
Enhancing RNA Velocity Inference with Single-Cell Aggregation Techniques

Juan P. Bernal-Tamayo

King Abdullah University of Science and Technology
Thuwal, Saudi Arabia
juan.bernaltamayo@kaust.edu.sa

Abstract

In this work, we propose a method for improving gene velocity inference by incorporating metacells, to gene dynamic inference by generative models. Our approach is based on the VeloVAE model, a generative model that learns a low-dimensional representation of the dynamic parameters and latent that govern gene expression changes. We preprocess the data using SEACells to obtain the corresponding metacells for each cell, and then introduce the metacells as an additional input to the model by modifying the encoder architecture to be conditioned by the metacell. Our experiments on the Pancreatic Endocrinogenesis dataset demonstrate (hopefully) that our method improves the accuracy of gene velocity prediction.

1 Introduction

1.1 Cell differentiation

Proteins are the building blocks of cells and give them their unique properties and functions. The information to make these proteins is encoded in the DNA of every cell. However, only certain genes are expressed in each cell, depending on which parts of the DNA are accessible for transcription into RNA molecules. To produce a protein, the DNA is first transcribed into RNA, the RNA molecule must then be spliced to remove regions that don't contain any information, called introns, and retain only the exons, which contain the information of the genes. This RNA molecule is then transported to the ribosomes, where it is translated into a protein. The process of protein translation, depicted in figure 1 is fundamental to cellular function and plays a critical role in many biological processes. Understanding the levels of proteins in each cell is crucial for studying various biological processes, such as cellular differentiation or cancer development. Nonetheless, counting the number of proteins in a single cell or a group of cells is challenging due to the complex structures of proteins. Therefore, counting the RNA molecules that are subsequently translated into proteins serves as a proxy to estimate the number of proteins present in each cell, enabling a better understanding of the cells. Technologies such as single-cell RNA sequencing are widely used to obtain such information.

1.2 Single cell RNA sequencing (scRNA-Seq)

Single-cell RNA sequencing (scRNA-seq) technology Haque et al. [2017] has become the state-of-the-art approach counting RNA transcripts within individual cells, revealing the composition of different cell types and their specific functions. Since its first discovery in 2009, studies based on scRNA-seq provide massive information across different fields making exciting new discoveries in better understanding the composition and interaction of cells. By capturing the messenger RNA of individual cells, scRNA-seq provides high-resolution insights into cellular diversity, and gene expression regulation. This technology has also been used in building comprehensive cellular atlases

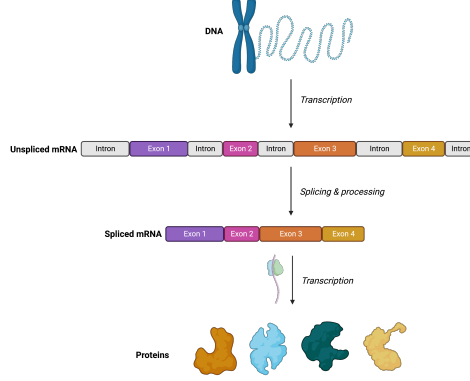


Figure 1: Original VeloVAE architecture.

across different organisms, which have revolutionized our understanding of cell types and functions in health and disease.

1.3 RNA velocity

In principle, because of the destructive nature of scRNA-Seq, each cell can only be sequenced once and any information on their future state is virtually unavailable. Manno et al. [2018] proposed Velocyto as a mathematical model 1, as well as a method for fitting the parameters of the model to find the rate of change in the counts of each gene and a timestamp for each cell, allowing estimation of gene velocity and pseudo-time, enhancing downstream analyses of scRNA-Seq experiments. The dynamic parameters' information and the solution of the system of equations can be used to extrapolate the model and predict the fate of a system of cells, as well as to interpolate the model to tweak certain parameters and favor certain cellular fates over others.

$$\begin{aligned}\frac{du(t)}{dt} &= \alpha - \beta u(t) \\ \frac{ds(t)}{dt} &= \beta u(t) - \gamma s(t)\end{aligned}\tag{1}$$

Where $u(t)$ and $s(t)$ are the respective unspliced and spliced abundances of a gene. And the dynamic parameters are, α , the transcription rate, β , the splicing rate, and γ , the degradation rate. Starting from an initial condition (u_0, s_0) , this system of equations can be solved analytically, obtaining

$$\begin{aligned}u(t) &= u_0 e^{-\beta t} + \frac{\alpha}{\beta} (1 - e^{-\beta t}) \\ s(t) &= s_0 e^{-\gamma t} + \frac{\alpha}{\gamma} (1 - e^{-\gamma t}) + \frac{\alpha - \beta u_0}{\gamma - \beta} (e^{-\gamma t} - e^{-\beta t})\end{aligned}\tag{2}$$

1.4 Metacells

Metacells are groups of cells derived from single-cell sequencing data that represent distinct and granular cell states. They are useful for identifying rare cell types, detecting outlier cells, and revealing gene expression dynamics during cell differentiation. Metacells can be identified through algorithms such as MetaCell by Baran et al. [2019], Metacell-2 by Ben-Kiki et al. [2022], and SEACells by Persad et al., which allow efficient decomposition of scRNA-seq datasets into small and cohesive groups of cells.

2 Problem formulation

Accurate estimation of gene velocity parameters can be challenging, especially for large datasets. Recently, generative models such as VeloVAE by Qiao and Huang and VeloVI by Gayoso et al. have been proposed as a way to fit the dynamic parameters of the gene velocity model, but their

performance can be further improved. In this work, we propose to preprocess scRNA-seq datasets using SEACells to obtain metacells, which represent highly granular and distinct cell states. We then introduce an extra input layer that takes both the cell and its corresponding metacell as input, conditioning the encoder part of the model with the additional metacell information. Our hypothesis is that this additional information will provide the encoder with more knowledge in the generation of the latent distribution, resulting in improved dynamic parameters, thus, improved gene velocity estimation.

3 Model

In this section, we introduce a new implementation of the gene velocity model using a generative model approach called VeloVAE. Our aim is to leverage the additional information provided by metacells, obtained from the SEACells algorithm, to improve the quality of the latent space representation. To achieve this, we propose an extra input layer for the metacells and an encoder layer that creates an embedding of the metacells, which is then concatenated to the two existing encoder layers. We will present the architecture of the model and the loss function, which remains the same as in the original VeloVAE since we are only adding an additional encoding layer and not new outputs or latent space variables.

3.1 Architecture

3.1.1 Original VeloVAE architecture

The original architecture, depicted in figure 2, takes a dataset $D \in \mathbb{R}^{n_c \times 2n_g}$, where n_c and n_g are the number of cells and genes respectively. The factor of 2 in the number of genes comes from the fact that both the spliced and unspliced counts of the genes are used. The input layer has a size of $2n_g$. There are two encoding layers with sizes $N_1 = 500$ and $N_2 = 250$. Each layer consists of a Linear layer followed by a batch normalization layer and a dropout layer with a dropout probability of $p = 0.2$. To obtain the latent space, the output of the second layer is mapped onto the corresponding latent dimensions. The latent variable for the dynamic parameter has a dimension of $\dim(c) = 5$, and the latent time variable has a dimension of $\dim(t) = 1$. The decoder takes the latent spatial variable and decodes it using two layers symmetrical to those of the encoder. The difference is that the output dimension is n_g instead of $2n_g$ because only one transcription rate, ρ_i , per gene is obtained. The splicing rate β and the degradation rate γ are trainable parameters of the decoder, of size n_g . With the time and all the parameters, the ODE system (equation 1) is constructed and solved analytically (equation 2). From this solution, and the parameters $\theta = (\rho, \beta, \gamma)$ we obtain the function $F(t; \theta)$ from which we can reconstruct the spliced, $\hat{s}(t)$, and unspliced, $\hat{u}(t)$ counts.

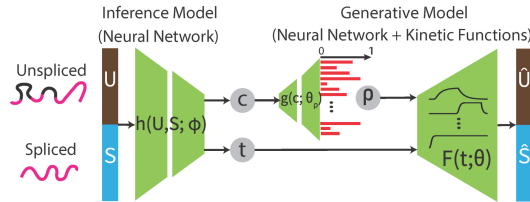


Figure 2: Original VeloVAE architecture.

3.1.2 Modified architecture: original idea

The proposed architecture of our model, shown in figure 3 takes the VeloVAE architecture and adds an additional $2n_g$ input layer for the spliced and unspliced counts of the metacells, then using two linear layers of dimension $N_1 = 500$ and $N_2 = 250$, embeds the metacells into spaces of the same dimension of the intermediate layers of the encoder, then the metacells are concatenated to the corresponding cells, before being passed to the batch normalization and dropout layers, which have to be enlarged to accept double the number of dimensions.

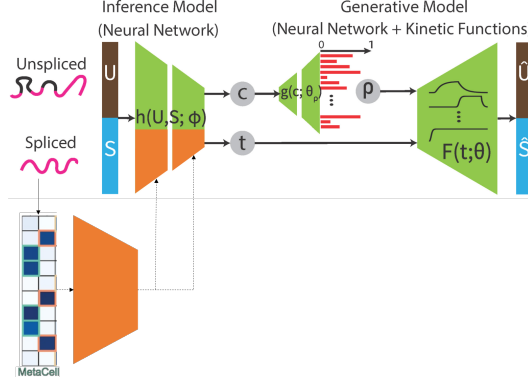


Figure 3: Proposed modified architecture.

3.1.3 Modified architecture: sampling the model

Both the original VeloVAE and the proposed modified architectures suffer from a problem: they are not suitable for sampling a new cell from the latent distribution. This is because each of the latent variables, i.e., the cell variable c and the time variable t , are sampled independently. Consequently, the sampled latent variable could correspond to a cell at the end of the differentiation process while the sampled time could correspond to an earlier cell. To address this issue, we further modified the architecture by introducing a conditioning latent variable z , the modified architecture is shown in figure 4. This variable is passed forward through independent linear layers to obtain the latent variable c and the latent time t . With this modification, we can generate a cell by first sampling from the distribution of z , then passing this sample forward to obtain the parameters of the distribution of c and t , sampling independently from these distributions, and finally generating the cell using the decoder.

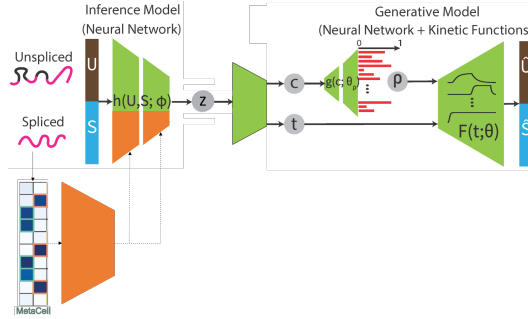


Figure 4: Architecture modified to include the conditioning variable.

3.2 Loss function

The loss function is identical to the one used in VeloVAE, the ELBO. Since the distributions are assumed to be normal, this can be decomposed into the reconstruction loss and the KL divergence between the posterior and the prior distributions

$$ELBO = \sum_{i=1}^{n_g} E_{q(c,t|x_i)} [\log p(x_i|c,t)] - KL(q(c,t|x_i)||p(c,t))$$

Assuming independence in the latent variables we obtain

$$ELBO = \sum_{i=1}^{n_g} E_{q(c,t|x_i)} [\log p(x_i|c,t)] - KL(q(c|x_i)||p(c)) - KL(q(t|x_i)||p(t))$$

3.2.1 Conditioning latent variable

After adding the conditioning latent variable, the new *ELBO* becomes

$$ELBO = \sum_{i=1}^{n_g} E_{q(c,t|x_i)} [\log p(x_i|z, c, t)] - KL(q(z|x_i)||p(z)) - KL(q(c|x_i, z_i)||p(c)) - KL(q(t|x_i, z_i)||p(t))$$

Where we assume that given a sample of the conditioning latent variable z_i , the dynamic latent variable c , and the latent time are independent.

3.3 Prior distributions

4 Datasets

4.1 Pancreatic endocrinogenesis

The pancreatic endocrinogenesis generated by Bastidas-Ponce et al. [2019] dataset has been extensively used for velocity analysis due to its well-defined cell types and the understanding of the transitions between them. This dataset contains approximately 3700 cells. After filtering and pre-processing, the number of genes has been reduced to around 2000 to 3000 from the original 23000 sequenced genes. In figure 5a a UMAP projection of the dataset is depicted with the different cell types.

4.2 Zebrafish pigment development

The cells on this dataset were collected during an active period of post-embryonic development, when many of these cell types are migrating and differentiating as the animal transitions into its adult form. This study also explores the role of thyroid hormone, on the development of these different cell types. Such developmental and other dynamical processes are especially suitable for dynamical analysis. This dataset contains approximately 4200 cells. After filtering and preprocessing, the number of genes has been reduced to around 3000 from the original 16900 sequenced genes. In figure 5b a UMAP projection of the dataset is depicted with the different cell types.

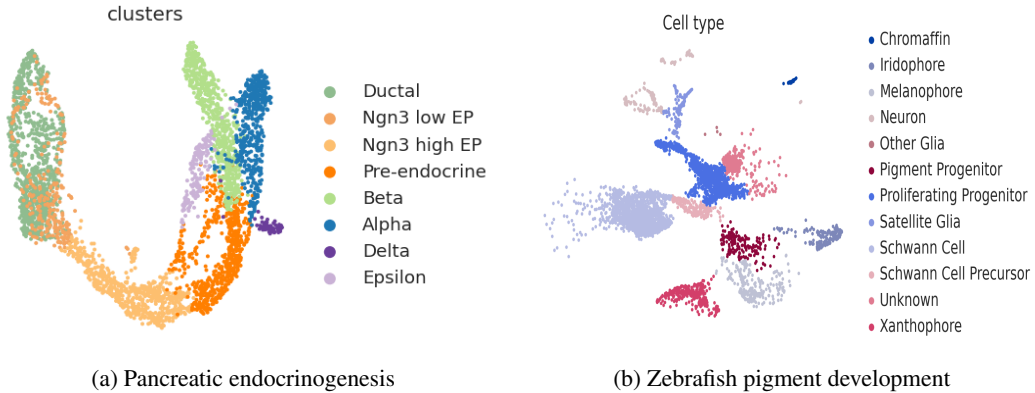


Figure 5: UMAP projection of the datasets classified by cell type.

5 Evaluation

5.1 In-cluster coherence

This is an evaluation metric proposed in VeloVAE evaluates the coherence of velocity direction for neighboring cells within a cluster or cell-type.

A k-NN graph is constructed among all cells using $k = 30$. For each cell in a cluster, the average cosine similarity between the inferred velocity vectors is computed. This is done for all the neighbors

that also belong to the same cluster, resulting in the in-cluster coherence measure. The in-cluster coherence for cluster K is defined as

$$ICC(K) = \frac{1}{|K|} \sum_{x \in K} \frac{1}{|\mathcal{N}(x) \cap K|} \sum_{y \in \mathcal{N}(x) \cap K} S_{\cos}(v_x, v_y). \quad (3)$$

Where $S_{\cos}(v_x, v_y)$ is the cosine similarity between the velocities, v_x and v_y .

$$S_{\cos}(x - y, v_x) = \cos(\angle(v_x, v_y)) = \frac{(v_x) \cdot v_y}{\|v_x\| \|v_y\|}$$

5.2 Cross-boundary correctness

It measures how accurately the model captures the patterns of cell development across multiple cell types that are known to be related to each other. This metric evaluates the ability of the model to capture the biological processes underlying cell differentiation. It measures the alignment of cell velocities within a cluster with the vector pointing towards the descendant cell types. Essentially, it indicates how well the model is able to capture the direction of cell differentiation. The Cross-boundary correctness for the transition from cell-type K_1 to cell-type K_2 is defined as:

$$CBC(K_1 \rightarrow K_2) = \frac{1}{|K_1|} \sum_{x \in K_1} \frac{1}{|\mathcal{N}(x) \cap K_2|} \sum_{y \in \mathcal{N}(x) \cap K_2} S_{\cos}(y - x, v_x). \quad (4)$$

5.3 Time score

The time score is a metric that measures the average probability of the cells in a progenitor cell type to have an earlier time than its corresponding neighbors in the descendant cell type. It provides a way to evaluate how well the model is able to capture the temporal ordering of the differentiation process. The time score for the transition from cell-type K_1 to cell-type K_2 is defined as:

$$T_{score}(K_1 \rightarrow K_2) = \frac{1}{|K_1|} \sum_{x \in K_1} \sum_{y \in \mathcal{N}(x) \cap K_2} P(t(x) < t(y)) \quad (5)$$

This probability can be calculated as

$$P(X \leq 0) = \int_{-\infty}^0 p(x) dx$$

Where

$$X \sim \mathcal{N}(t(y) - t(x), \sqrt{std(x)^2 + std(y)^2})$$

Where, for each cell x the time $t(x)$ and standard deviation $std(x)$ are the parameters of the latent time distribution.

5.4 Other metrics

Other several metrics commonly used to evaluate how well generative models can reconstruct samples I used to compare the models are:

- **Mean absolute error**

$$\frac{1}{N} \sum_{i=1}^N |\hat{u}_i - u_i| + |\hat{s}_i - s_i|$$

- **Mean square error**

$$\frac{1}{N} \sum_{i=1}^N (\hat{u}_i - u_i)^2 + (\hat{s}_i - s_i)^2$$

- **log-likelihood**

$$\frac{1}{N} \sum_{i=1}^N \sum_{j=1}^{n_g} -\frac{(\hat{u}_i^j - u_i^j)^2}{2\sigma_{u^j}^2} - \frac{(\hat{s}_i^j - s_i^j)^2}{\sigma_{s^j}^2} - \log(\sigma_{u^i}) - \log(\sigma_{s^j}) - \log(2\pi)$$

Where \hat{u}_i and \hat{s}_i are the respective reconstructed spliced and unspliced counts of cell $x_i = (u_i, s_i)$, (u_i^j, s_i^j) is the expression of the gene j in cell i , and σ_{u^j} and σ_{s^j} are parameters of the VAE corresponding to the standard deviation of each gene in the data.

6 Results

6.1 Metrics

After modifying the entire VeloVAE architecture to include the metacells and the conditioning latent variable, we implemented the new metrics and evaluated the complete model. The average results for all metrics in both datasets are displayed in table 1. Moreover, table 2 shows the results for the Cross-boundary Correctness and Time score metrics for the pancreatic endocrinogenesis dataset, while table 3 shows the corresponding results for the zebrafish pigment development regarding the cell-type transitions.

Table 1: Comparison of VeloVAE and VeloVAE+SEACells on the Pancreatic and Zebrafish datasets

Metric	Pancreas		Zebrafish	
	VeloVAE	VeloVAE + Metacells	VeloVAE	VeloVAE + Metacells
In-cluster coherence	0.725	0.811	0.832	0.945
Cross-boundary correctness	0.118	0.237	0.172	0.174
Time score	0.372	0.479	0.411	0.601
LL test	2657.383	2498.271	1404.154	1059.517
LL train	2685.372	2515.599	1473.136	1022.742
MAE test	0.142	0.152	0.178	0.237
MAE train	0.140	0.152	0.173	0.245
MSE test	0.236	0.221	0.227	0.467
MSE train	0.240	0.222	0.217	0.527
Training time	468 s	878 s	710 s	716. s

Table 2: Comparison of Cross-boundary correctness and Time score

Transition	Cross-boundary correctness		Time score	
	VeloVAE	VeloVAE + Metacells	VeloVAE	VeloVAE + Metacells
Ductal → Ngn3 low EP	0.051	-0.129	0.563	0.575
Ngn3 low EP → Ngn3 high EP	0.819	0.771	0.890	0.917
Ngn3 high EP → Pre-endocrine	0.603	0.773	0.765	0.923
Pre-endocrine → Alpha	0.356	0.485	0.589	0.797
Pre-endocrine → Beta	0.612	0.624	0.487	0.700
Pre-endocrine → Delta	0.407	0.450	0.926	0.270
Pre-endocrine → Epsilon	0.611	0.720	0.875	0.618

From tables 2 and 3, we can observe that while the cross-boundary correctness results were not consistently better with the modified VeloVAE architecture for the zebrafish dataset, the Time score showed a general improvement for both datasets. This suggests that the modified architecture may not only improve the accuracy of cell-type transitions but also increase the latent time inference accuracy, which is important in capturing the dynamics of biological processes. Therefore, the modified VeloVAE architecture may have potential for further applications in analyzing developmental processes in single-cell transcriptomics data.

Table 3: Comparison of Cross-boundary correctness and Time score

Transition	Cross-boundary correctness		Time score	
	VeloVAE	VeloVAE + Metacells	VeloVAE	VeloVAE + Metacells
Prolif. Prog. → Schwann Cell Prec.	0.212	0.149	0.226	0.762
Prolif. Prog. → Pigment Progenitor	0.192	0.177	0.500	0.696
Prolif. Prog. → Satellite Glia	-5.564×10^{-3}	0.219	5.517×10^{-1}	0.810
Schwann Cell Prec. → Schwann Cell	0.460	0.208	0.184	0.587
Satellite Glia → Neuron	0.112	0.074	0.497	0.648
Pigment Prog. → Melanophore	0.161	0.257	0.497	0.594
Pigment Prog. → Iridophore	0.169	0.268	0.576	0.293
Pigment Prog. → Xanthophore	0.073	0.037	0.333	0.417

6.2 Qualitative behavior

Besides the evaluation results we can observe the comparison between the velocity fields (figures 6 and 7) of each of the models in the UMAP projection. For the pancreatic endocrinogenesis dataset the main observation we can make is regarding the transitions from pre-endocrine to all the endocrine cell types. With the original VeloVAE architecture (figure 6a), the velocity flow crosses different endocrine cells to arrive at other endocrine cells. For example, we can observe a false transition from Beta to Alpha. In the case of the modified architecture (figure 6b), all the transitions are strictly from pre-endocrine to each of the endocrine cell types. This result is corroborated by the increase in the cross-boundary correctness score in all the pre-endocrine to endocrine transitions. Another interesting result that can be observed in the velocities is the difference in velocity flow between the Ductal/Ngn3 low EP region for the original VeloVAE model (figure 6a) and the modified MetaVAE architecture (figure 6b). In the original model, the velocity flow goes directly from the Ductal/Ngn3 low EP region to the Ngn3 high EP cell type, while in the modified architecture, the velocity swirls around before transitioning. This behavior is more biologically plausible because the Ductal and Ngn3 low EP cell types have similar gene expression distributions and are difficult to separate in the UMAP. The transition occurs from Ductal to Ngn3 low EP, creating a velocity field that seems to point inwards. This also explains why both models have a very low Cross-boundary correctness score in this transition.

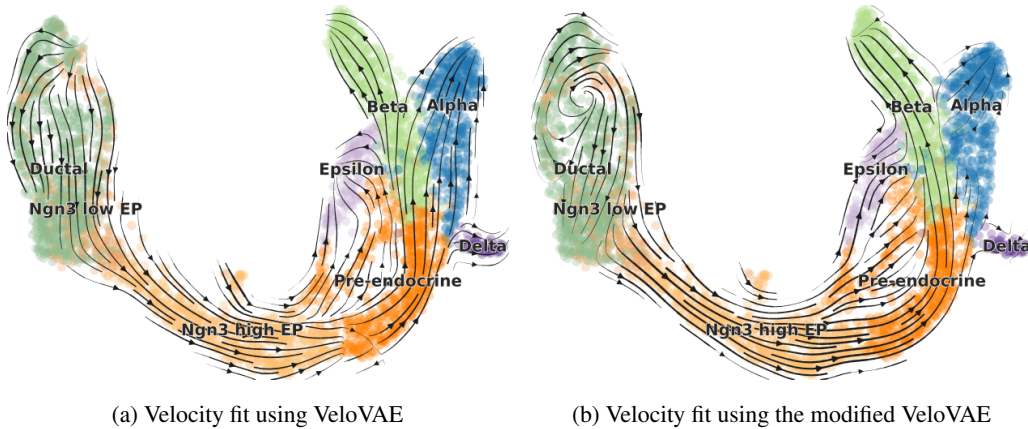


Figure 6: Comparison between velocity fit for pancreatic endocrinogenesis

On the other hand, for the zebrafish pigment dataset the cross-boundary correctness scores did not improve as much, and this can also be noticed in the velocity embedding in figure 7.

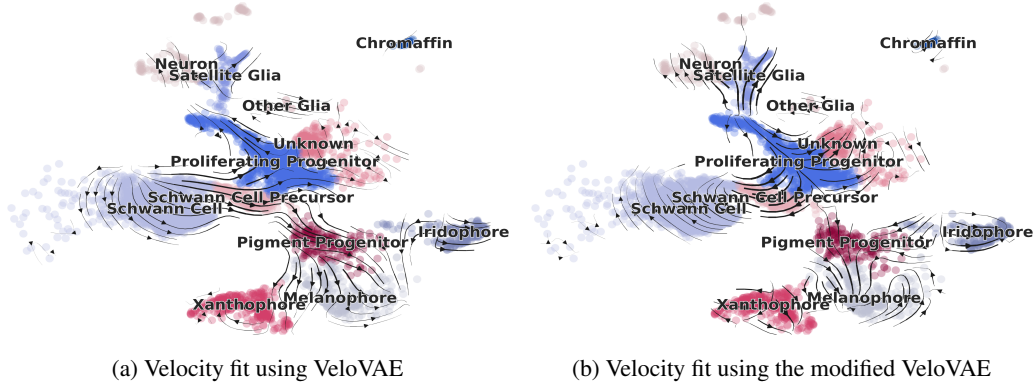


Figure 7: Comparison between latent times for zebrafish pigment development

6.3 Generative properties

To evaluate the generative properties of our model, I randomly sampled 25% of the data points and passed them through the model. However, instead of generating 1 sample per datapoint, I generated 4 samples for each point. In figures 8 and 9 we can see a UMAP embedding of the original distribution of the data in blue, and in orange or green the embedding of the corresponding sampled distribution. Unfortunately, the results were not very good. I repeated the experiment with 100% of the data points. However, even with the entire dataset, the models were not able to match the distribution of the original datapoints. Although my model achieved better scores with the evaluation metrics, it still performed worse in reproducing the original distribution, most likely because of the extra latent variable causing the sampling to stay around the metacells.

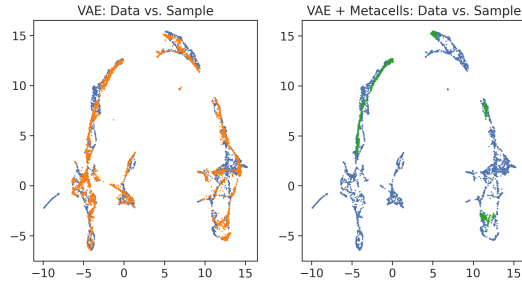


Figure 8: Cell sampling by VeloVAE (left) and Modified VeloVAE (right) for the pancreatic endocrinogenesis dataset

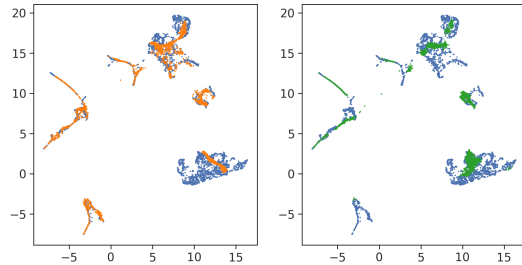


Figure 9: Cell sampling by VeloVAE (left) and Modified VeloVAE (right) for the zebrafish pigment development dataset

To gain insights into why this behavior was occurring, I generated data using the metacells as inputs from both the original data and my model. The results shed light on the problem. Figures 10 and 11 demonstrate that when the metacells are reconstructed from sampling, they are constrained to a

limited part of the distribution, rather than being spread across the entire distribution as expected. This suggests that the modified VeloVAE model's sampling problem may be caused by over-fitting around the metacells that constrains the sampling to a smaller region. The extra sampling step may be causing the over-fitting, and since this new latent variable is not necessary to guarantee dependence on the sampled latent state and latent time variables, a good starting point for future research is to remove this extra latent variable.

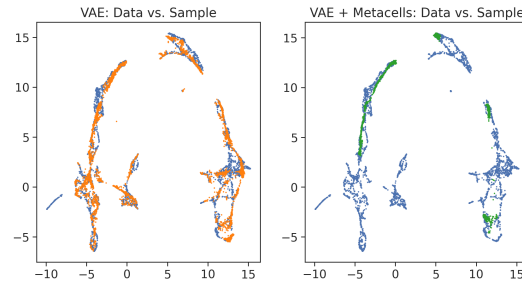


Figure 10: Cell sampling from metacells by VeloVAE (left) and Modified VeloVAE (right) for the pancreatic endocrinogenesis dataset

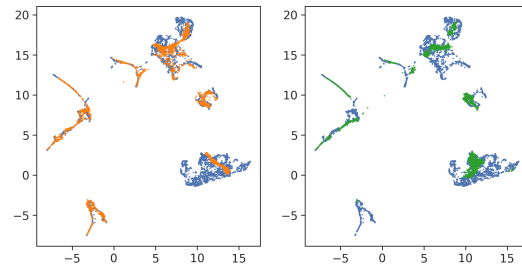


Figure 11: Cell sampling from metacells by VeloVAE (left) and Modified VeloVAE (right) for the zebrafish pigment development dataset

References

- Yael Baran, Akhiad Bercovich, Arnau Sebe-Pedros, Yaniv Lubling, Amir Giladi, Elad Chomsky, Zohar Meir, Michael Hoichman, Aviezer Lifshitz, and Amos Tanay. Metacell: Analysis of single-cell rna-seq data using k-nn graph partitions. *Genome Biology*, 20, 10 2019. ISSN 1474760X. doi: 10.1186/s13059-019-1812-2.
- Aimée Bastidas-Ponce, Sophie Tritschler, Leander Dony, Katharina Scheibner, Marta Tarquis-Medina, Ciro Salinno, Silvia Schirge, Ingo Bartscher, Anika Böttcher, Fabian Theis, Heiko Lickert, and Mostafa Bakhti. Massive single-cell mrna profiling reveals a detailed roadmap for pancreatic endocrinogenesis. *Development*, 1 2019. ISSN 1477-9129. doi: 10.1242/dev.173849.
- Oren Ben-Kiki, Akhiad Bercovich, Aviezer Lifshitz, and Amos Tanay. Metacell-2: a divide-and-conquer metacell algorithm for scalable scrna-seq analysis. *Genome Biology*, 23, 12 2022. ISSN 1474760X. doi: 10.1186/s13059-022-02667-1.
- Adam Gayoso, Philipp Weiler, Mohammad Lotfollahi, Dominik Klein, Justin Hong, Aaron Streets, Fabian J Theis, Nir Yosef, Chan Zuckerberg Biohub, and San Francisco. Deep generative modeling of transcriptional dynamics for rna velocity analysis in single cells. doi: 10.1101/2022.08.12.503709.
- Ashraf Haque, Jessica Engel, Sarah A. Teichmann, and Tapio Lönnberg. A practical guide to single-cell rna-sequencing for biomedical research and clinical applications, 8 2017. ISSN 1756994X.
- Gioele La Manno, Ruslan Soldatov, Amit Zeisel, Emelie Braun, Hannah Hochgerner, Viktor Petukhov, Katja Lidschreiber, Maria E Kastri, Peter Lönnberg, Alessandro Furlan, Jean Fan, Lars E Borm,

- Zehua Liu, David Van Bruggen, Jimin Guo, Xiaoling He, Roger Barker, Erik Sundström, Gonçalo Castelo-Branco, Patrick Cramer, Igor Adameyko, Sten Linnarsson, and Peter V Kharchenko. Rna velocity of single cells. *Nature*, 560:494–498, 2018. ISSN 0028-0836. doi: 10.1038/s41586-018-0414-6. URL <https://dx.doi.org/10.1038/s41586-018-0414-6>.
- Sitara Persad, Zi-Ning Choo, Christine Dien, Ignas Masilionis, Ronan Chaligné, Tal Nawy, Chrysothemis C Brown, Itsik Pe’er, Manu Setty, and Dana Pe’er. Seacells: Inference of transcriptional and epigenomic cellular states from single-cell genomics data. doi: 10.1101/2022.04.02.486748. URL <https://doi.org/10.1101/2022.04.02.486748>.
- Chen Qiao and Yuanhua Huang. Representation learning of rna velocity reveals robust cell transitions. *BIOPHYSICS AND COMPUTATIONAL BIOLOGY COMPUTER SCIENCES*. doi: 10.1073/pnas.2105859118/-/DCSupplemental. URL <https://doi.org/10.1073/pnas.2105859118>.