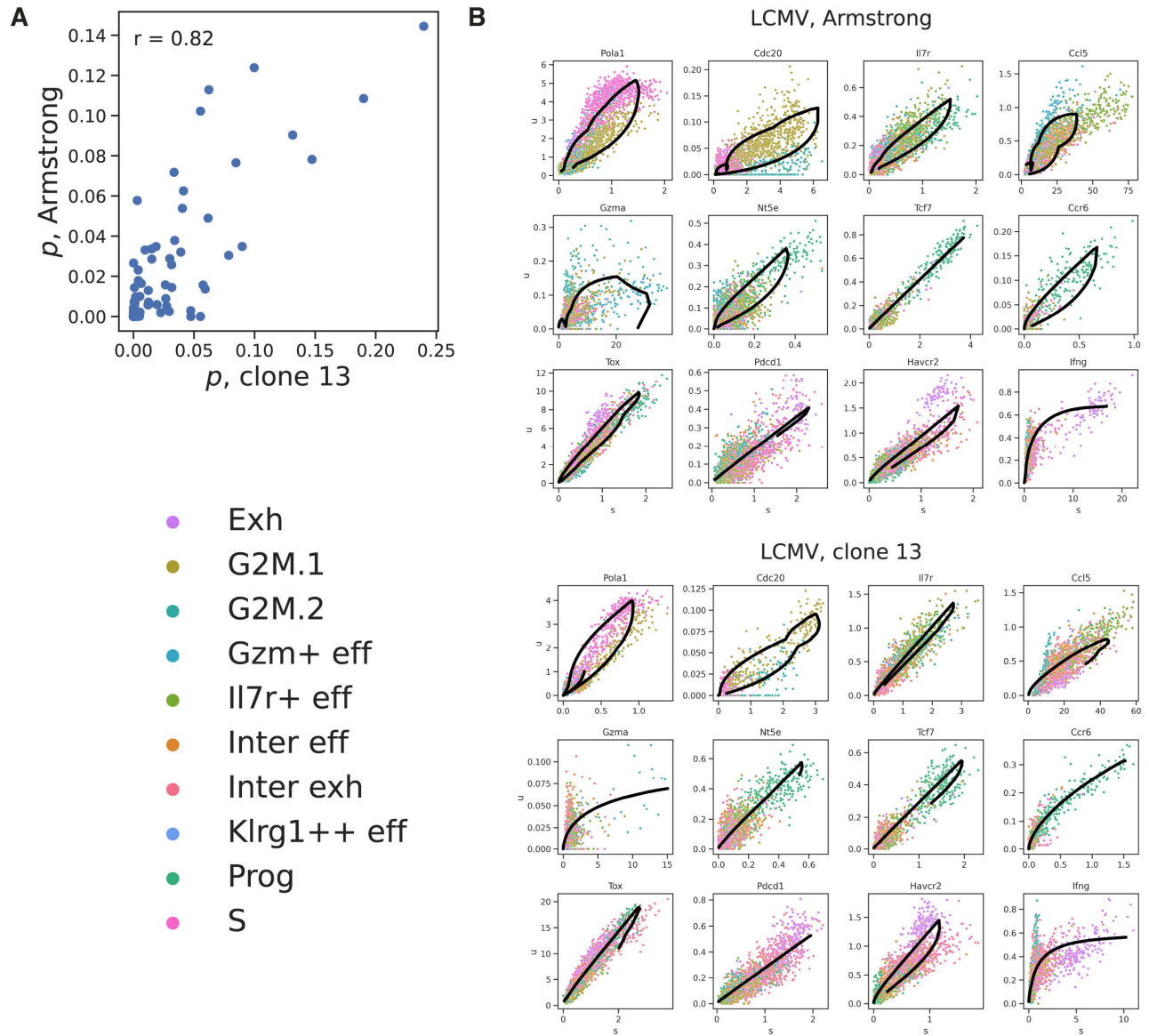
**Supplementary Figure 3. Velocity loadings and gene selection for archetypal components in the mouse embryonic brain dataset.**

(A) UMAP embedding colored by velocity loadings for archetype A5. Loadings are shown as the  $l_2$ -norms of the inferred velocity component associated with A5 (see **Methods**).

(B) Top genes contributing to archetype A5, selected as described in **Methods**.

(C) Same as panel A, for archetype A8.

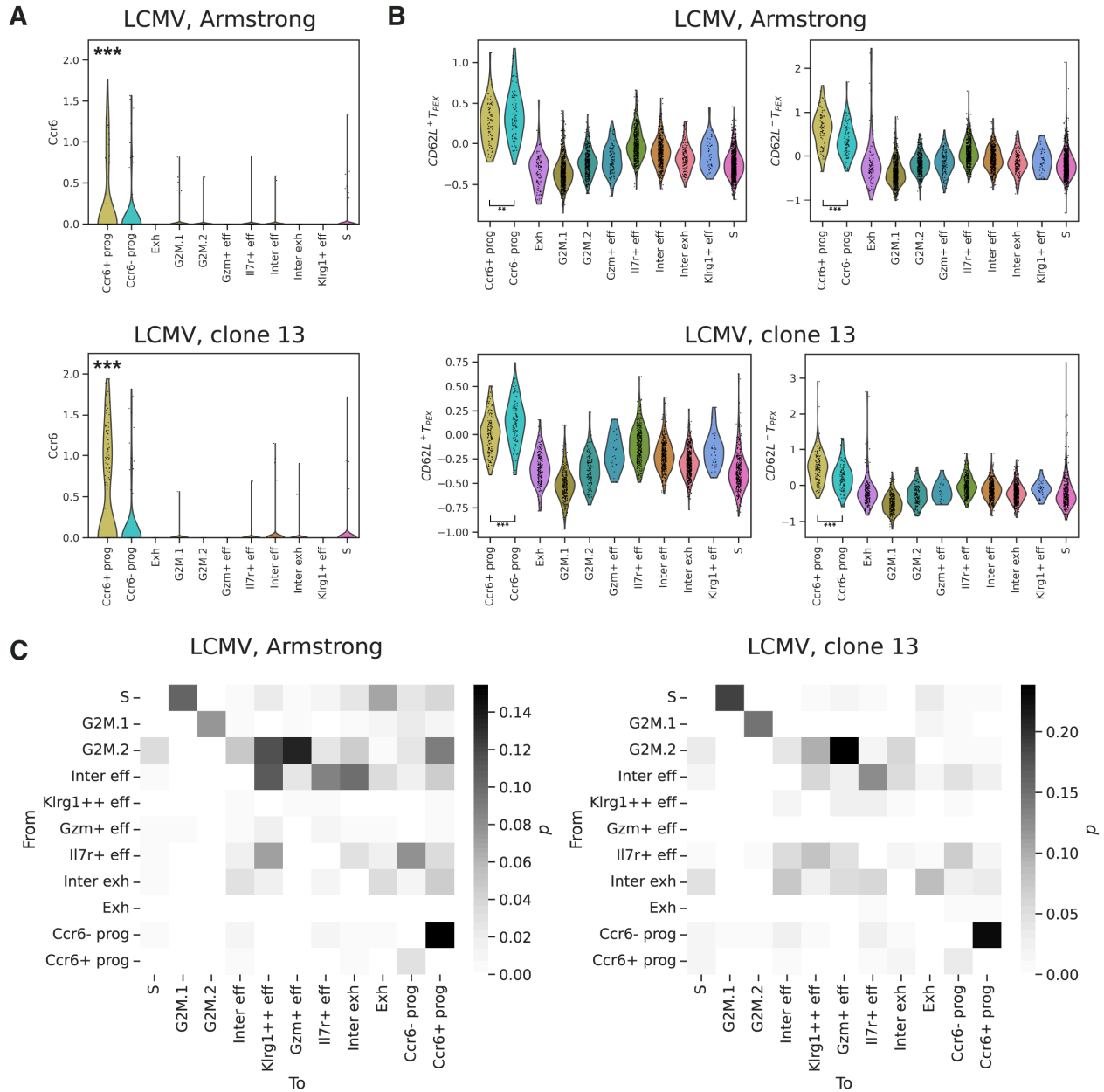
(D) Same as panel B, for archetype A8.



**Supplementary Figure 4. Additional details of the analysis of the CD8 T cell dataset.**

(A) Scatterplot comparing ArchVelo transition probabilities between the two LCMV infection conditions. Each point corresponds to a cluster pair.  $x$ -axis: LCMV clone 13;  $y$ -axis: LCMV Armstrong. Pearson correlation coefficient  $r = 0.82$ .

(B) Phase plots for selected genes in LCMV Armstrong and clone 13. Black curve: ArchVelo fit.

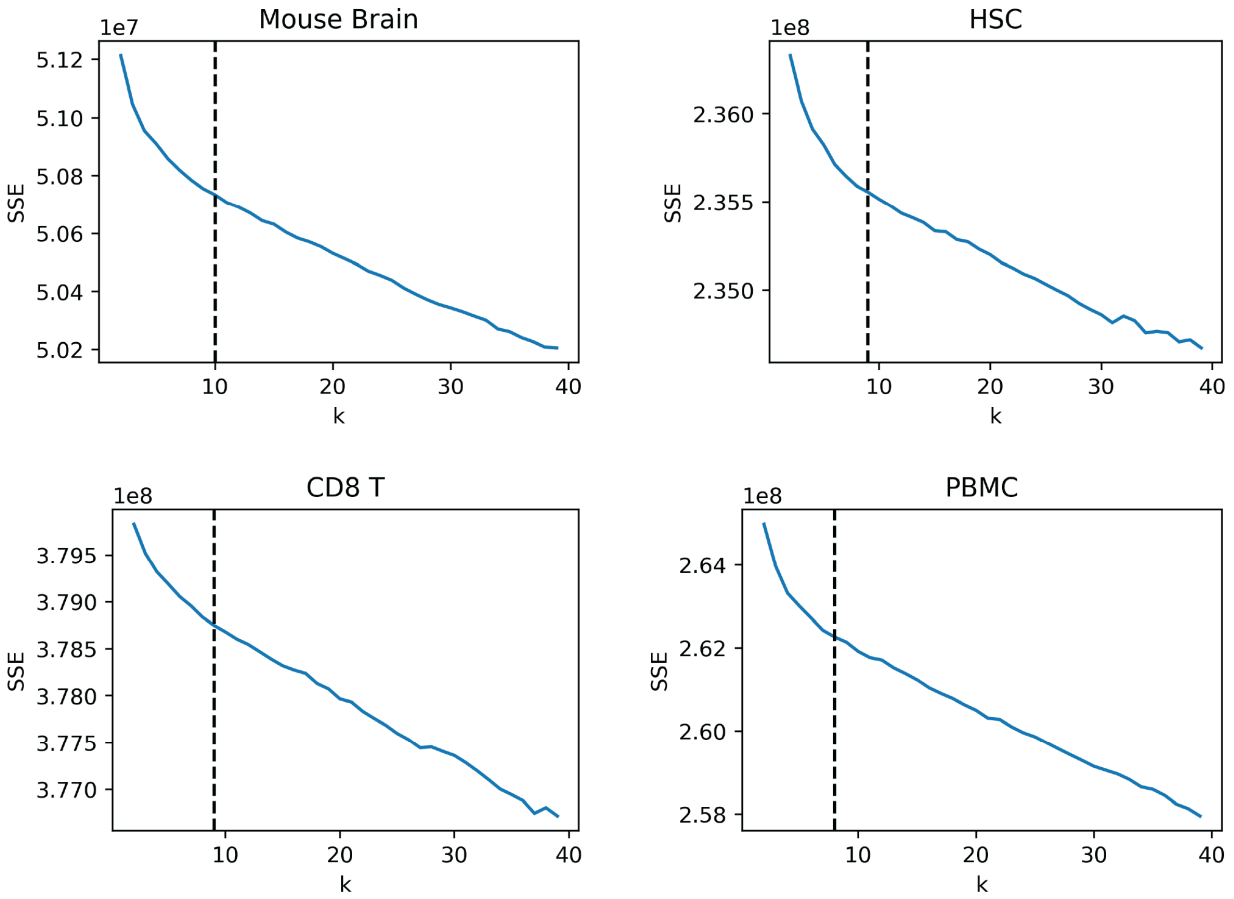


**Supplementary Figure 5. Additional details of the analysis of  $Ccr6^+$  and  $Ccr6^-$  CD8 T cell progenitor subsets.**

(A) Violin plots show the distribution of normalized *Ccr6* expression (log<sub>1p</sub> transformed) across the identified subpopulations in both the LCMV Armstrong and LCMV Clone 13 samples. Differential gene expression analysis (Wilcoxon test, adjusted for multiple hypothesis testing) was performed comparing cells from each subpopulation against all other cells in each of the two samples. Significant *Ccr6* overexpression was observed exclusively in the  $Ccr6^+$  progenitors (\*\*\*,  $p < 0.001$ ); comparisons for other subpopulations did not reach statistical significance ( $p > 0.05$ ).

(B) Violin plots of gene scores for genes differentially expressed in the  $CD62L^+$  and  $CD62L^-$  precursor exhausted T ( $T_{PEX}$ ) cell clusters reported in a previous study (Tsui et al. Nature 2022, PMID: 35978192), separately for our cell clusters in LCMV Armstrong (top) and clone 13 (bottom). The  $CD62L^+$   $T_{PEX}$  signature is enriched in our  $Ccr6^-$  progenitors, and  $CD62L^-$   $T_{PEX}$  signature is enriched in our  $Ccr6^+$  progenitors. \*\*,  $p < 0.01$ ; \*\*\*,  $p < 0.001$ , one-sided Mann-Whitney U test.

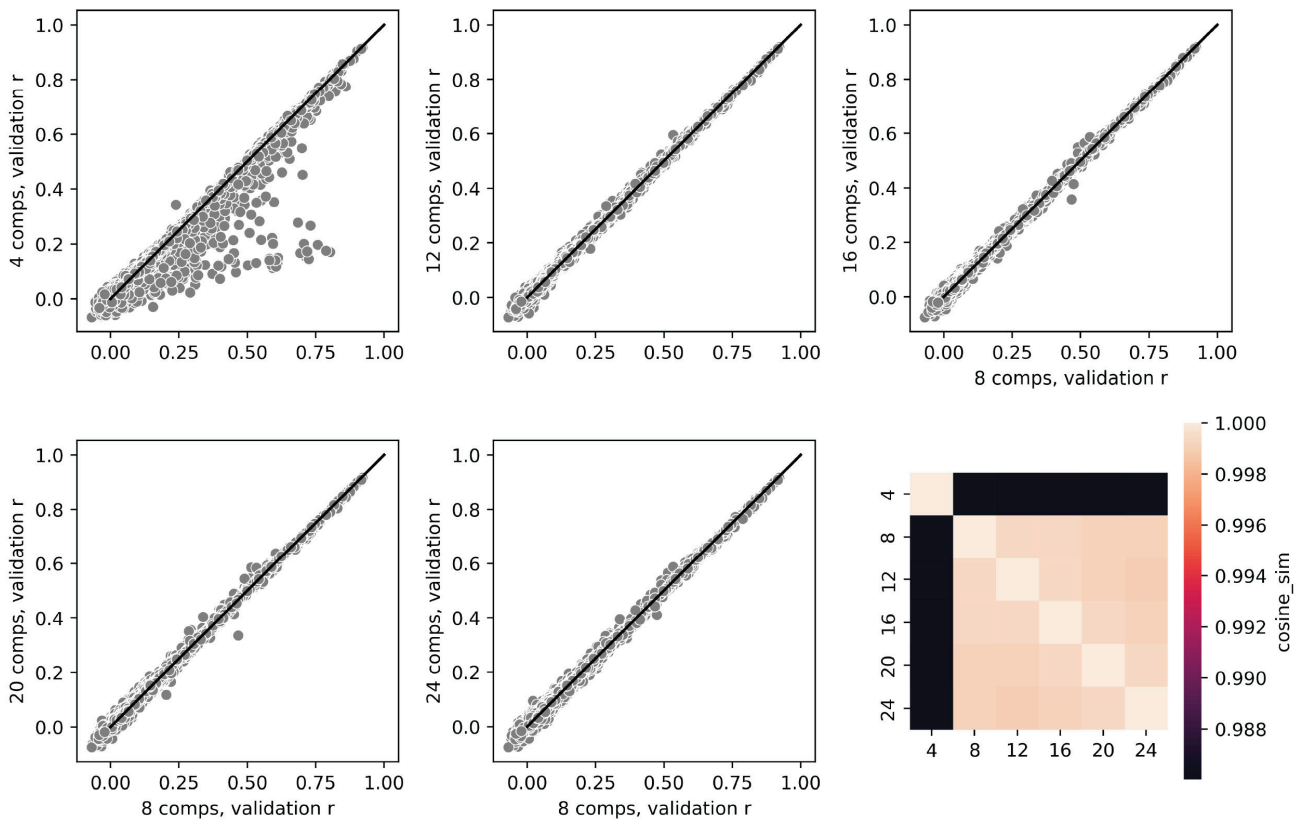
(C) Heatmaps of ArchVelo-inferred transition probabilities (Methods) between clusters including the  $Ccr6^+$  and  $Ccr6^-$  CD8 T cell progenitors, shown separately for LCMV Armstrong (left) and LCMV clone 13 (right).



**Supplementary Figure 6. Selection of the optimal number of archetypes using the elbow method.**

For each dataset used in this study, the plots show the sum of squared errors (SSE) of the archetypal reconstruction as a function of the number of archetypes,  $k$ . The vertical line indicates the optimal  $k$  based on the “elbow” of the SSE curve: PBMC,  $k = 8$ ; mouse embryonic brain:  $k = 10$ ; HSC,  $k = 9$ ; CD8 T cells,  $k = 9$ .

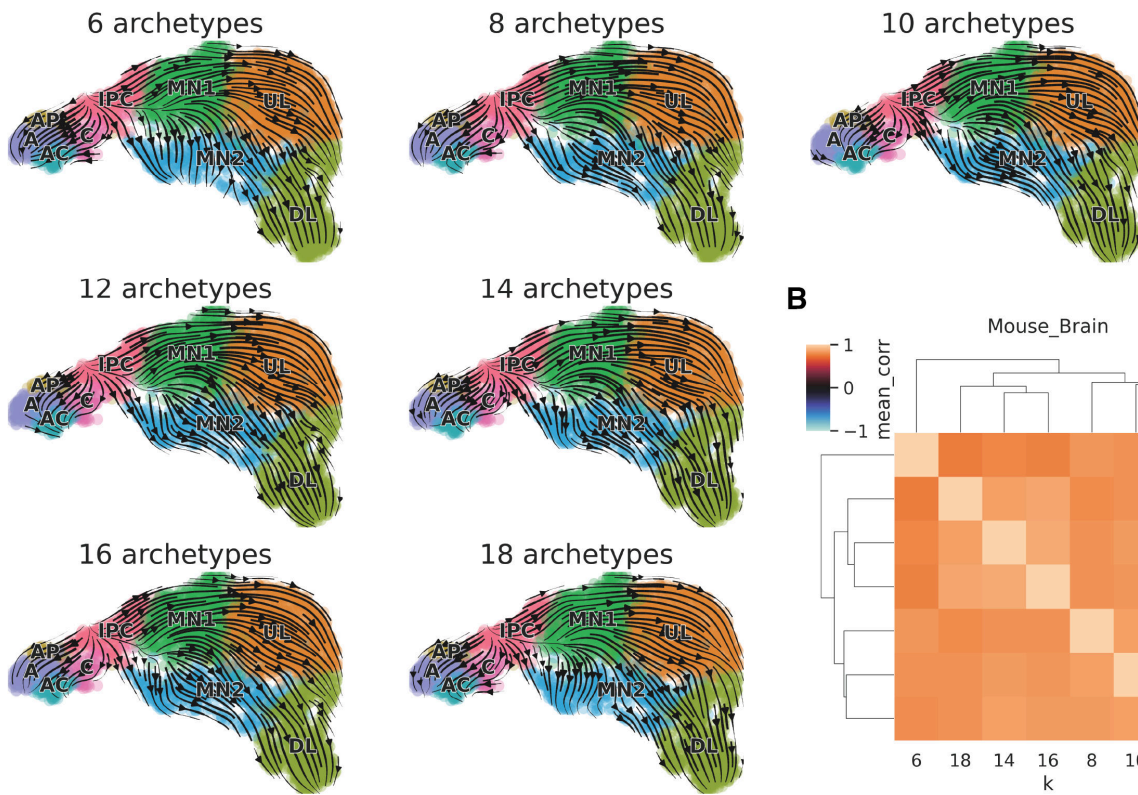
## Regression results, comparison over k



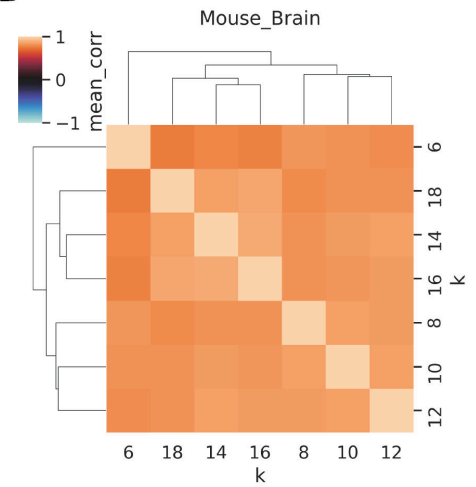
**Supplementary Figure 7. Robustness of scATAC-to-scRNA gene expression predictions to the number of chromatin accessibility archetypes.**

Archetypal analysis was performed on the scATAC-seq component of the PBMC multi-omics scATAC+RNA-seq dataset using varying numbers of archetypes  $k = 4, 8, 12, 16, 20, 24$ . These scATAC-seq archetypes were used as features in a ridge regression model to predict scRNA-seq gene expression in a cross-validation setting (as in **Fig. 1F,G**). The scatter plots show the mean gene-wise Pearson correlations  $r$  between predicted and observed expression over validation folds, comparing the baseline model ( $k = 8$ ,  $x$ -axis) against models with other  $k$  values ( $y$ -axis); black line shows the  $y = x$  diagonal. The heatmap shows cosine similarity between these results for different number of archetypes.

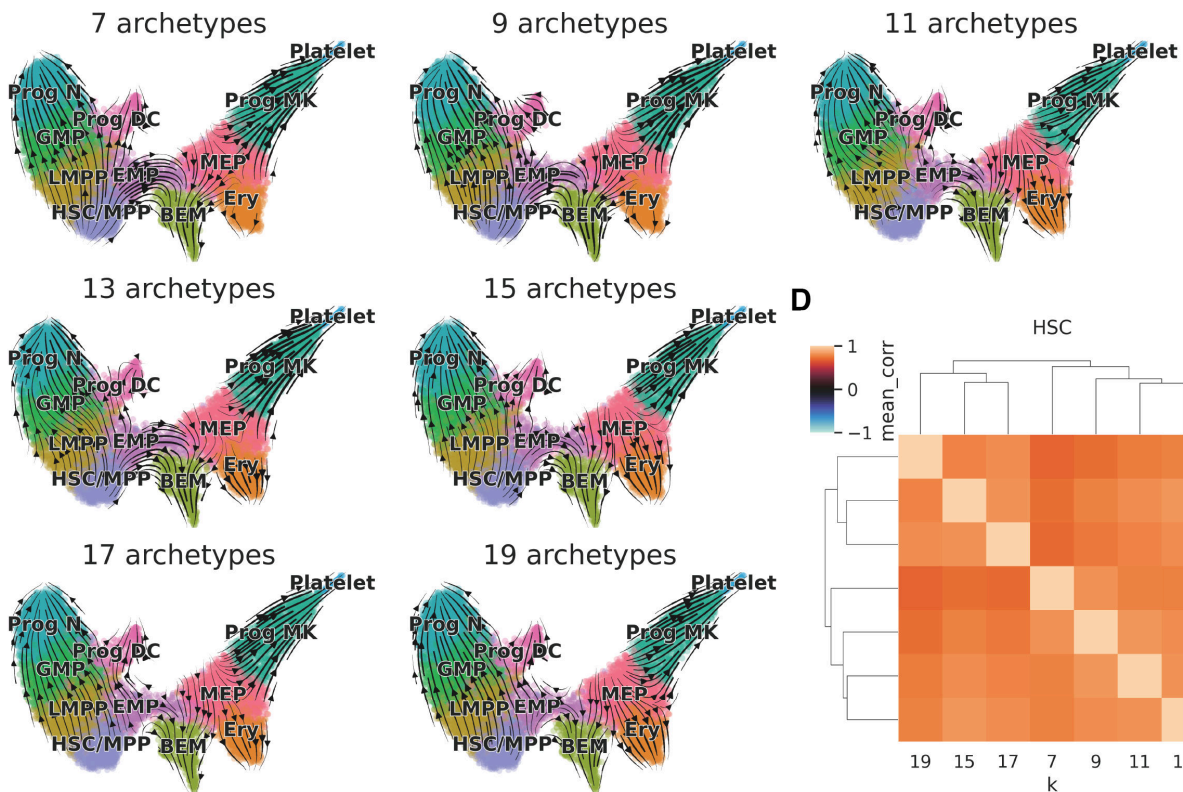
A



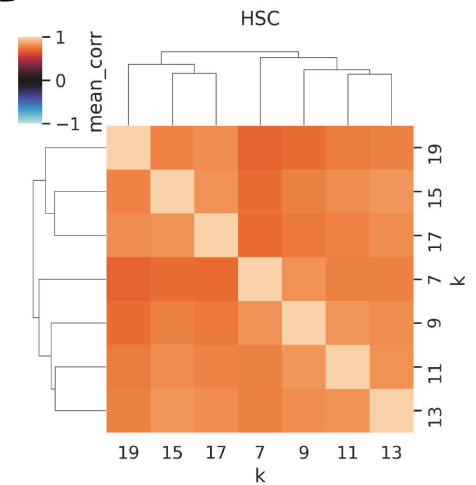
B



C



D



**Supplementary Figure 8. Robustness of the ArchVelo results to the choice of the number of archetypes.**

- (A) UMAP plots overlaid with ArchVelo velocity streams inferred using varying numbers of archetypes  $k = 6, 8, 10, 12, 14, 16, 18$  for the mouse embryonic brain dataset.
- (B) Heatmap showing correlations between velocities (mean over genes) for different values of  $k$  shown in panel A.
- (C) Same as panel A, for the HSC dataset and  $k = 7, 9, 11, 13, 15, 17, 19$ .
- (D) Same as panel B, for the HSC dataset.