

Supplement: ‘Dhaka: Variational Autoencoder for Unmasking Tumor Heterogeneity from Single Cell Genomic Data’

Sabrina Rashid,¹ Sohrab Shah^{2,3,4}, Ziv Bar-Joseph⁵, Ravi Pandya^{6*}

¹ Computational Biology Department, Carnegie Mellon University, Pittsburgh, USA

² Department of Computer Science, University of British Columbia, Vancouver, Canada

³ Department of Molecular Oncology, BC Cancer Agency, Vancouver, Canada

⁴ Department of Pathology and Laboratory Medicine, University of British Columbia, Vancouver, Canada

⁵ Machine Learning Department and Computational Biology Department, Carnegie Mellon University, Pittsburgh, USA

⁶ Microsoft Research, Redmond, USA

*To whom correspondence should be addressed; E-mail: ravip@microsoft.com

1 Appendix

1.1 Software

We have developed a python package for the Dhaka variational autoencoder using the Keras module [1]. The package is released as open source (<https://github.com/MicrosoftGenomics/Dhaka>). Since this is a probabilistic encoding of the genomic data, often we need to do multiple warm starts of the encoder to select the best encoding. For example, if we are interested in identifying clusters, from each projected encoding we will compute the silhouette score, and select the encoding that maximizes the score [2]. We have used multiple warm starts only for the synthetic data analysis. We did not use multiple warm starts for the copy number and gene expression data. The number of warm starts is a user parameter for the package (5 in case of synthetic dataset). The Dhaka package can also perform gene selection, if needed. We have three options for selecting informative genes for analysis.

- Coefficient of variation (CV) score: CV of gene i with expression profile $g_i \in R^{1 \times m}$ is defined as $CV_i = std(g_i)/mean(g_i)$. Here m is the total number of cells.
- Entropy En : $En_i = -sum(p_i \log_2(p_i))$. Here p_i is the estimated histogram from g_i .
- Average expression value \bar{A} : This is simply the average expression value of a particular gene across all cells.

The gene selection criteria and number of genes to be included in the analysis are both user parameters. We have used gene selection for the three RNA-Seq gene expression datasets, (5000 genes with \bar{A} criteria). Dhaka is robust to the drop out events, therefore we did not have to model the drop out events separately. The other parameters of the package are the number of the latent dimensions, learning rate, batch size, number of epochs, and clip norm of the gradient ¹. We have used the following values for these parameters, learning rate = 0.0001, batch size =50, clip norm =2, and number of epochs =5. We used the same parameter values for all the datasets analyzed in the paper. A sensitivity analysis on these parameters on the simulated dataset are summarized in Table S1. For each of the parameters we considered three possible values keeping the other parameter values fixed at the aforementioned values. We can see that the most sensitive parameters are the batch size and number of epochs. If we reduce the batch size to 20 we see a significant decrease in the performance, we can see similar effect in increasing the number of epochs in training. The learning rate and clip Norm parameters have relatively stable ARI scores.

¹ Gradients will be clipped when their L2 norm exceeds this value. This parameter is used for the stability of the gradient descent algorithm.

Learning rate	Batch size	Clip norm	Number of epochs	ARI
0.0001	50	2	5	.73
0.001	50	2	5	.71
0.01	50	2	5	.65
0.0001	100	2	5	.72
0.0001	20	2	5	.22
0.0001	50	1	5	.74
0.0001	50	.5	5	.65
0.0001	50	2	10	.67
0.0001	50	2	20	.26

Table S1. Sensitivity analysis of the Dhaka hyperparameters

The input to the method is log2 transformed TPM counts. While it is standard practice to scale the data between [0,1] for sigmoid output activation, we have not used it in our method. We tested raw expression counts, log2 transformed counts, and scaled counts [0,1] as inputs to our structure and observed that the best correlation scores were found when log2 transformed counts were fed directly to the network. The software also provides option to calculate relative expression profiles (section 1.4) to be used as inputs.

In addition to sigmoid activation, we have also included the option of ReLU output activation. We have shown the results with ReLU activation function in Fig. S1. As can be seen, for the Oligodendroglioma and Glioblastoma datasets we obtain similar correlation scores for ReLU and sigmoid activation (± 0.03). However for the melanoma and simulated datasets the scores for ReLU are lower (a difference > 0.05), which means that these correlations are similar to the ones we obtain for the methods we compared to. The choice of the output activation function (sigmoid or ReLU) is a user defined parameter.

1.2 Simulated data generation

To generate the simulated data we sampled expression counts $E_{i,j}$ for gene i and cell j from a negative binomial distribution parameterized by $NB(\lambda_i, \omega)$. The λ_i parameter is specific to each gene and learned from a gamma distribution parameterized by $\gamma(\eta_k, \zeta_k)$, where k stands for the index of the cluster the cell belongs too. The ω parameter is kept constant. By tuning ω we can control the amount of noise in the simulated RNA-seq counts. For each gene i , we randomly select a fraction of cells to be dropped out. The fraction, f_i is sampled from uniform distribution $u(0.2, 0.9)$. To insert completely noisy genes without any cluster specific expression profile, we randomly select a subset of genes in each cell and replace their cluster specific expression values with values sampled from uniform distribution parameterized by $u(0, \max(E_{i,j}))$.

1.3 Runtime comparison

We have compared the runtime of the proposed method and all the comparing methods for the simulated dataset. For each method the time to compute the three dimensional projection given fixed parameters is reported in Table S2. All the methods were performed on a 32 GB 3.4 GHz Windows machine. As can be seen from Table S2, although PCA, MAGIC, and NMF are faster than Dhaka, they have much lower ARI score.

1.4 Relative gene expression

The relative gene expression $Er_{i,j} = E_{i,j} - \text{mean}(E_{i,1,\dots,n})$. Here i and j correspond to gene and cell, respectively.

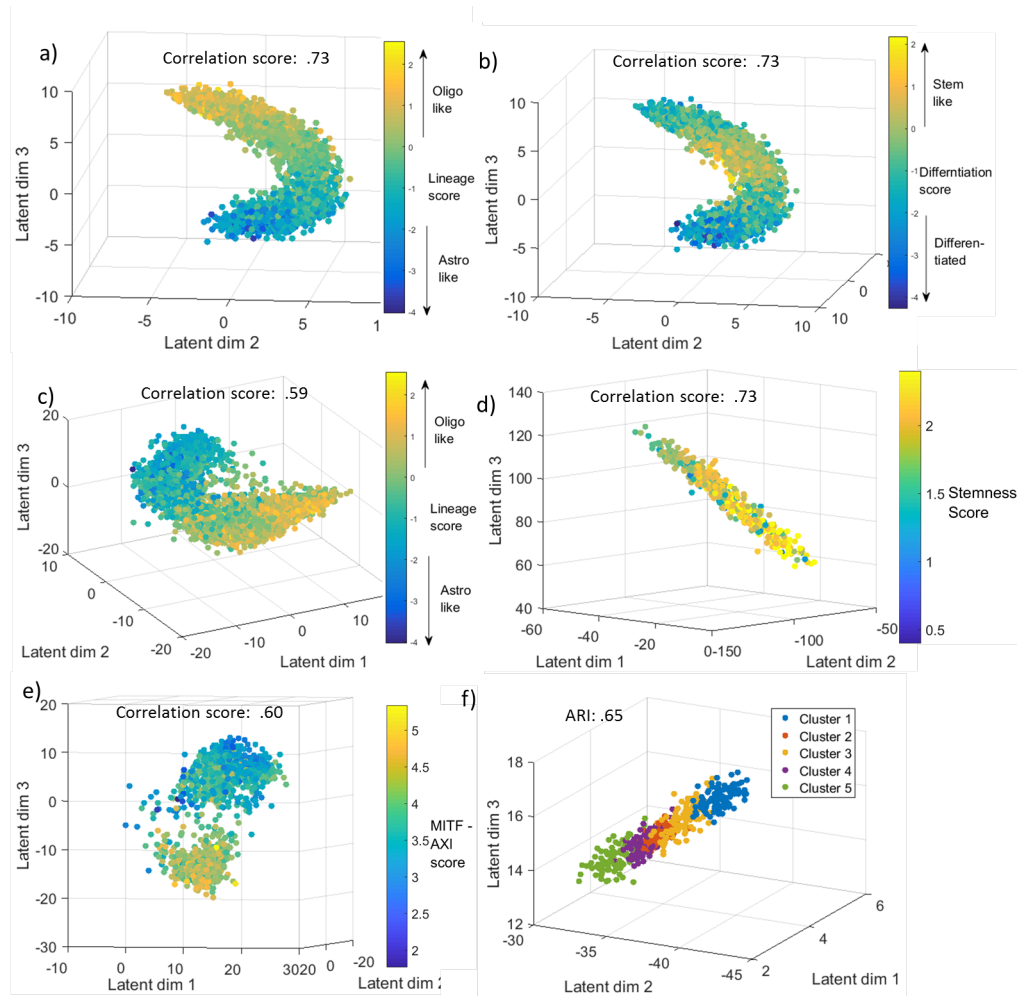


Fig. S1. Performance of Dhaka with ReLU activation function at the output layer. a) and b) Oligodendrogloma dataset with 265 signature genes. c) Oligodendrogloma dataset with auto selected genes. d) Glioblastoma dataset e) Melanoma dataset f) Simulated dataset.

	Dhaka	PCA	t-SNE	ZIFA	SIMLR	MAGIC	NMF	Autoencoder	scVI
ARI	.73	.16	.27	.58	.69	0	0.17	0.72	0.19
Runtime (s)	3.43	.20	3.57	501.23	17.18	0.64	0.15	3.42	16.15

Table S2. Runtime comparison on simulated dataset

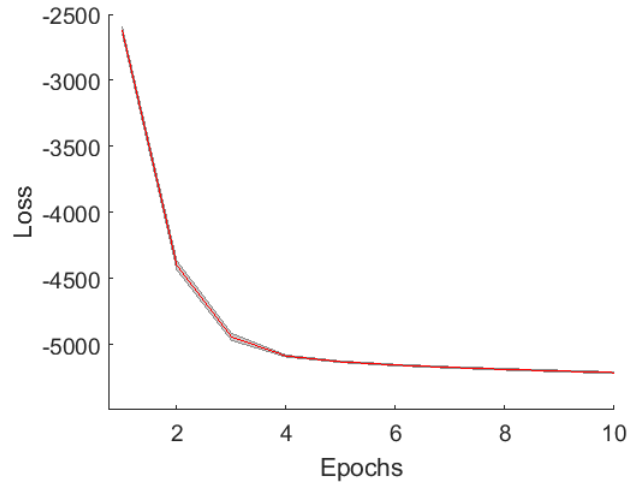


Fig. S2. Loss function plot from 50 independent trials of the Dhaka method on the Oligodendrogloma dataset. The low standard error (marked by the shaded region) in the loss function plot demonstrates the robustness of the Dhaka method.

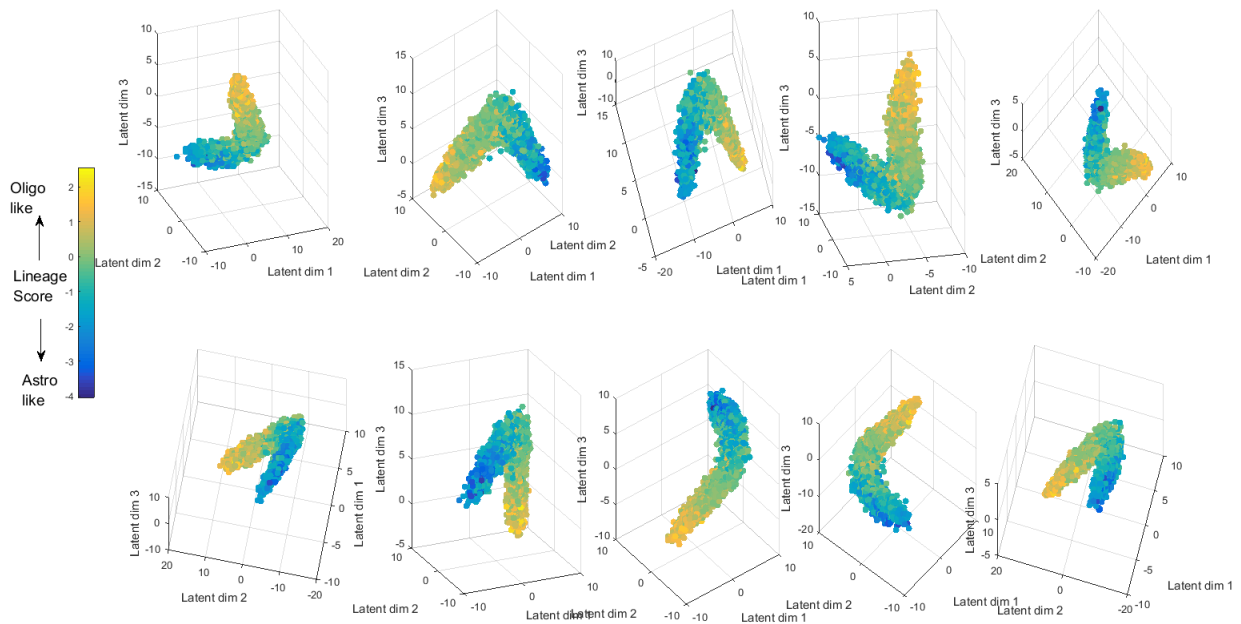


Fig. S3. Dhaka projection from 10 independent runs on Oligodendrogloma dataset with signature genes. Same v-structure is captured in all of the runs.

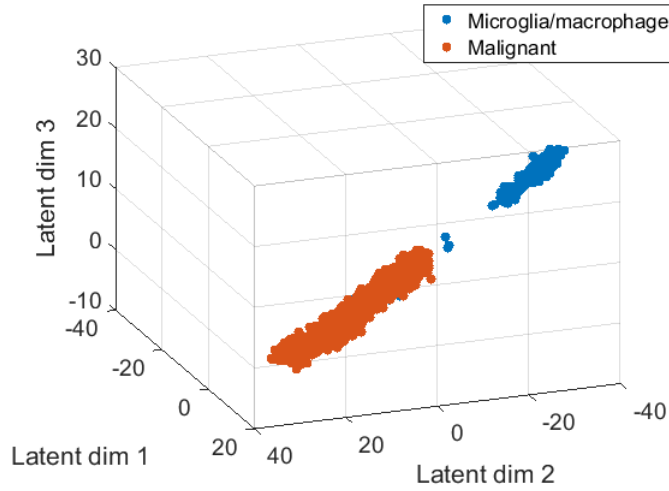


Fig. S4. Oligodendroglioma dataset with 5000 auto-selected genes. Dhaka projection separating malignant cells from non-malignant microglia/macrophage cells.

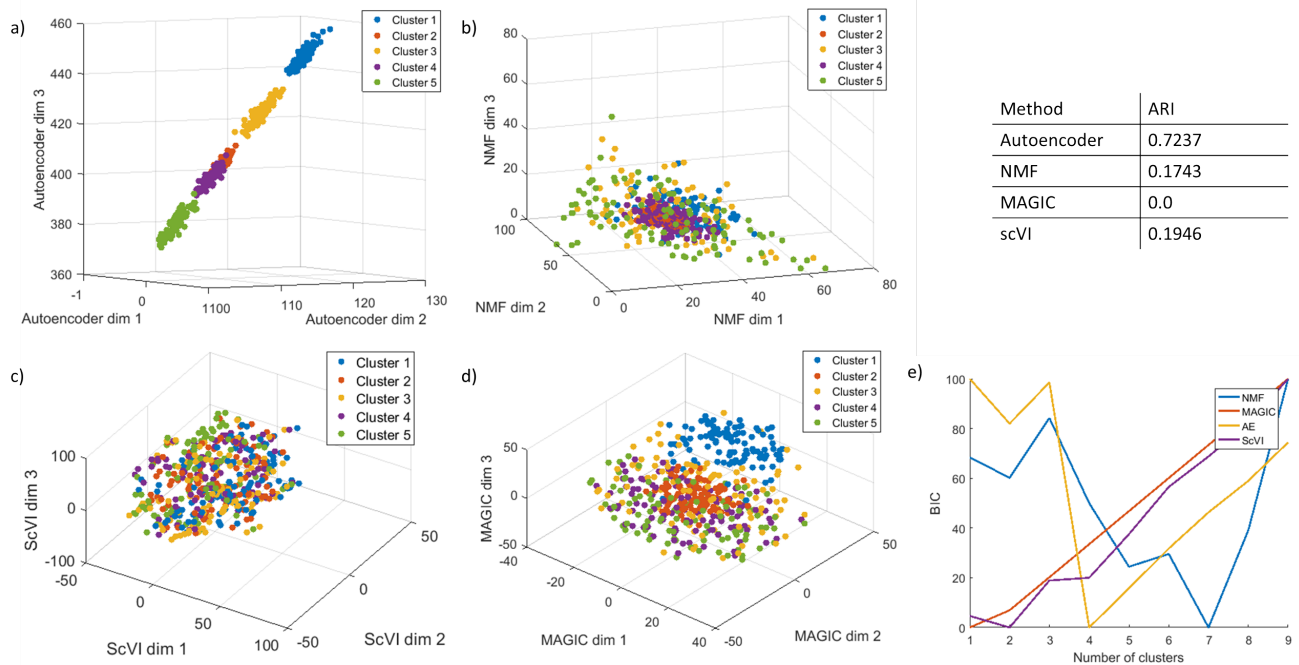


Fig. S5. Performance of regular Autoencoder, NMF, MAGIC, and scVI on simulated dataset. a) Autoencoder, b) NMF, c) MAGIC, d) scVI. The colors correspond to the ground truth cluster ids. e) Plot of BIC calculated from fitting Gaussian Mixture Model to the 3D projection of the data to estimate number of clusters. The number with lowest BIC is considered as the estimated number of clusters in the data.

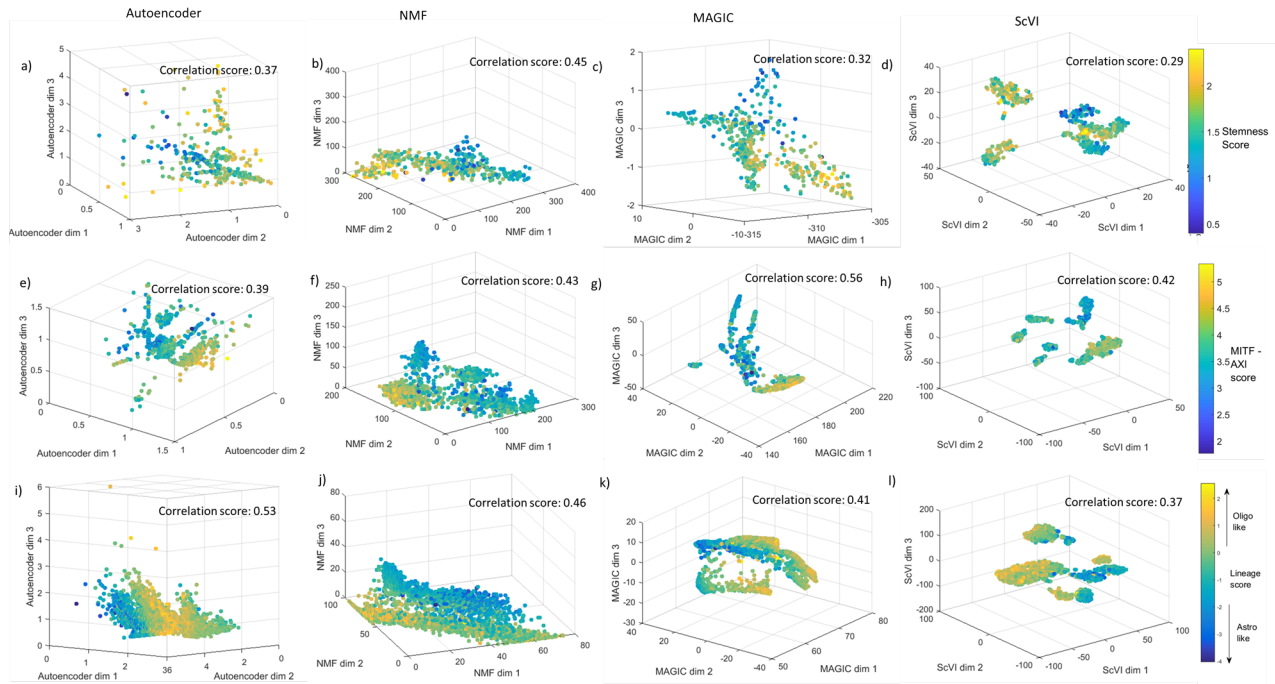


Fig. S6. Performance of regular Autoencoder, NMF, MAGIC, and scVI on single cell RNA-seq datasets. a)-d) Glioblastoma, e)-h) Melanoma, i)-l) Oligodendrogloma. The Spearman rank correlation scores of the scoring metric (a)-d) Stemness score, e)-f) MITF-AXL score, i)-l) lineage and differentiation score) and the learned projections are also shown in the figures.

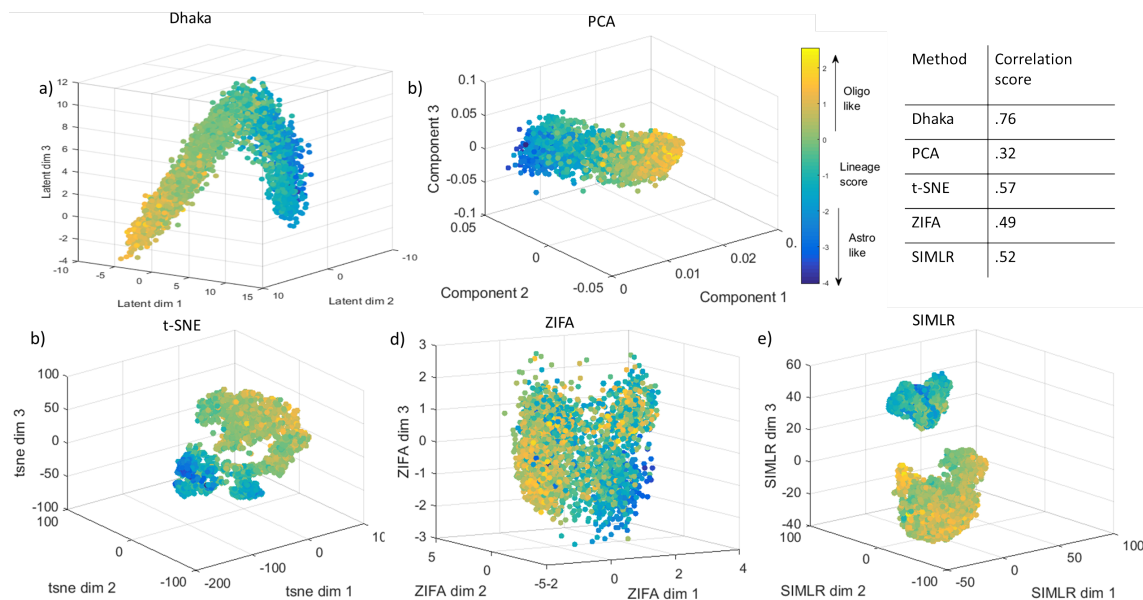


Fig. S7. Comparison of Dhaka with PCA, t-SNE, ZIFA, and SIMLR on Oligodendrogloma dataset with 265 signature genes. a) Dhaka, b) PCA, c) t-SNE, d) ZIFA, e) SIMLR projections colored by lineage score. The Spearman rank correlation scores with the scoring metric (lineage and differentiation score) and Dhaka and the reported method projections are shown in tabular form. We can clearly see that Dhaka preserves the original scoring metric the best.

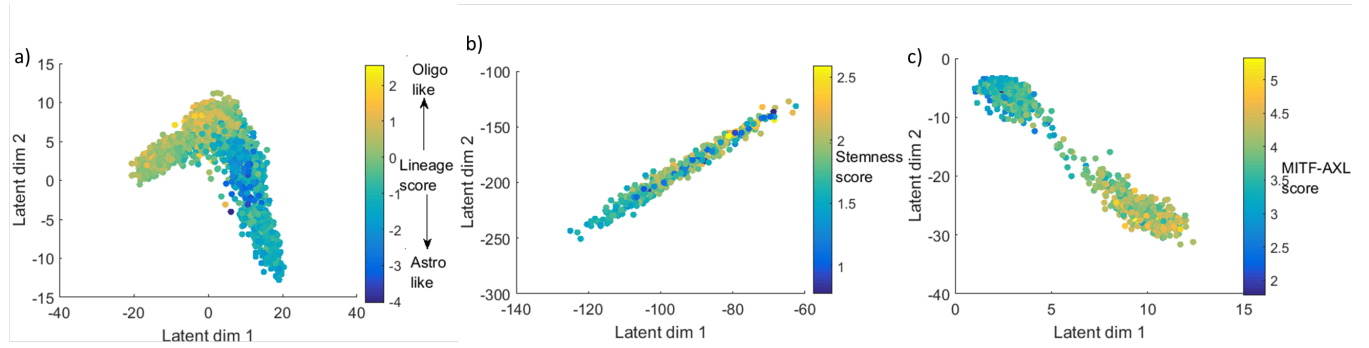


Fig. S8. Two dimensional projections obtained from Dhaka method. a) Oligodendroglioma b) Glioblastoma c) Melanoma. We can see that the v-structure of Oligodendroglioma dataset is preserved in 2D projection as well (Correlation score 0.61). Even though the linear trajectory is preserved in the Glioblastoma dataset, the correlation score decreased from 3D case to 2D case (0.61 from 0.72). For Melanoma dataset, we still see a good separation of the two tumor subpopulation based on MITF-AXL score. The correlation score decreases slightly from 0.68 in 3D projection to 0.66 in 2D projection.

1.5 Correlation score computation

Latent projections from different methods have very different axis range. Hence instead of computing Pearson correlation, Spearman rank correlation is chosen as performance evaluation metric. In clustering evaluation, metrics like nearest neighbor preservation, adjusted Rand index, and silhouette score are typically used as performance metric. However, in this paper the primary focus is to uncover evolutionary trajectory of tumor populations instead of just clustering. Hence, a correlation score with lineage/differentiation metrics are a better indicator of the algorithm performance in latent projection computation.

For 2D scoring metric such as in case of Oligodendroglioma ($X = \text{Lineage score}$, $Y = \text{Differentiation score}$), we computed 2D projections for all the competing methods along with Dhaka. Let us consider, $T \in \mathbb{R}^{n \times 2}$ is the true scoring metric, and $P \in \mathbb{R}^{n \times 2}$ is the learned 2D projection. Then the Spearman correlation coefficient will be a $CC = 2 \times 2$ matrix. Here $CC(i, j) = \text{Corr}(T(:, i), P(:, j))$. Then the overall correlation score is computed as $CS = \max(\sqrt{\frac{1}{2} \sum_i CC(i, i)^2}, \sqrt{\frac{1}{2} \sum_{i \neq j} CC(i, j)^2})$.

In case of 1D scoring metric such as Stemness score in Glioblastoma and MITF-AXL score in Melanoma, we computed correlation coefficient of 1D scoring metric with each dimension of the 3D projection, i.e., $CC \in \mathbb{R}^{1 \times 3}$. Then CS is simply set as $CS = \max(CC)$.

1.6 Dropout analysis

To evaluate robustness of the method to drop out genes, we introduced additional dropouts in the Oligodendroglioma dataset. Each gene in the dataset have fraction of dropout ranging from 0 to 0.8 (Please see Supporting Fig. S9a). To introduce additional dropouts, for each gene we randomly select a subset of cells, for example in the 20% case we randomly select 20% of the cells for each gene. Then we make the expression value of these cells for that gene 0. After artificially introducing 20% more drop out, we should expect all the genes will have minimum drop out fraction of .2. Hence we see that now the histogram is shifted to the right from 0 to 0.2. We can see similar scenario for 30% and 50% case. Note that since we are randomly selecting cells to be dropped out, some of the selected cells might already have 0 expression values. Hence when we introduce 50% additional dropout, the histogram actually shifts to 0.4 instead of 0.5. As can be seen from the Fig. S9 up to

30% additional dropout, Dhaka can still retain the v-structure. At 50%, we lose the v-structure, but the method can still separate oligo-like and astro-like cells even with this highly sparse data.

1.7 t-SNE, ZIFA, SIMLR, scVI, Autoencoder, MAGIC implementation details

To compute t-SNE projection, first the input data was pre-processed through PCA. PCA reduced the dimensions to 50. Then t-SNE was performed on the reduced dimension dataset. The perplexity parameter was set at 20. We have tested perplexity parameters {10, 20, 30, 40, 50} and 20 gave the best projection (in terms of correlation score) for our discussed datasets. We have used MATLAB implementation of basic t-SNE (<https://lvdmaaten.github.io/tsne/>).

The SIMLR MATLAB implementation was used in the paper (<https://github.com/BatzoglouLabSU/SIMLR>). The input read count datasets were log10 transformed as required by the package. The SIMLR package requires an estimated number of clusters to compute the similarity matrix. The estimated number of clusters were computed using *Estimate_Number_of_Clusters_SIMLR.m* function of the package.

The python implementation of ZIFA package was used from <https://github.com/epierson9/ZIFA>.

The python implementation of scVI package was used from <https://github.com/YosefLab/scVI>. The original package reduces the data dimension to 10. The dimension was further reduced to 3 using t-SNE (as mentioned in the package).

We have used the same NN structure from Dhaka (without the custom variational loss layer) for the regular Autoencoder implementation using Keras.

The MAGIC MATLAB implementation was used from <https://github.com/KrishnaswamyLab/MAGIC> with the default parameters used in the package.

SIMLR and scVI also use t-SNE for dimensionality reduction following their similarity matrix estimation. Since t-SNE is stochastic in nature, we get slightly different output in each run. Hence, t-SNE, SIMLR, and scVI were ran 10 times on each dataset and average ARI/Correlation score were reported in the result section.

Except for SIMLR, for all other methods including NMF and PCA log2 transformed TPM counts were used as inputs.

1.8 Analysis of Astrocytoma data

The Astrocytoma dataset contains a total of 6341 cells with about 23K genes, among which 5097 are malignant cells. Astrocytoma is another type of brain tumor and this is a followup dataset from the Oligodendroglioma. Hence, we performed same analysis as for Oligodendroglioma. The non-malignant microglia/macrophage cells were clearly separated from the malignant cells (Fig. S12a) in this dataset too. The authors did not compute differentiation and lineage metric for this dataset, but did mention that most of the cells fall in the intermediate state. When we fed the expression profile of the malignant cells to Dhaka, it correctly placed most of the cells near the bifurcation point of the v-structure (Fig. S12b). For reference, we have also showed the Oligodendroglioma cells in the same plot colored by their differentiation score.

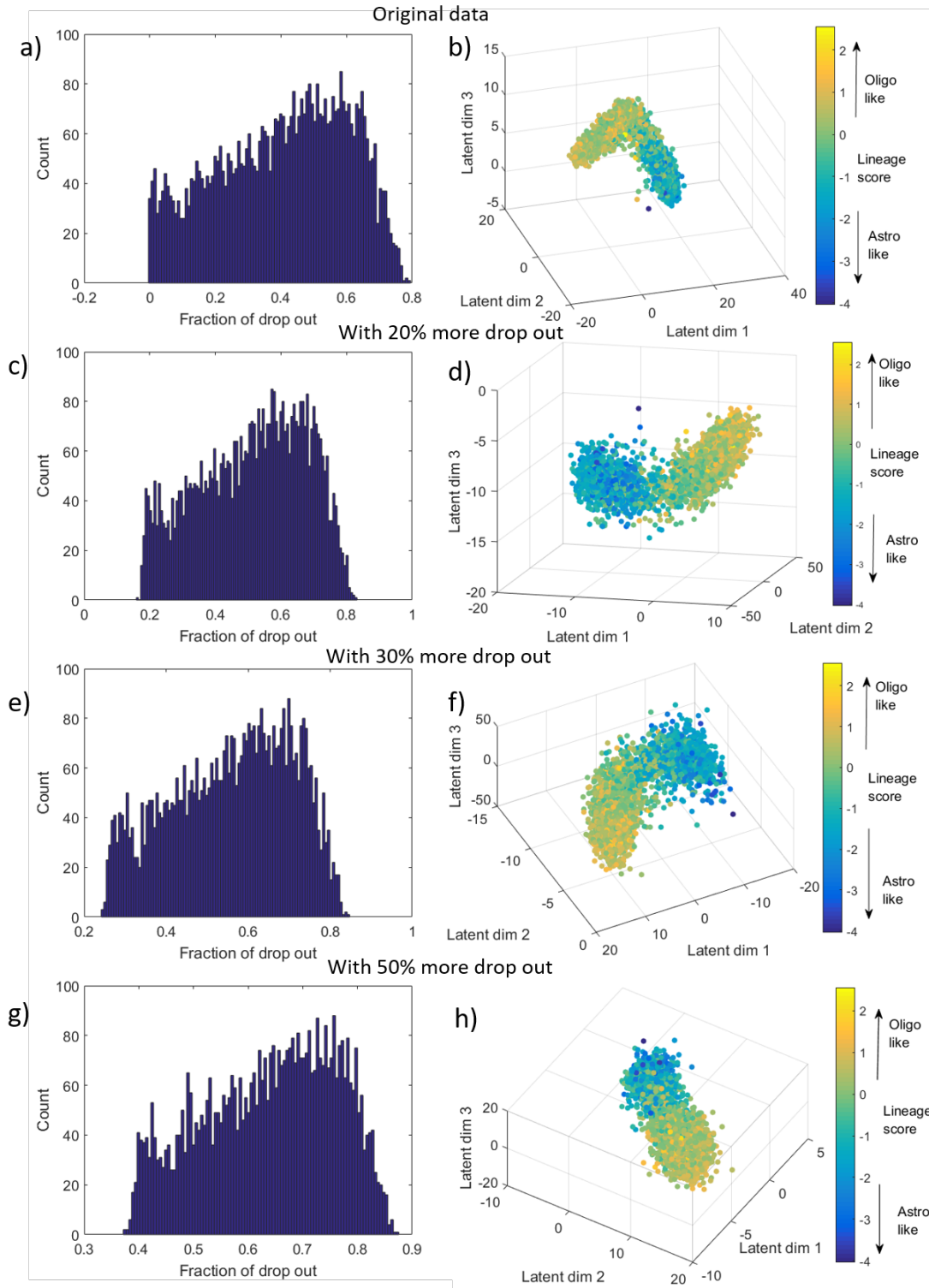


Fig. S9. Robustness analysis with Oligodendrogloma. a,c,e,g Histogram of drop out fraction in each gene after forcing 0%, 20%, 30%, and 50% more genes to be dropped out. b,d,f,h corresponding Dhaka projection of the data. We can see that up to 30%, Dhaka can correctly identify v-structure. Beyond that it loses the v-structure but still shows good separation between oligo-like and astro-like cells.

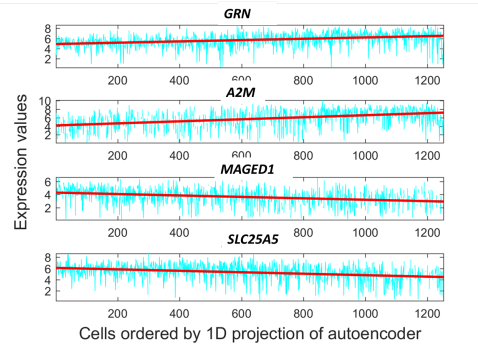


Fig. S10. Known marker genes for Melanoma MITF-AXL program.

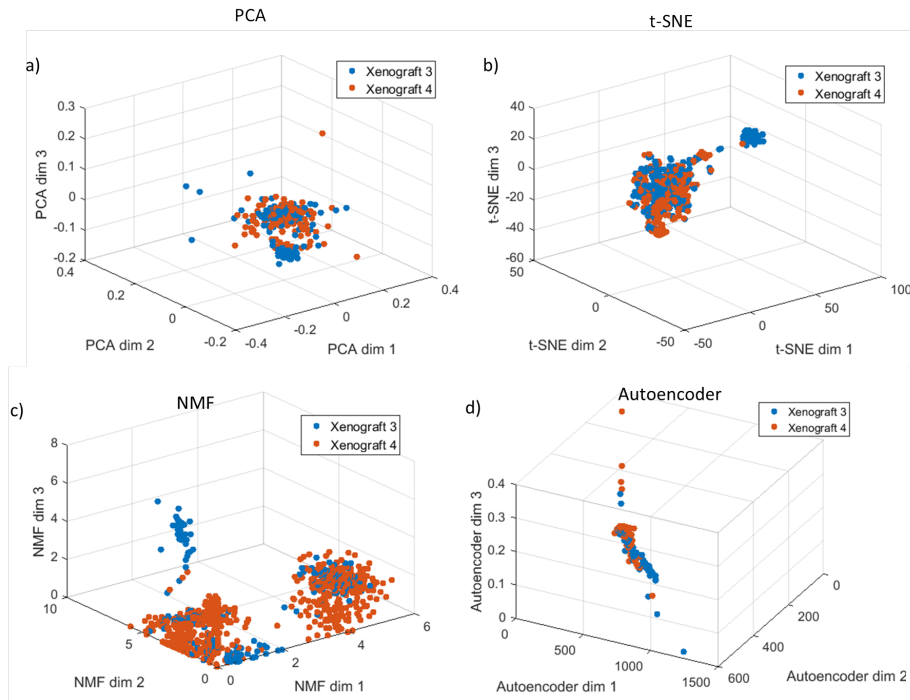


Fig. S11. Performance of PCA, t-SNE, NMF, and regular autoencoder on two xenograft breast tumor samples' copy number profile. We can see that the separation between major and minor cluster of xenograft 3 is not as distinct as Dhaka.

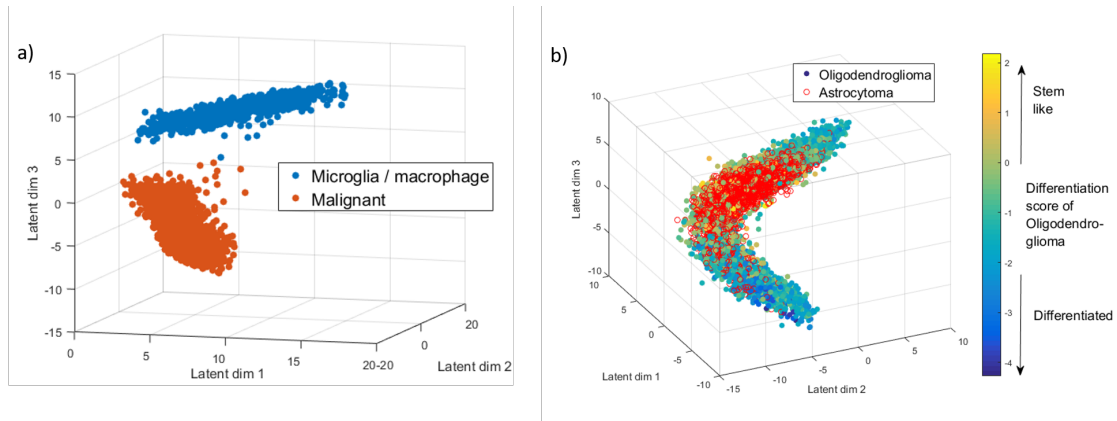


Fig. S12. Astrocytoma dataset. a) Dhaka output of Astrocytoma dataset with 5000 autoselected genes separating malignant cells from microglia/macrophage cells. b) Dhaka output from relative expression profile of malignant Astrocytoma cells (red) along with malignant Oligodendrogloma cells.

References

1. F. Chollet, “keras,” *GitHub repository*, 2015.
2. L. Kaufman and P. J. Rousseeuw, *Finding groups in data: an introduction to cluster analysis*, vol. 344. John Wiley & Sons, 2009.
3. C. Trapnell, D. Cacchiarelli, J. Grimsby, P. Pokharel, S. Li, M. Morse, N. J. Lennon, K. J. Livak, T. S. Mikkelsen, and J. L. Rinn, “The dynamics and regulators of cell fate decisions are revealed by pseudotemporal ordering of single cells,” *Nature biotechnology*, vol. 32, no. 4, pp. 381–386, 2014.

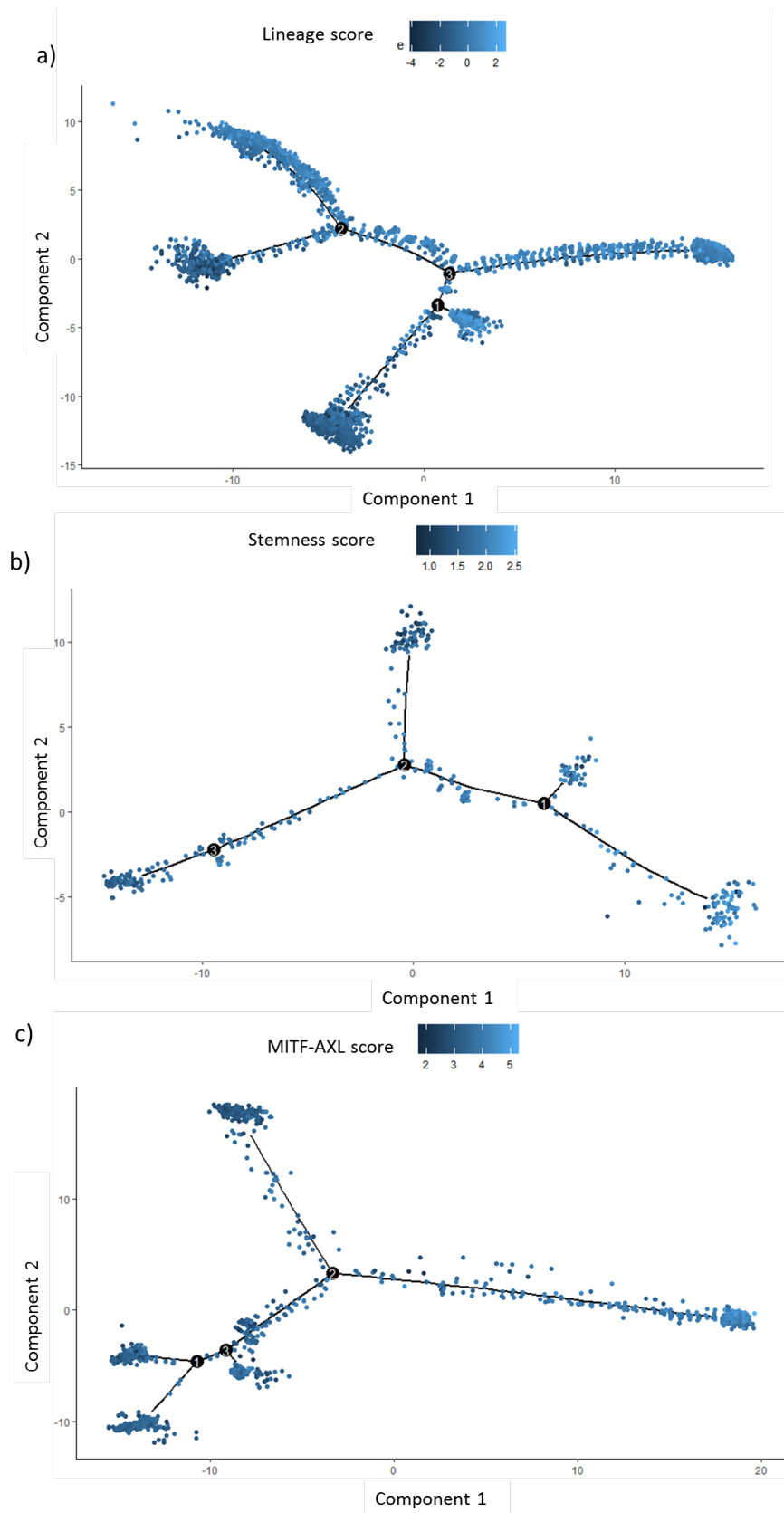


Fig. S13. Monocle [3] pseudotime ordering of a) Oligodendroglioma, b) Glioblastoma, c) Melanoma datasets. As can be seen from the figure, the pseudotime orderings from Monocle fail to capture the biological trend of the data. Monocle divides the cells in each dataset to multiple branches but none of the branches show significant correlation with the underlying tumor subpopulation and/or the scoring metrics. The ‘R’ package of Monocle was installed through bioconductor.