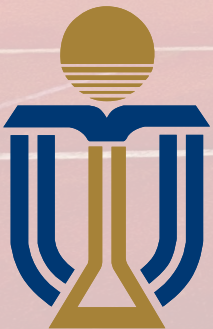


Foundation Models on Single Cell RNA Sequencing Data

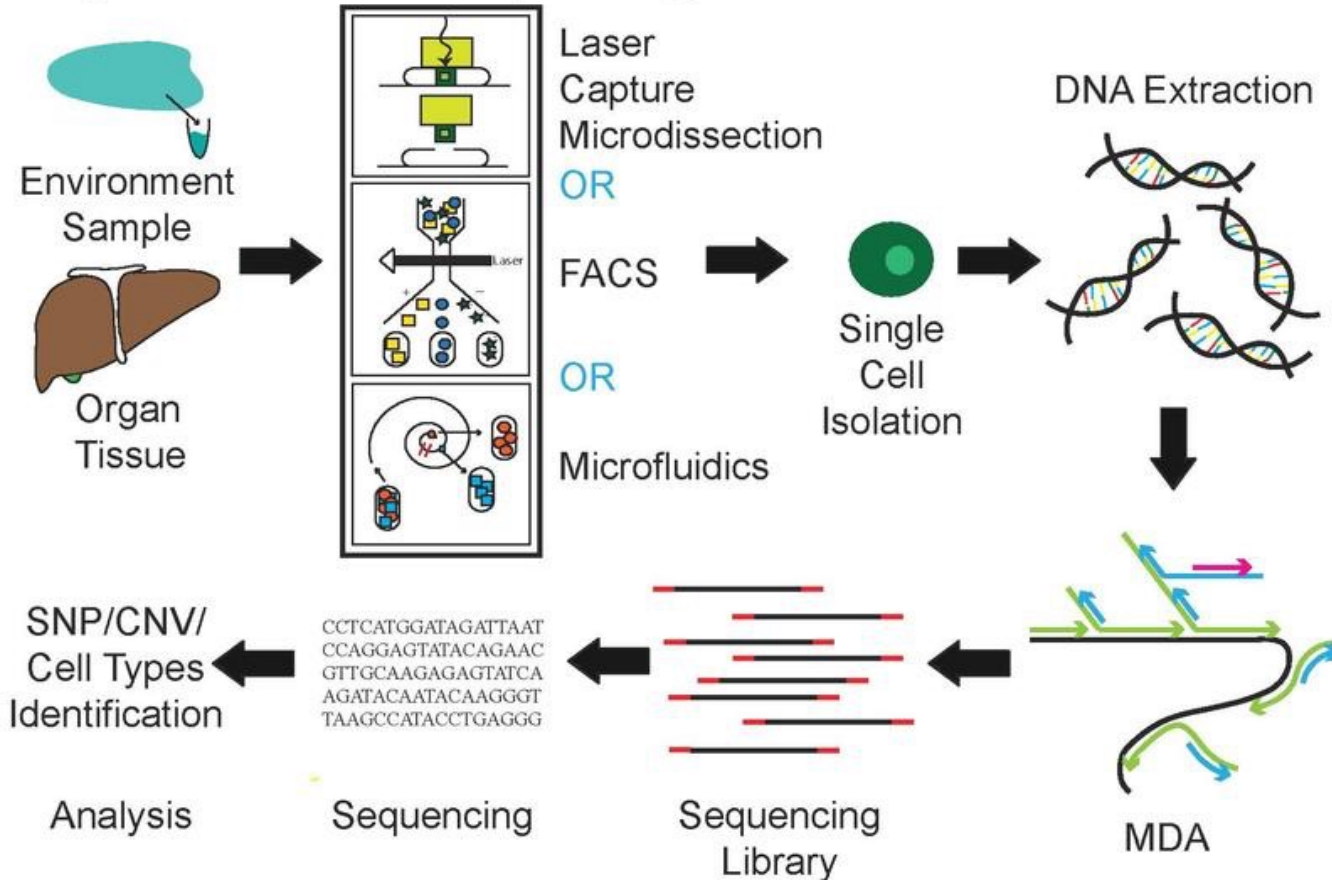
09/09/2023

Minghao WANG



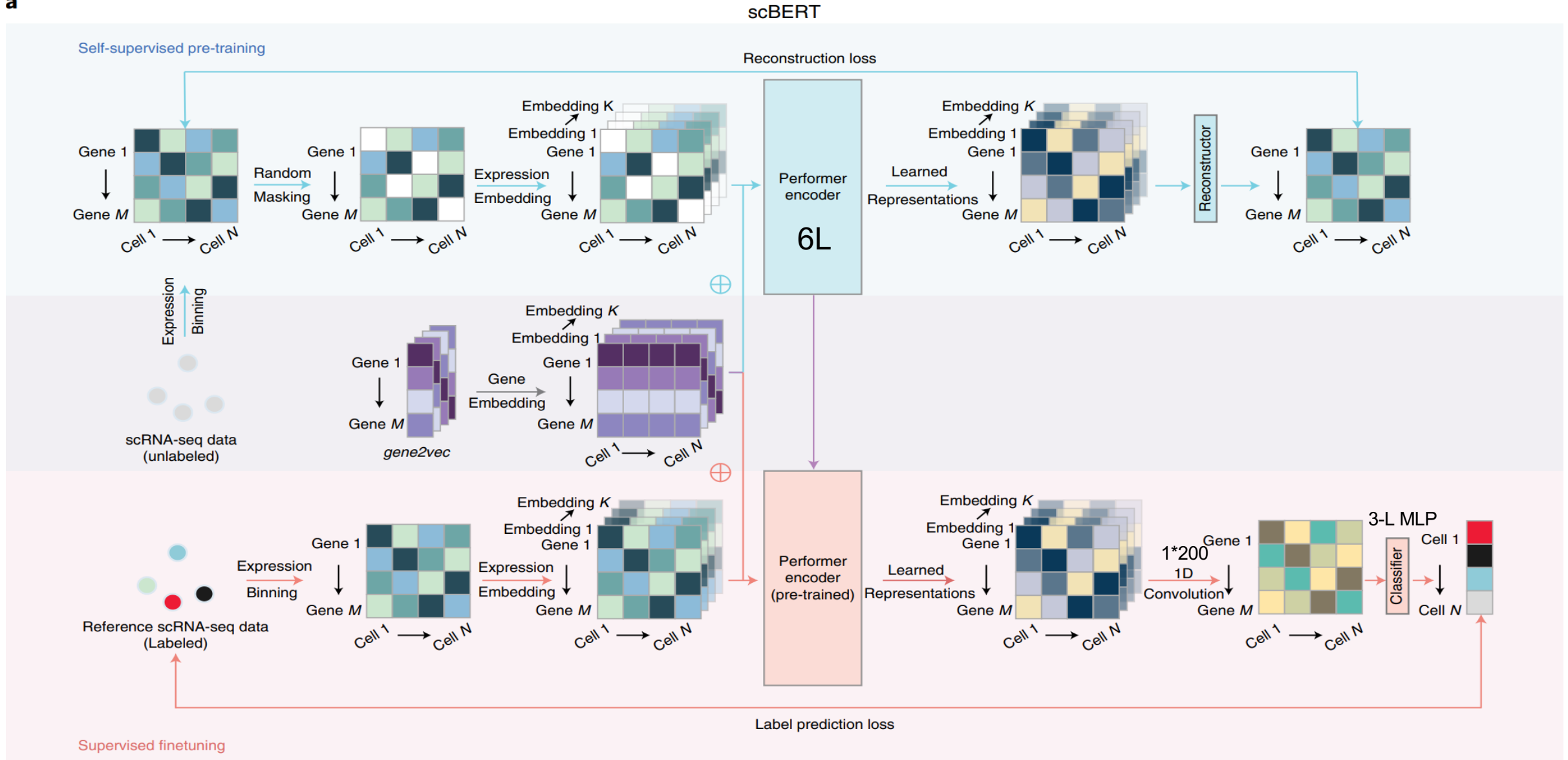
scRNA-seq provides information on which genes are expressed at the level of individual cells.

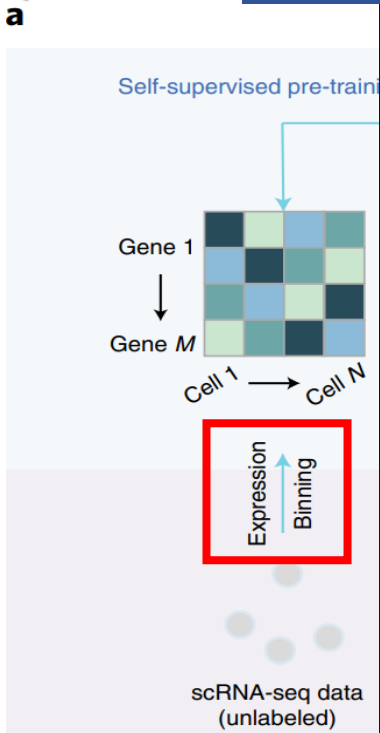
Single Cell Genome Sequencing Workflow



- Explore which cell types are present in tissues
- Identify unknown/rare cell types or states
- Elucidate gene expression changes during differentiation or across time or across states
- ...

| | orig.ident | nCount_RNA | nFeature_RNA | percent.mt | rename_sample | type |
|--------------|------------|------------|--------------|------------|---------------|------|
| MIR1302-10 | P1194.F | 10076.0 | 3081 | 7.423581 | P1194.F | WT |
| FAM138A | P1194.F | 1772.0 | 805 | 1.523702 | P1194.F | WT |
| OR4F5 | P1194.F | 4984.0 | 2143 | 4.133226 | P1194.F | WT |
| RP11-34P13.7 | P1194.F | 3292.0 | 1312 | 9.538275 | P1194.F | WT |
| RP11-34P13.8 | P1194.F | 6060.0 | 2278 | 3.036304 | P1194.F | WT |
| AC145205.1 | ... | ... | ... | ... | ... | ... |
| ... | P969 | 9172.0 | 2609 | 7.064980 | P969 | WT |
| ... | P969 | 29911.0 | 5913 | 5.613320 | P969 | WT |
| ... | P969 | 16849.0 | 4393 | 8.837320 | P969 | WT |
| ... | P969 | 10355.0 | 3377 | 9.570256 | P969 | WT |
| ... | P969 | 4188.0 | 2064 | 7.139446 | P969 | WT |



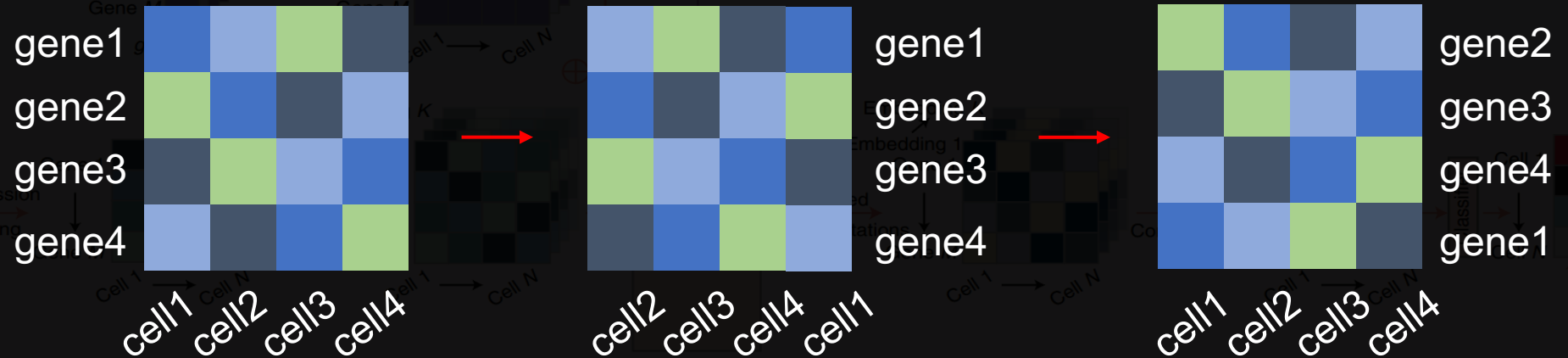


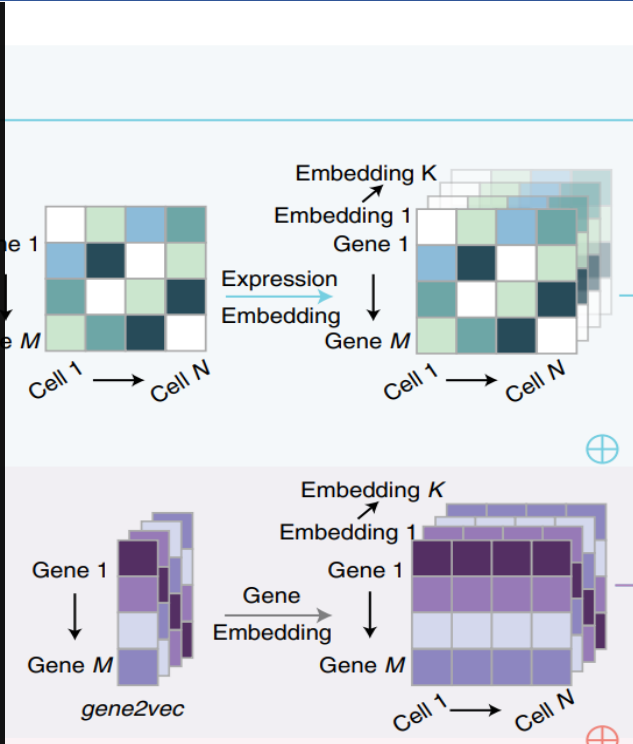
1. Expression Binning

$$x_j^{(i)} = \begin{cases} k, & \text{if } X_{i,j} > 0 \text{ and } X_{i,j} \in [b_k, b_{k+1}] \\ 0, & \text{if } X_{i,j} = 0 \end{cases}$$

2. Cell by Gene Matrix

What does each block mean? What if shuffling columns? What about rows?





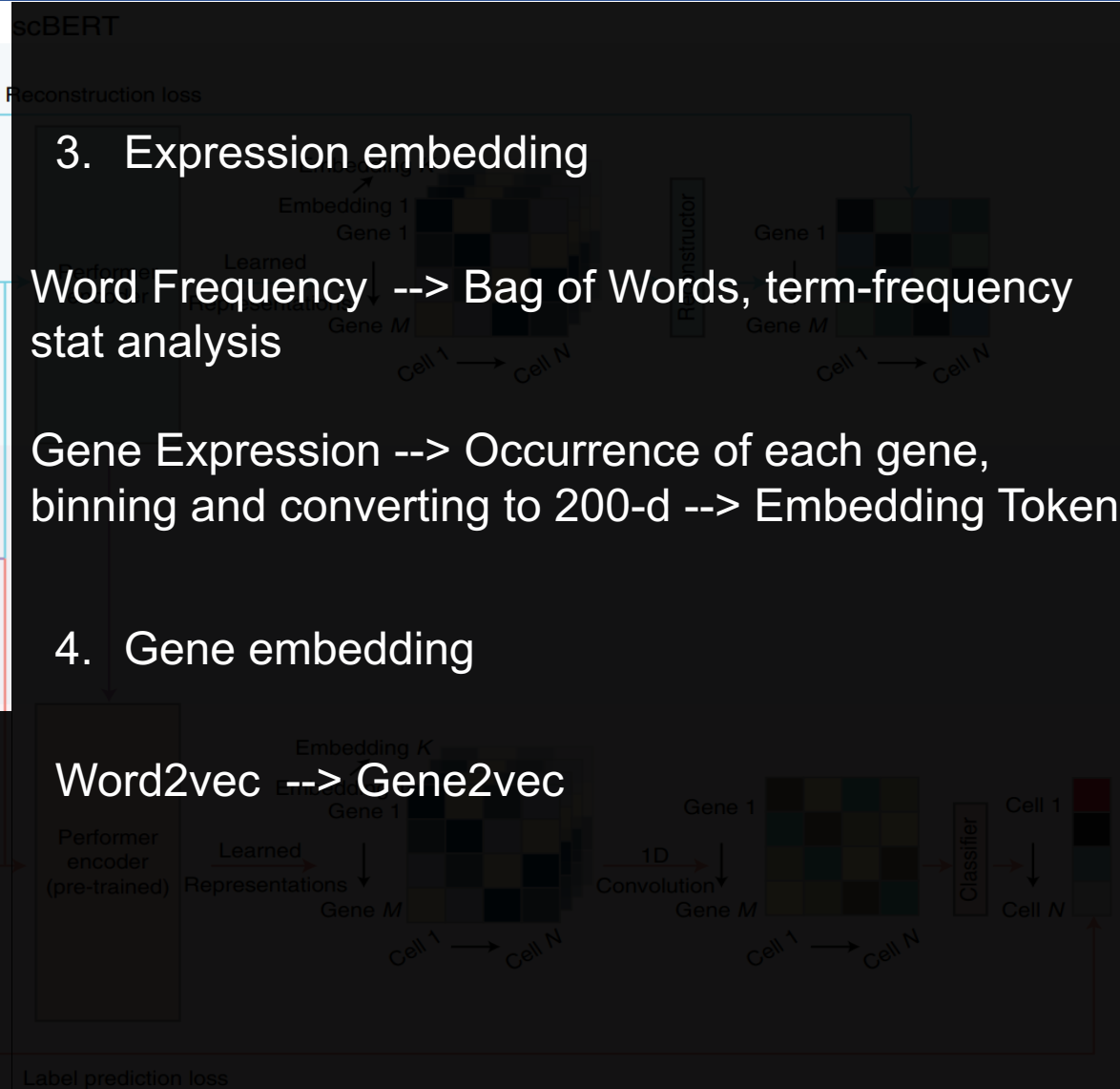
3. Expression embedding

Word Frequency --> Bag of Words, term-frequency stat analysis

Gene Expression --> Occurrence of each gene, binning and converting to 200-d --> Embedding Token

4. Gene embedding

Word2vec --> Gene2vec

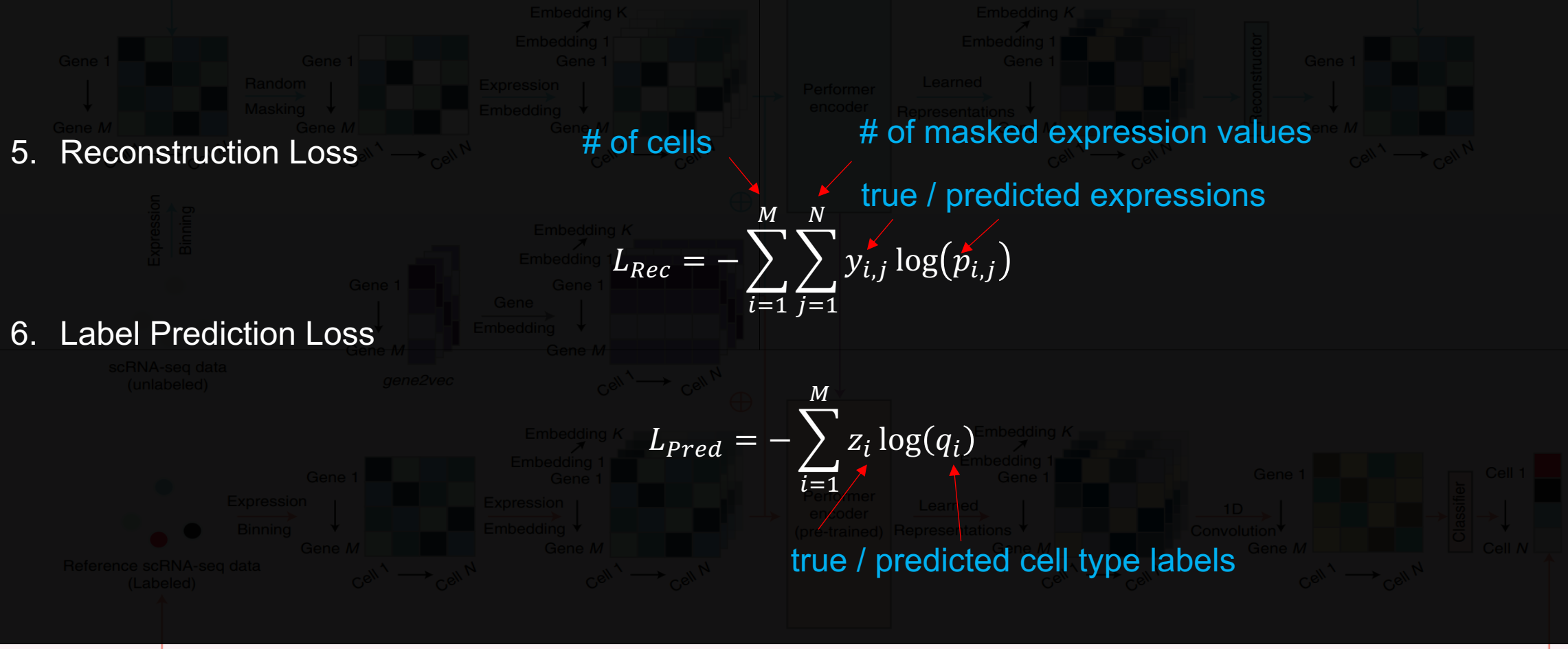


a

scBERT

Self-supervised pre-training

Reconstruction loss



Label prediction loss

Supervised finetuning

Input Embedding

$$\begin{bmatrix} C_1 \\ C_2 \\ \dots \\ C_N \end{bmatrix} \rightarrow \begin{bmatrix} g_1^{(2)} \\ g_2^{(2)} \\ \dots \\ g_M^{(2)} \end{bmatrix}$$

$$\begin{bmatrix} x_1 \\ x_2 \\ \dots \\ x_N \end{bmatrix} \quad \begin{bmatrix} id(g_1^{(2)}) \\ id(g_2^{(2)}) \\ \dots \\ id(g_M^{(2)}) \end{bmatrix} \quad \begin{bmatrix} t_{c,1} \\ t_{c,2} \\ \dots \\ t_{c,M} \end{bmatrix}$$

$$t_g^{(i)} = [id(g_1^{(i)}), id(g_2^{(i)}), \dots, id(g_M^{(i)})]$$

Gene Tokens

Maximum input length

Assign each gene a unique ID

Input gene (g) tokens of each cell i

$$t_c^{(i)} = [t_{c,1}^{(i)}, t_{c,2}^{(i)}, \dots, t_{c,M}^{(i)}]$$

Condition Tokens
(perturbations, ...)

scalar

$$x^{(i)} = [x_1^{(i)}, x_2^{(i)}, \dots, x_j^{(i)}, \dots, x_M^{(i)}]$$

Expression Values

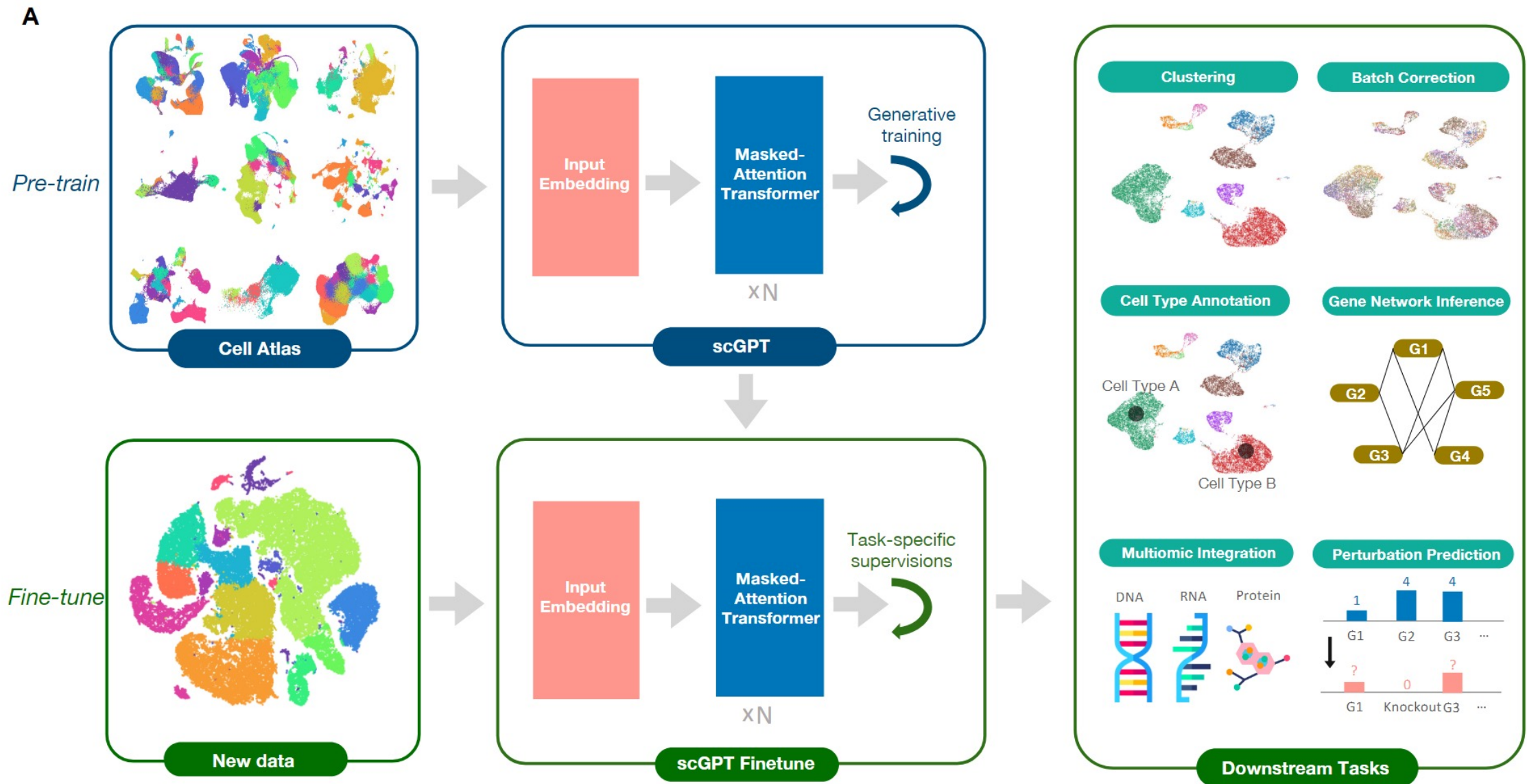
Value binning technique:

$$x_j^{(i)} = \begin{cases} k, & \text{if } X_{i,j} > 0 \text{ and } X_{i,j} \in [b_k, b_{k+1}] \\ 0, & \text{if } X_{i,j} = 0 \end{cases}$$

$$h^{(i)} = emb_g(t_g^{(i)}) + emb_x(x^{(i)}) + emb_c(t_c^{(i)})$$

Final Embedding

Fully Connected





Thank you.