

A mixture-of-experts deep generative model for integrated analysis of single-cell multiomics data

Kodai Minoura, Ko Abe, Hyunha Nam, Hiroyoshi Nishikawa, Teppei Shimamura

2021, Cell Reports Methods

2022.7.13

Motivation

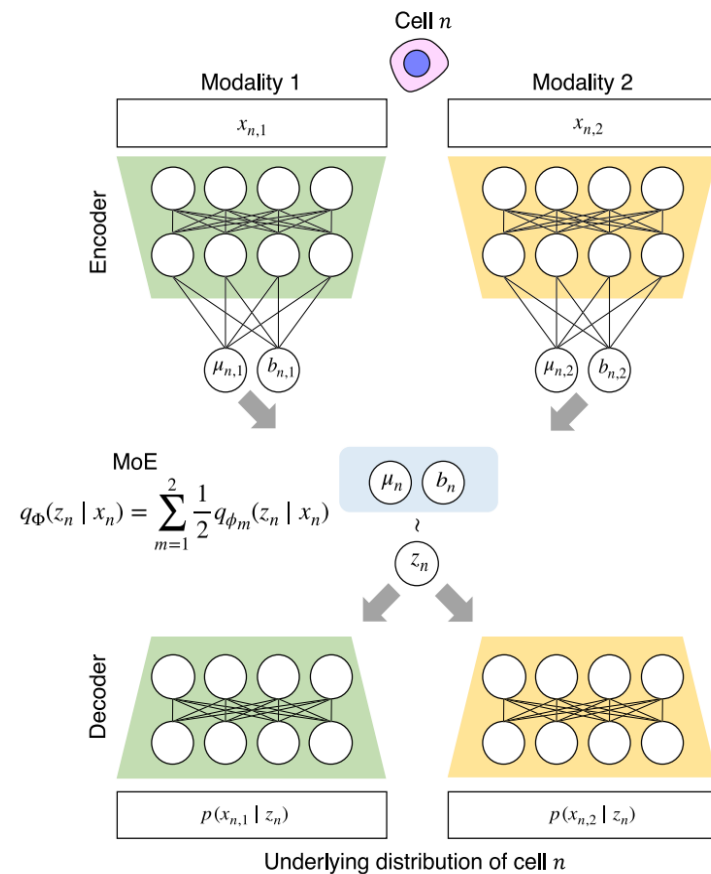
- ✓ Highly complex of single-cell multimodal data
- ✓ multimodal low-dimensional joint representations
- ✓ Cross-modal predictions is unsolved.
- ✓ “Black-box” nature of deep learning.

Highlights

- Learning low-dimensional joint representations from single-cell multi-omics data
- Detecting previously overlooked cell populations in single-cell multimodal data
- Accurately predicts missing modalities by crossmodal generation
- Pseudocell generation enables scMM to learn interpretable latent dimensions

Architecture

A



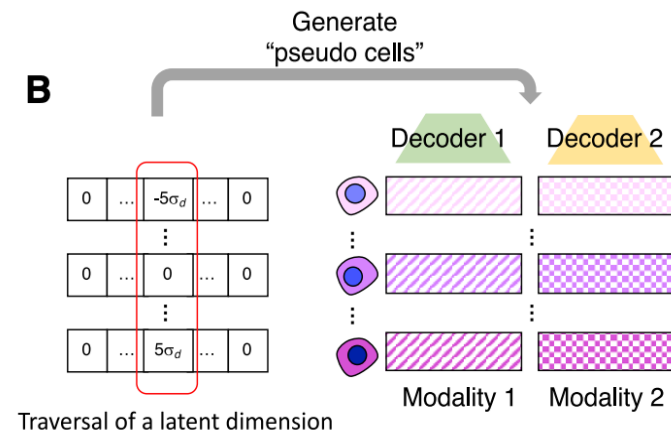
Multimodal latent variables

- Clustering
- Visualization

Data generation

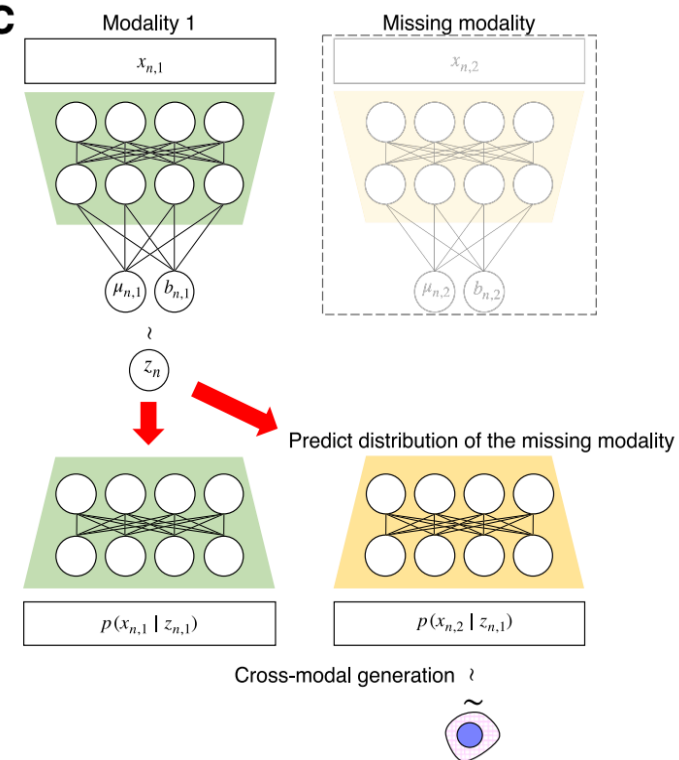
- Prediction of missing modality
- Integration of datasets

B

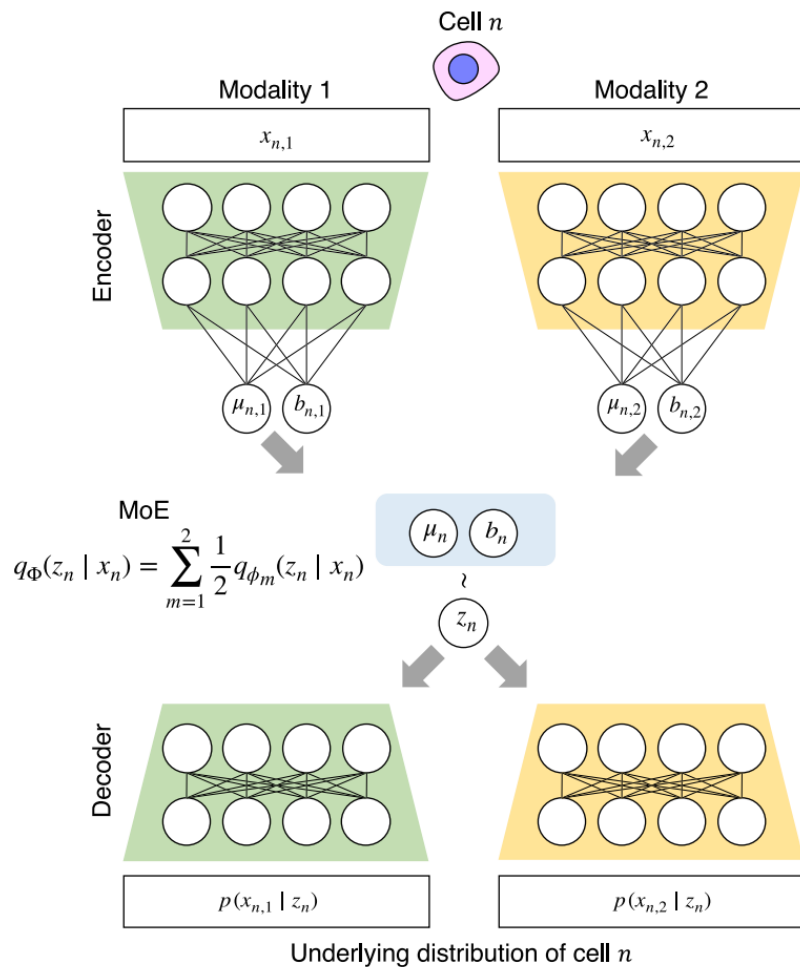


Calculate correlation

C



A



Multimodal latent variables

- Clustering
- Visualization

Data generation

- Prediction of missing modality
- Integration of datasets

Encoder: $q_{\phi_m}(\mathbf{z} | \mathbf{x}_m)$

Sampling:

1. Transcriptome and surface protein

Negative Binomial Distribution

$$X \sim \text{NB}(r ; P)$$

Non-negative counts with overdispersion (Gayoso et al., 2021)

2. Chromatin accessibility

zero-inflated negative binomial (ZINB) distribution

Extreme sparsity

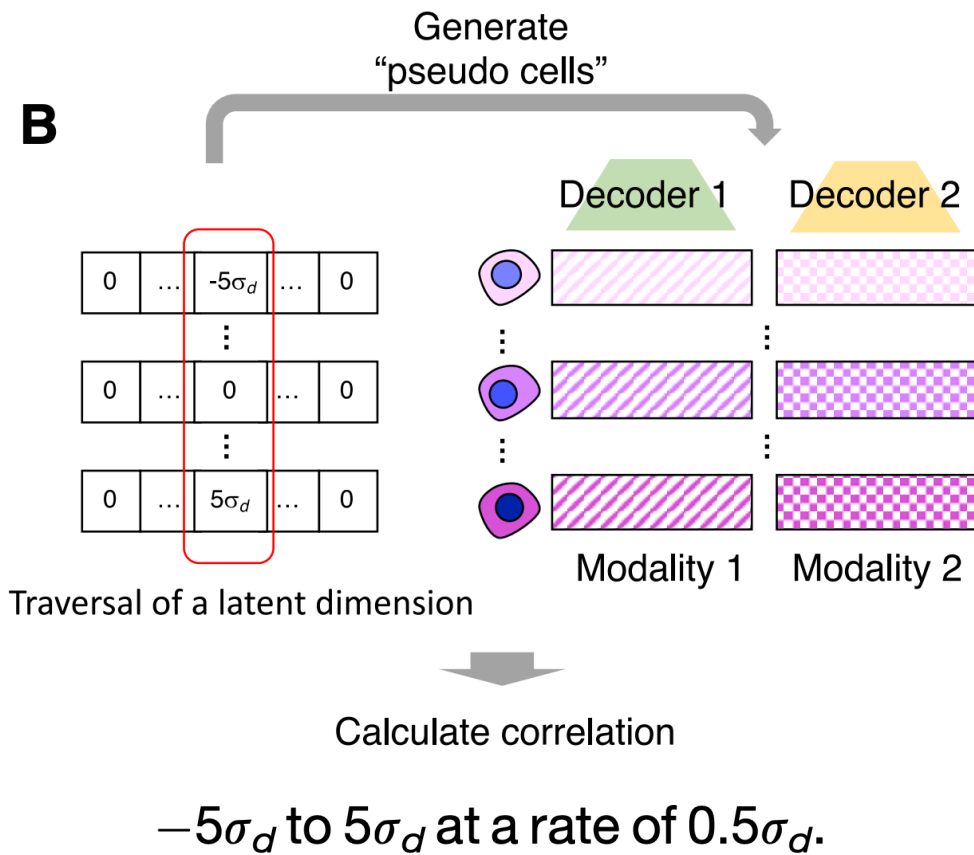
MOE (mixture-of-experts)

$$q_{\Phi}(z_n | x_n) = \sum_{m=1}^2 \frac{1}{2} q_{\phi_m}(z_n | x_n)$$

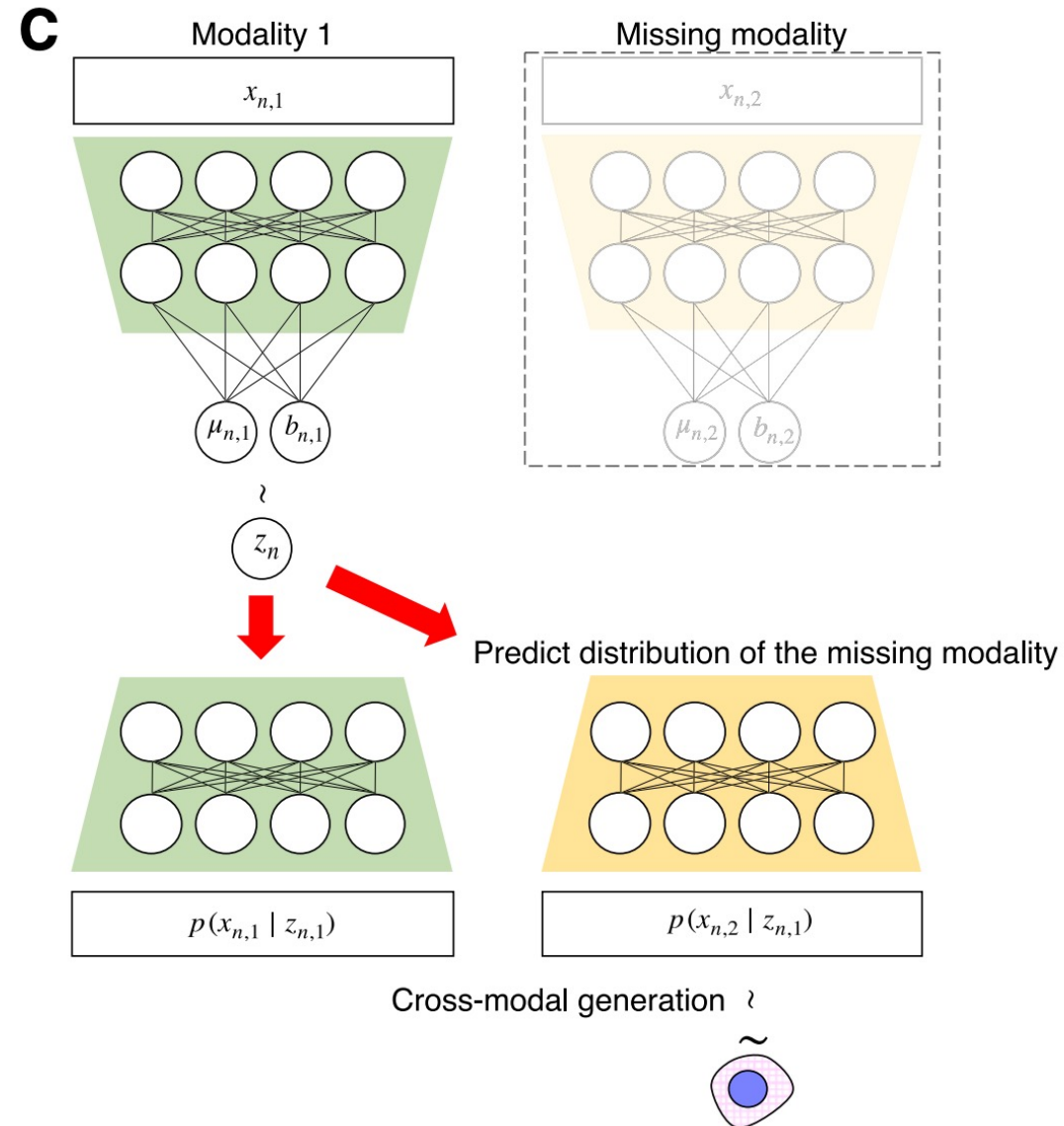
Decoder: $p_{\theta_m}(\mathbf{x}_m | \mathbf{z})$

$$\text{Loss} = \mathbb{E}_{\mathbf{z} \sim q_{\Phi}(\mathbf{z} | \mathbf{x})} [\log p_{\Theta}(\mathbf{x} | \mathbf{z})] - \text{KL}[q_{\Phi}(\mathbf{z} | \mathbf{x}) \parallel p(\mathbf{z})],$$

Interpretability



Cross-model Generation



Experiment

Transcriptome and surface protein

PMBC of CITE-seq dataset (接种疫苗患者的外周血单核细胞数据)

the transcriptome and 224 surface protein, and 54 cell populations
measurements for over 160,000 cells

BMNC if CITE-seq (接种疫苗患者的骨髓单核细胞数据)

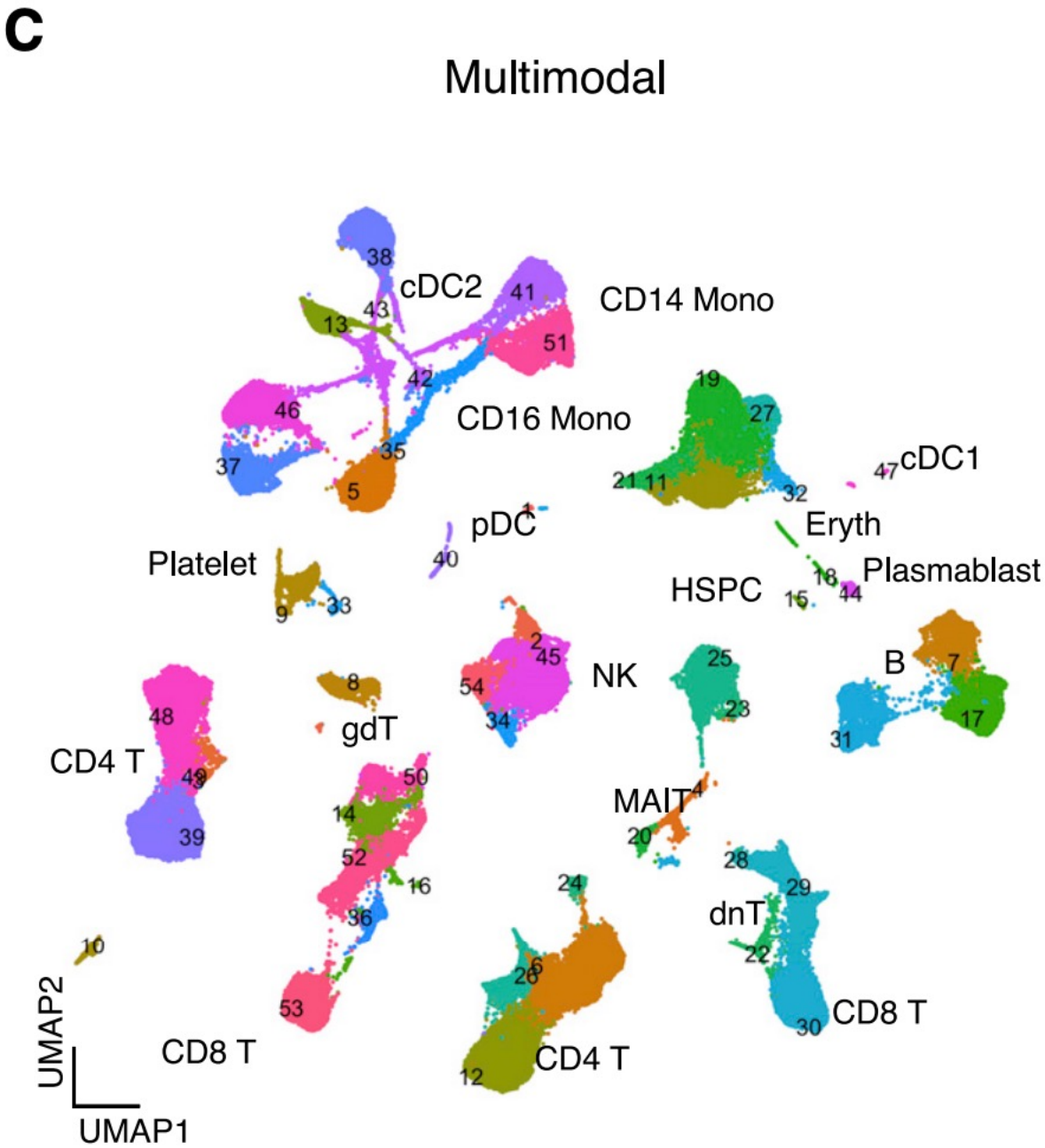
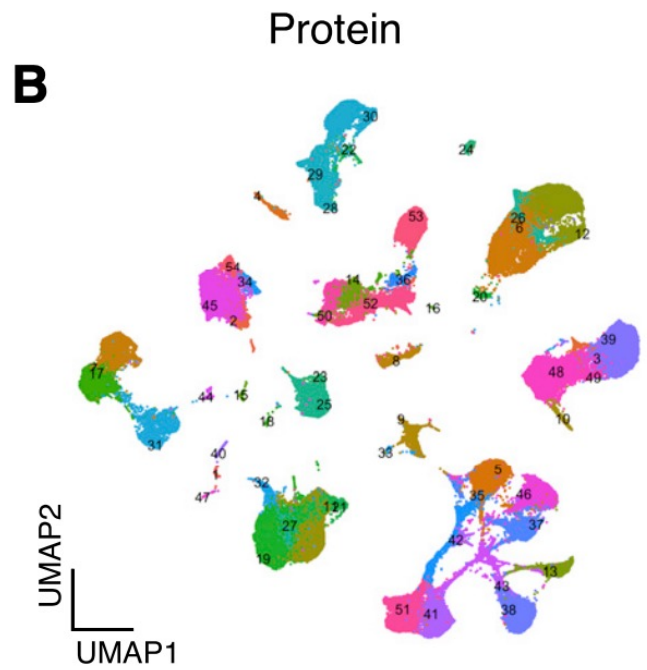
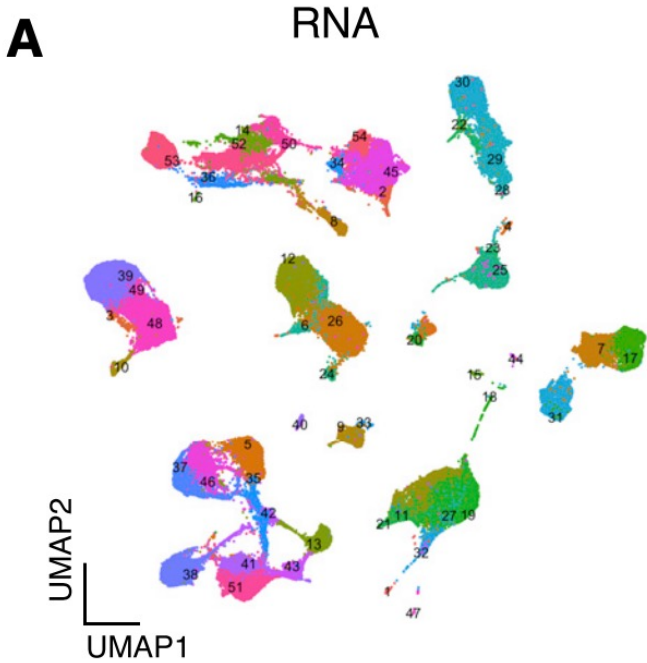
30,000 cells with transcriptome and 25 surface protein

Used for crossmodal estimation

Transcriptome and Chromatin accessibility

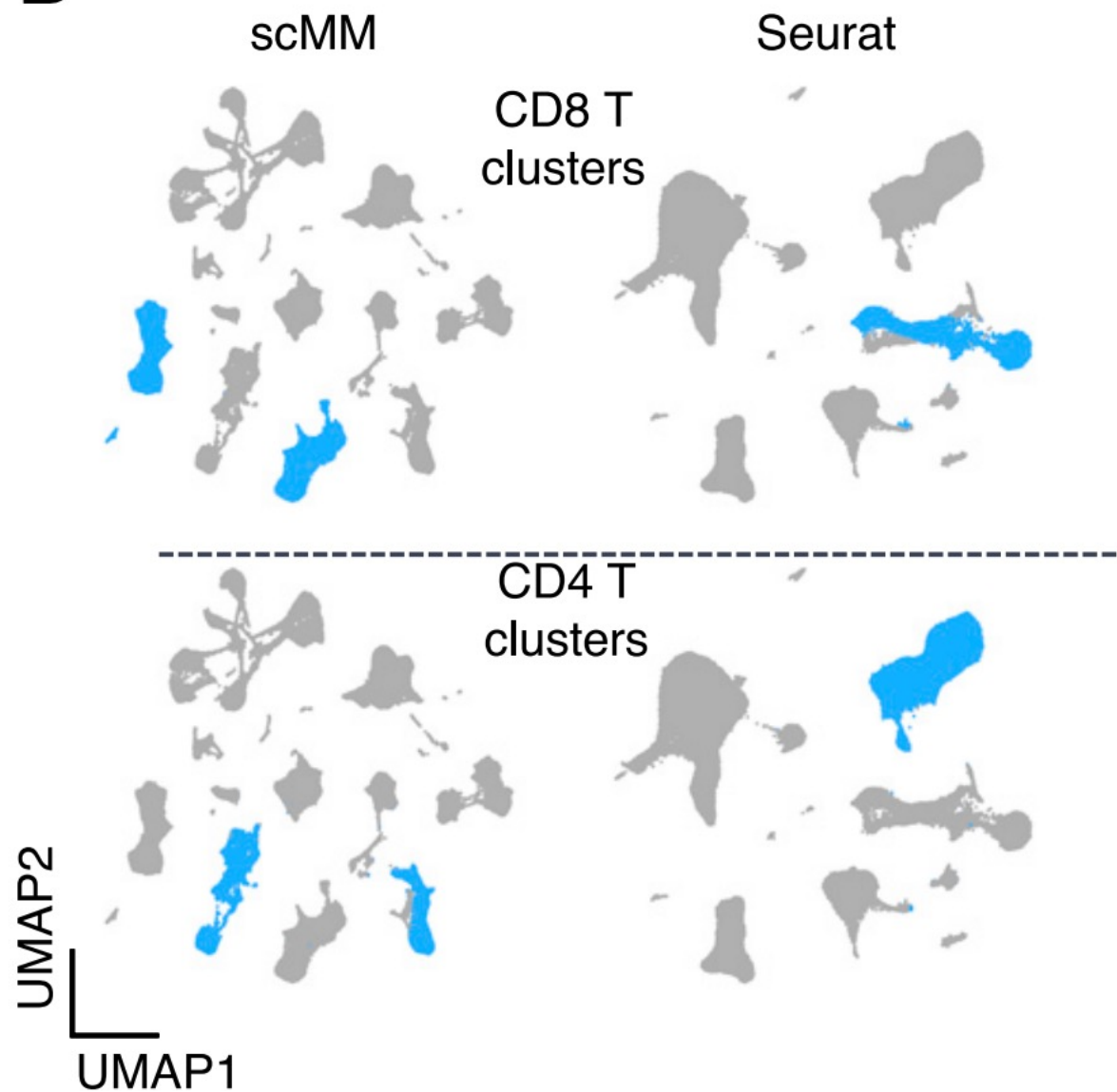
SHARE-seq dataset (小鼠皮肤单细胞转录组和染色质可及性)

UMAP
Visualization



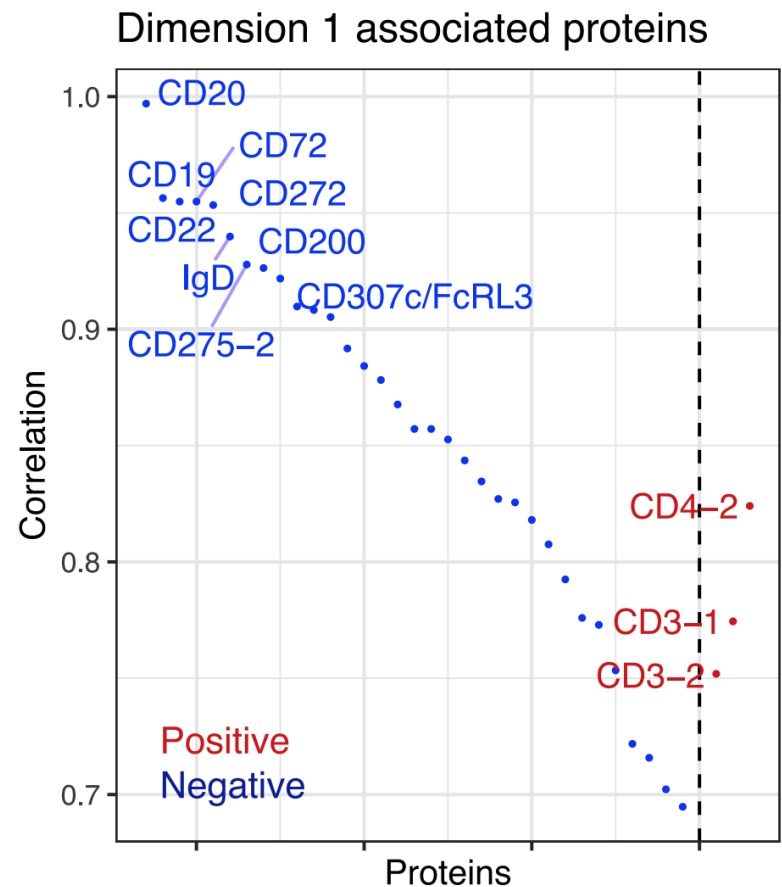
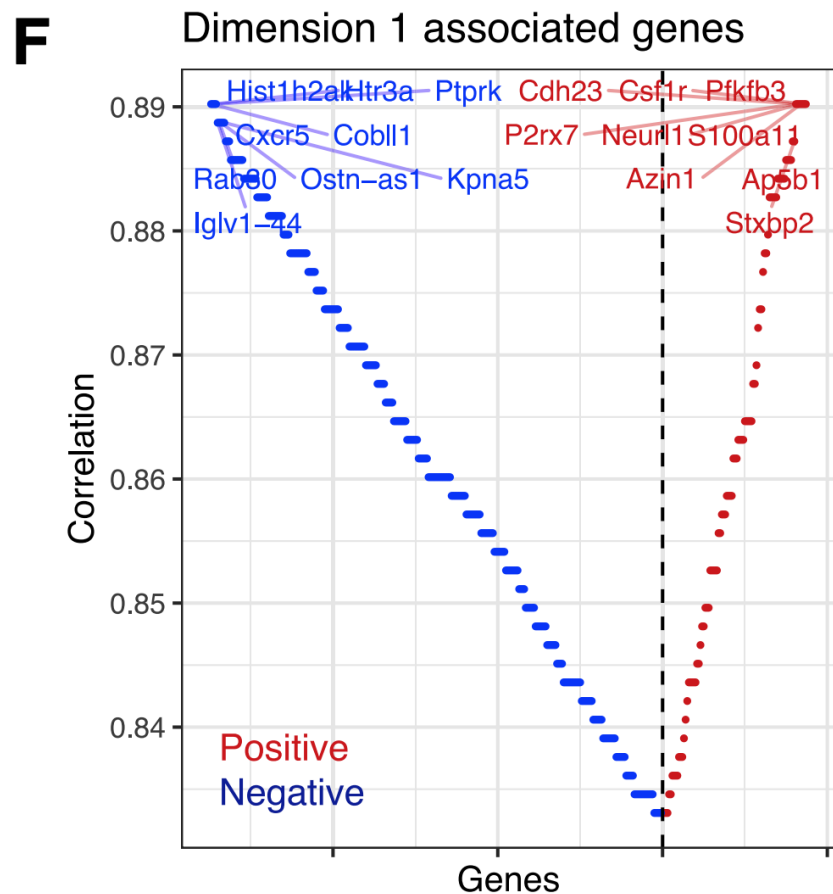
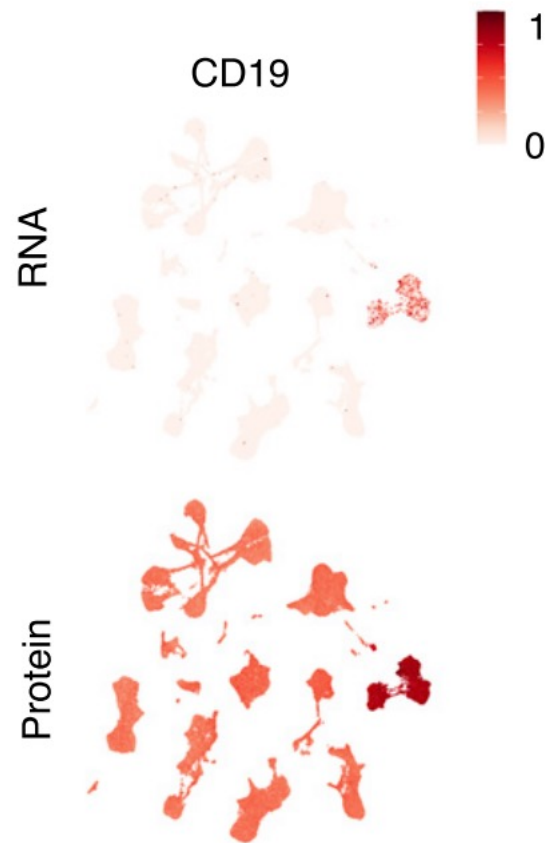
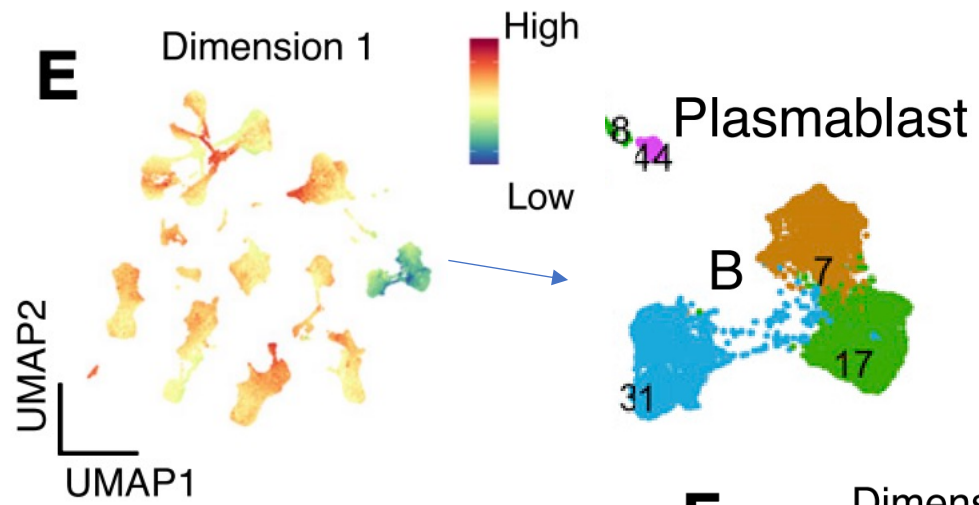
UMAP Visualization

D



MultiModal Heterogeneity Discover

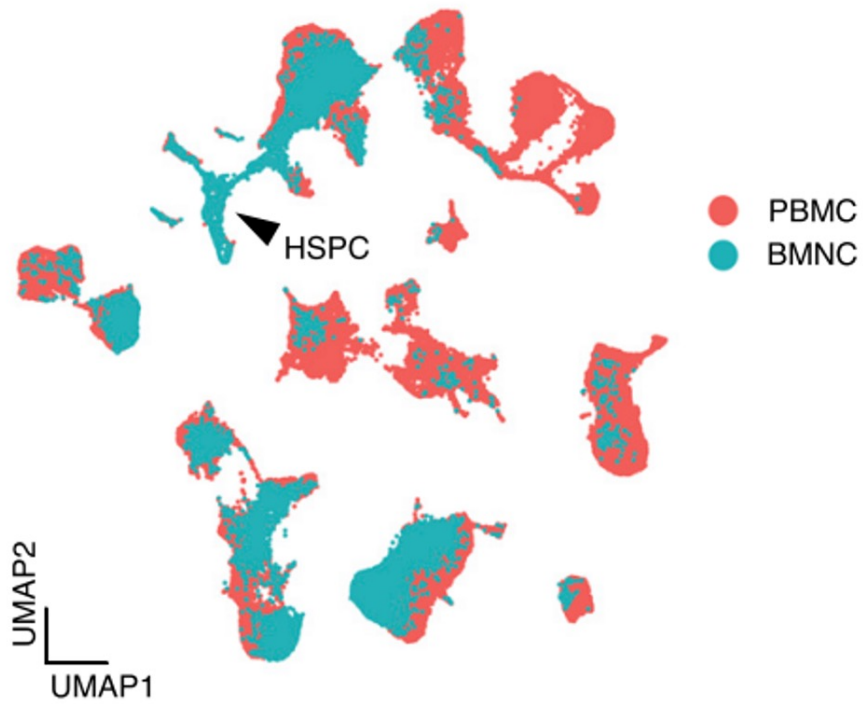
Interpretability Analysis



Crossmodal Prediction

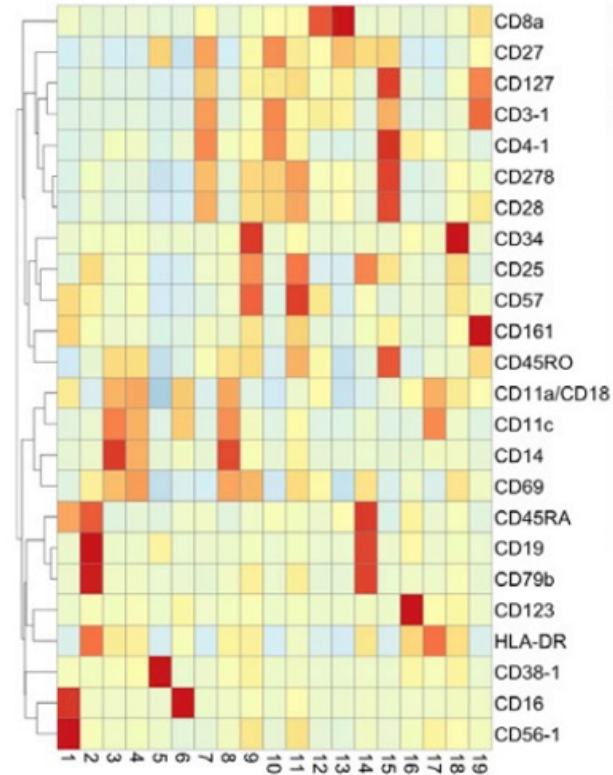
transcriptome-to-protein crossmodal estimation

A

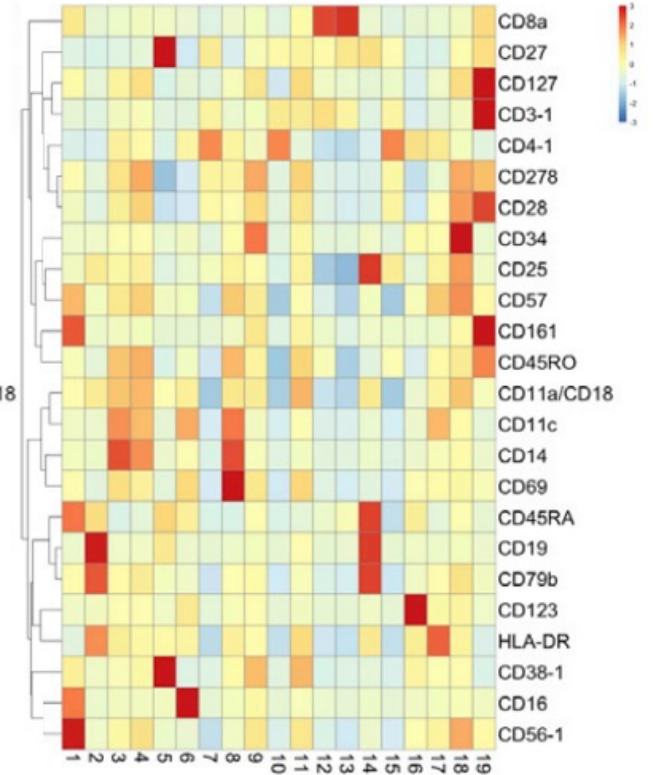


training with PBMC and test with BMNC

Original 24 proteins



Predicted 24 proteins



$$p(x_{n,2} | z_{n,1})$$

SHARE-seq dataset

