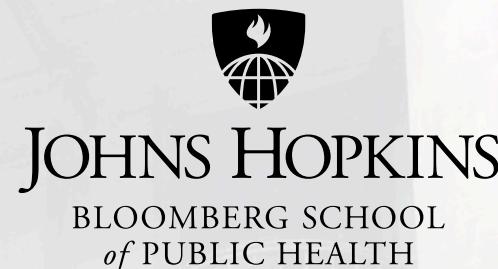


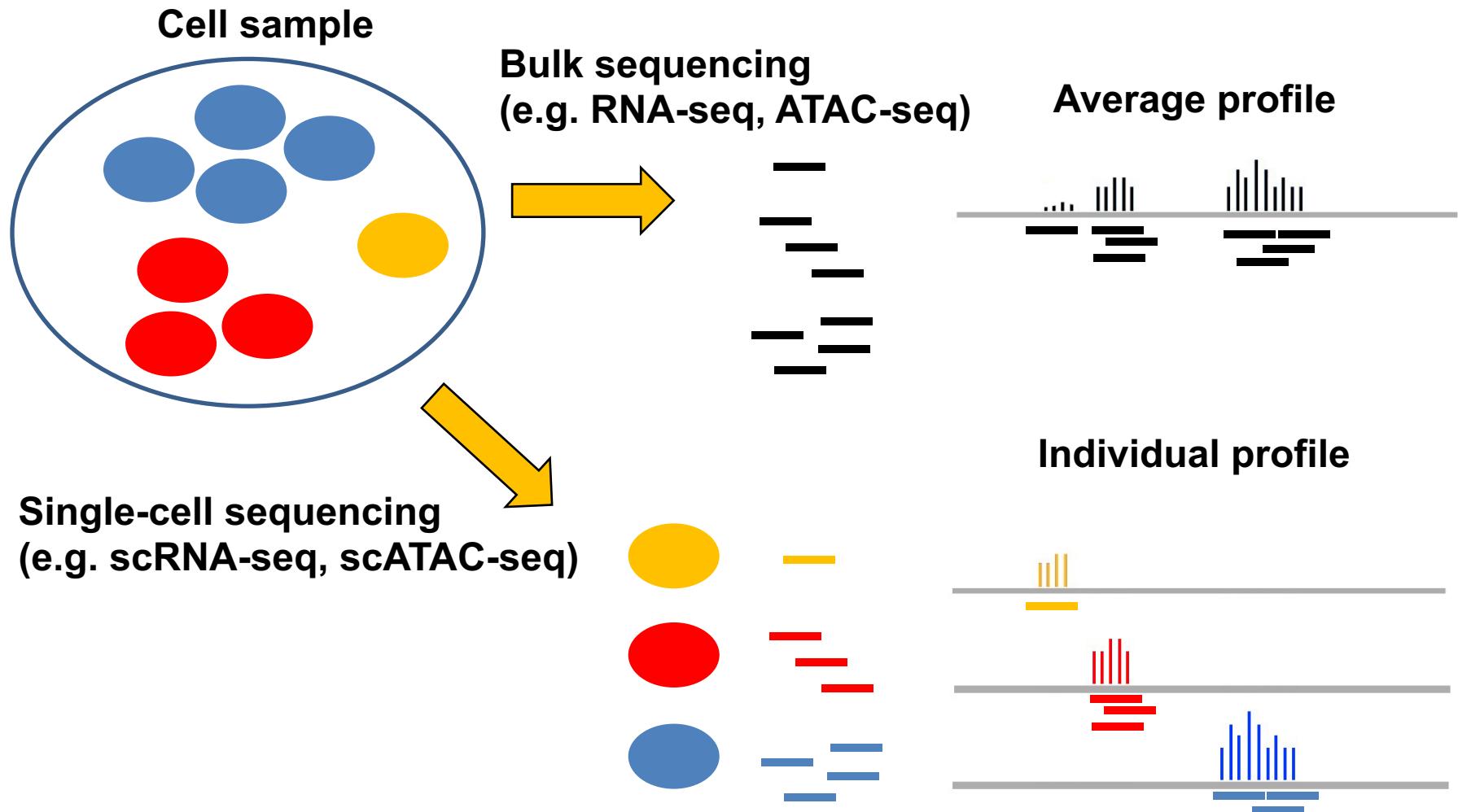
# Assign: a software for single-cell genomic data integration

Weiqiang Zhou

Assistant Scientist  
Department of Biostatistics



# From bulk to single cell



# Available single cell assays

## 10x genomics

- Single-cell RNA-seq
- Single-cell ATAC-seq
- Single cell immune profiling
- Single cell CNV

## Fluidigm

- Single-cell RNA-seq
- Single-cell ATAC-seq
- Single-cell DNA-seq



# Single cell multiomics

**Macaulay *et al.*** G&T-seq: parallel sequencing of single-cell genomes and transcriptomes. *Nature Methods* (2015). (**G&T-seq**: genomic DNA and gene expression)

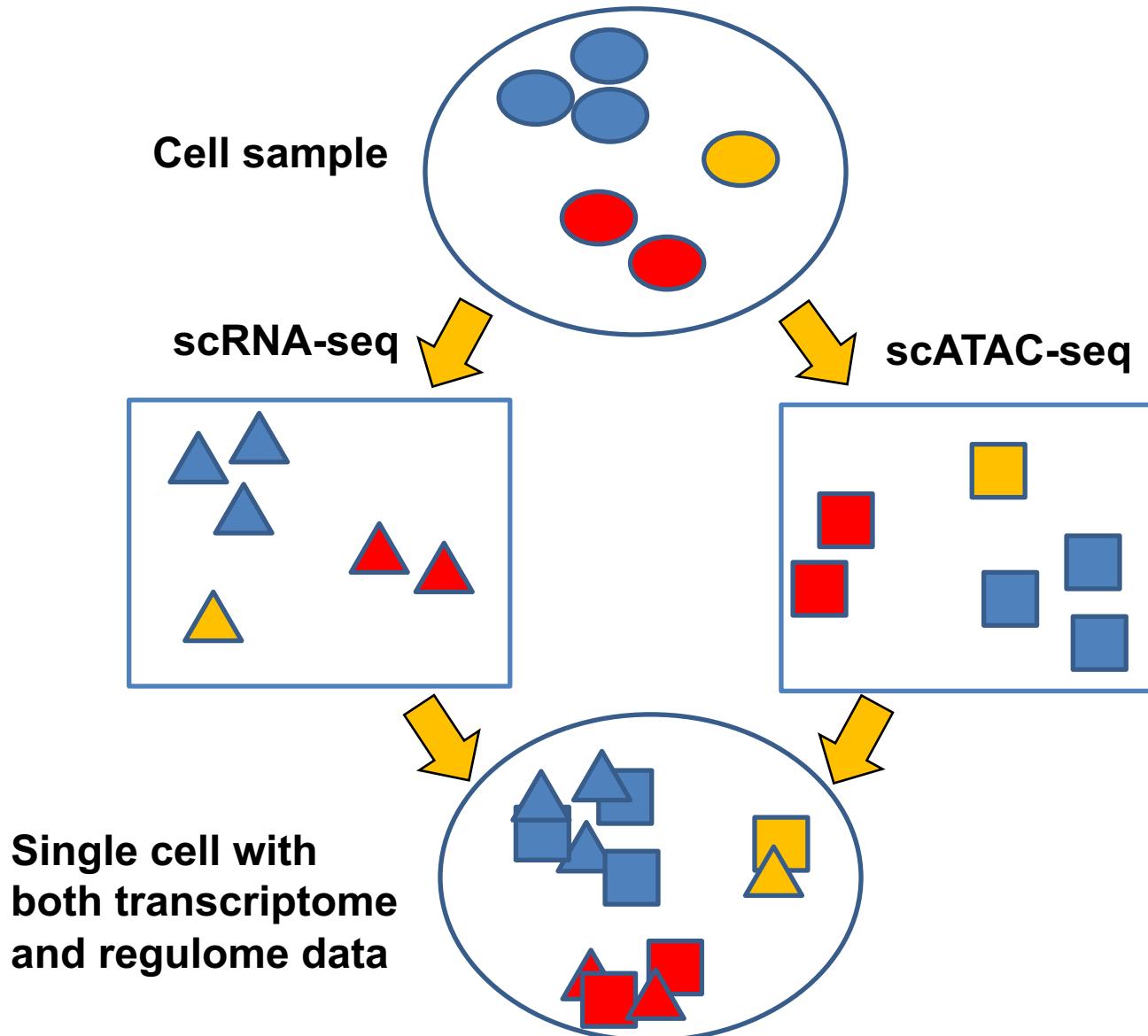
**Hu *et al.*** Simultaneous profiling of transcriptome and DNA methylome from a single cell. *Genome Biology* (2016). (**scMT-seq**: DNA methylation and gene expression)

**Cheow *et al.*** Single-cell multimodal profiling reveals cellular epigenetic heterogeneity. *Nature Methods* (2016). (**sc-GEM**: genotype, DNA methylation, and gene expression; use Fluidigm system)

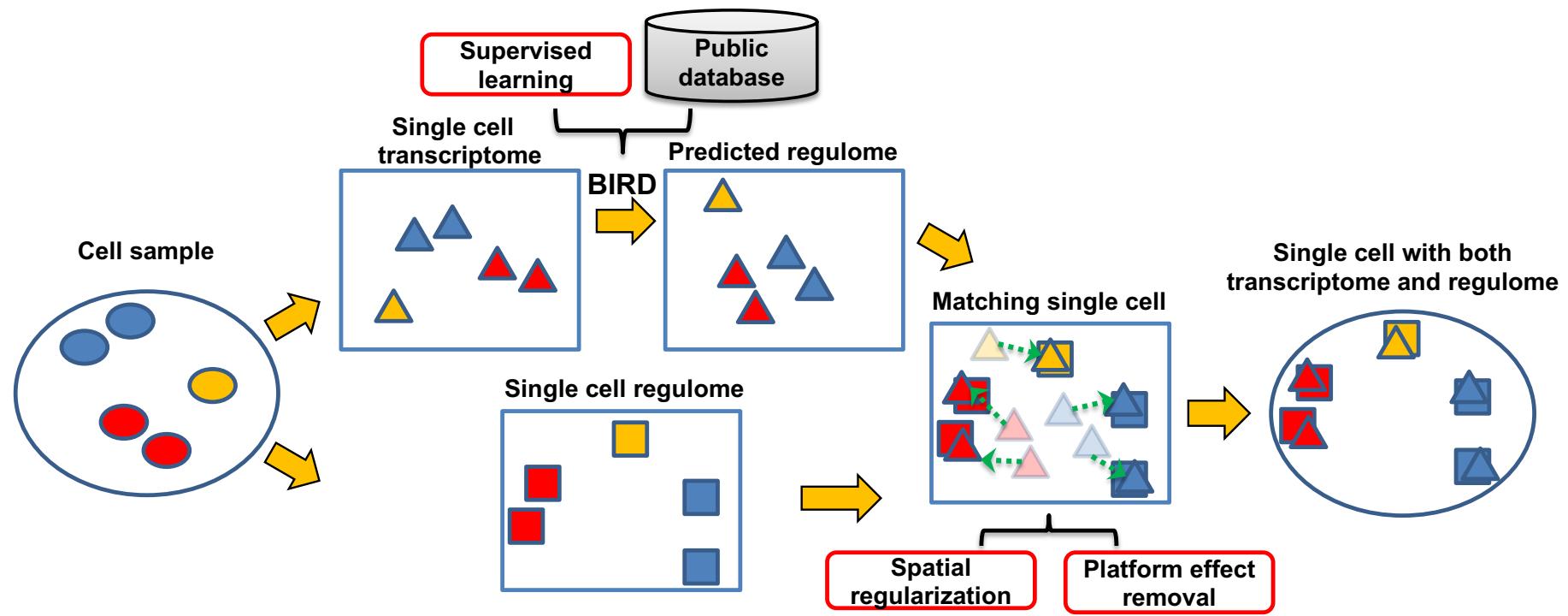
**Clark *et al.*** scNMT-seq enables joint profiling of chromatin accessibility DNA methylation and transcription in single cells. *Nature Communications* (2018). (**scNMT-seq**: chromatin accessibility, DNA methylation, and gene expression)

**Cao, J. *et al.*** Joint profiling of chromatin accessibility and gene expression in thousands of single cells. *Science* (2018). (**sci-CAR**: chromatin accessibility, and gene expression)

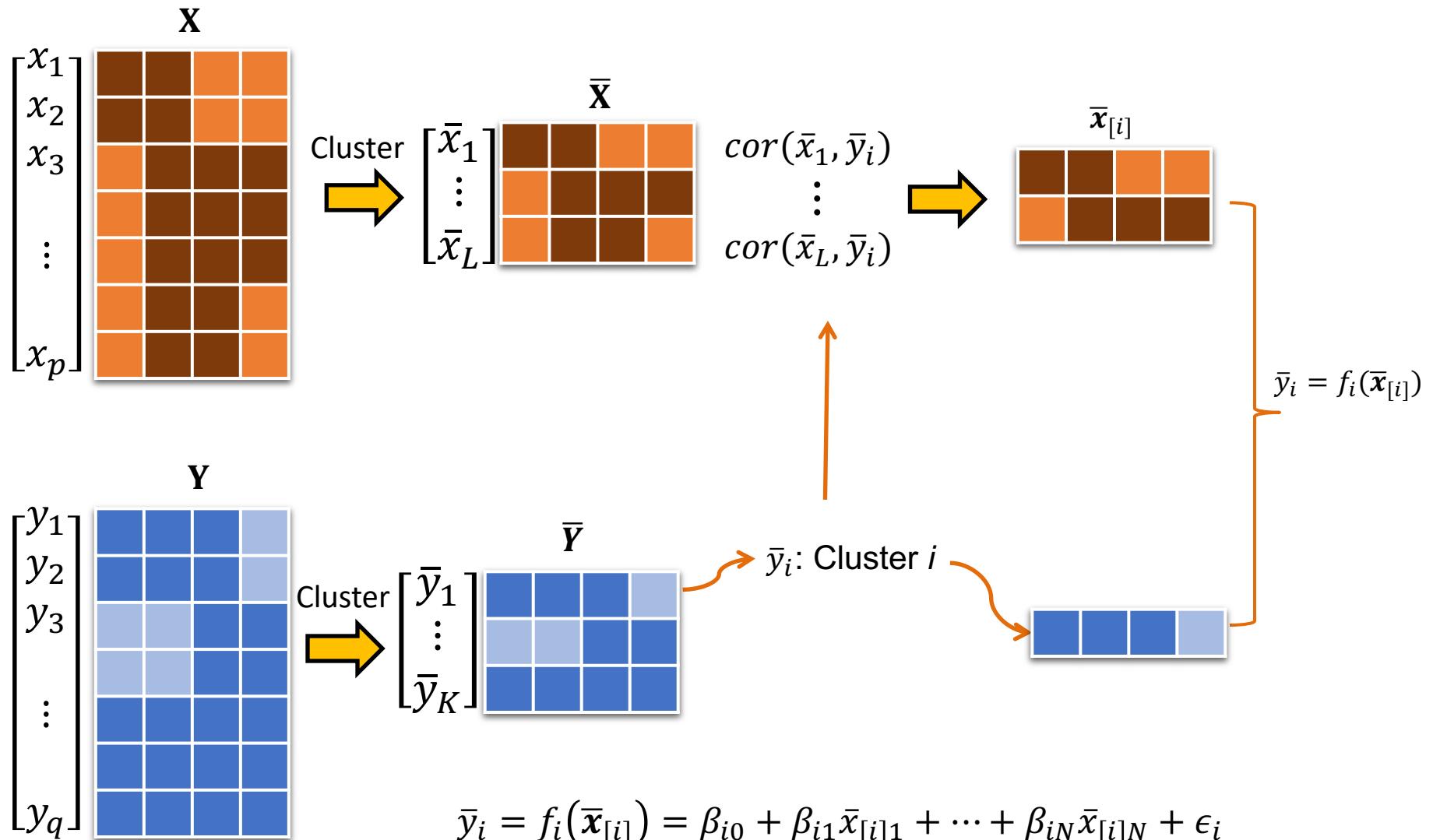
# Integrate multiple data types



# Overview



# BIRD



# Platform effect removal

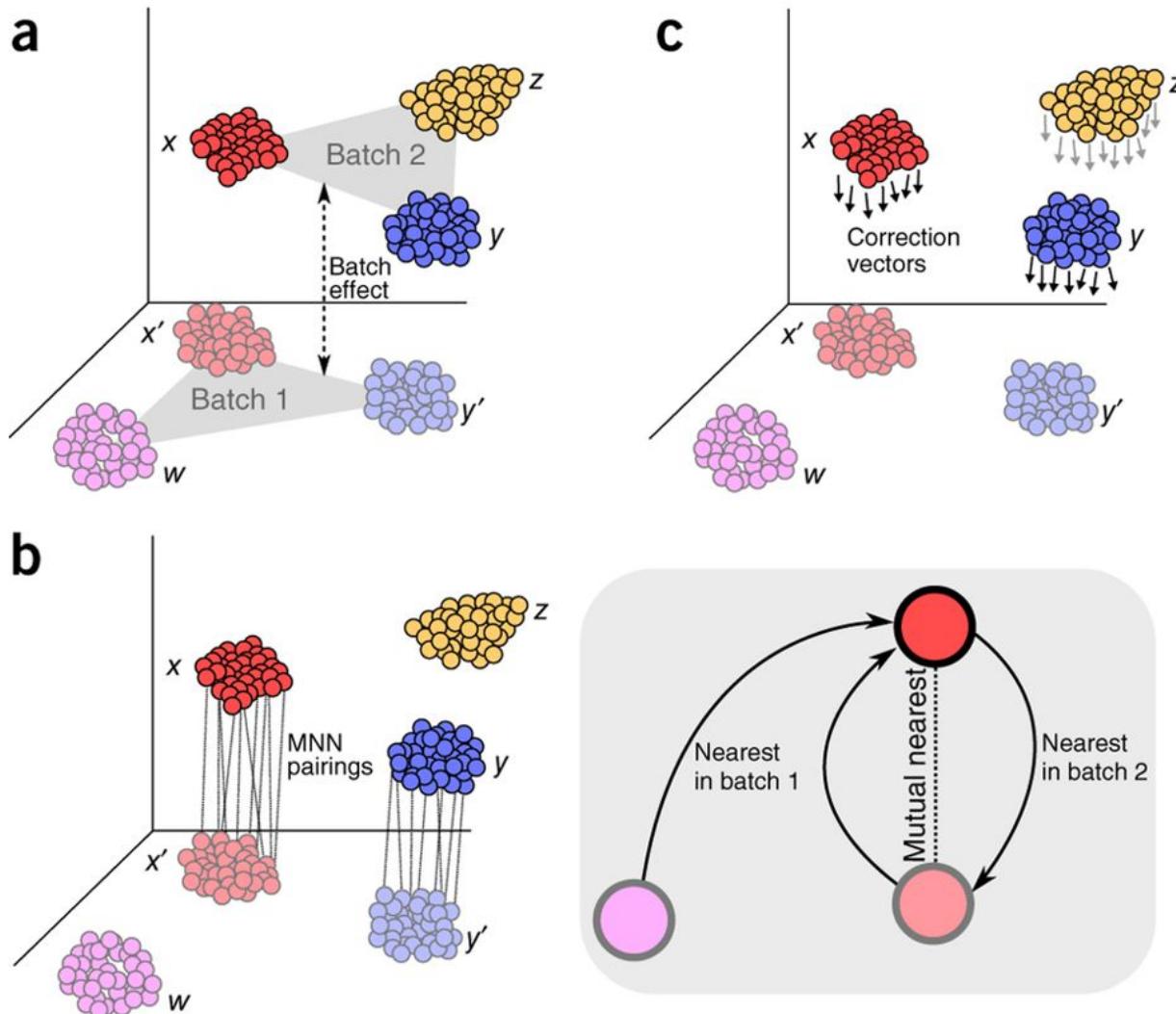


Figure obtained from Haghverdi *et al.* Batch effects in single-cell RNA-sequencing data are corrected by matching mutual nearest neighbors. *Nature Biotechnology* **36**, 421–427 (2018)

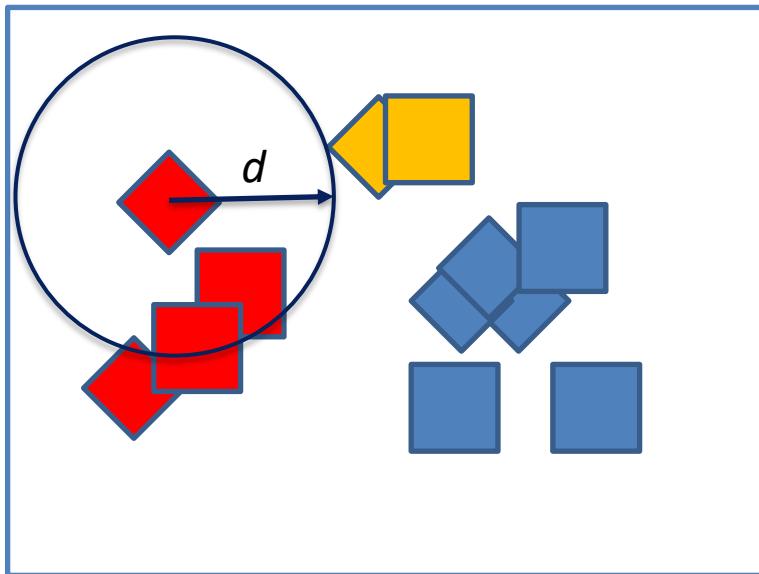


# Spatial regulation

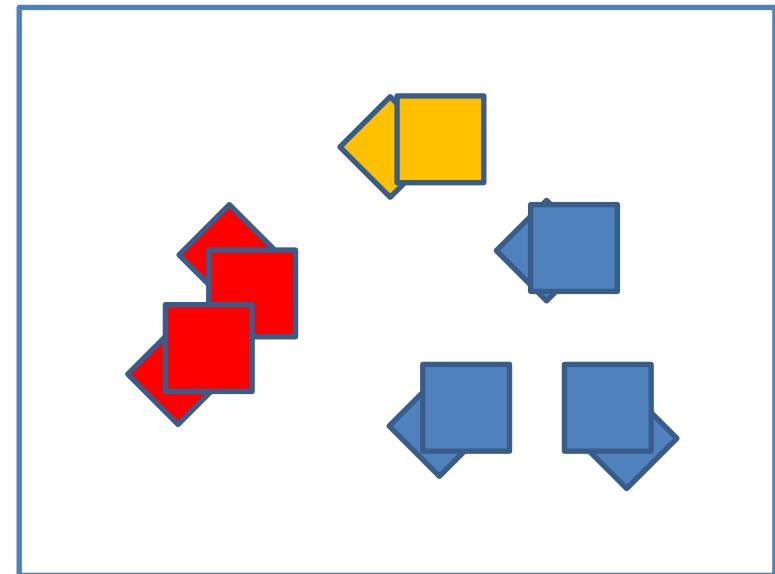
 Predicted scATAC

 True scATAC

Before optimization



After optimization

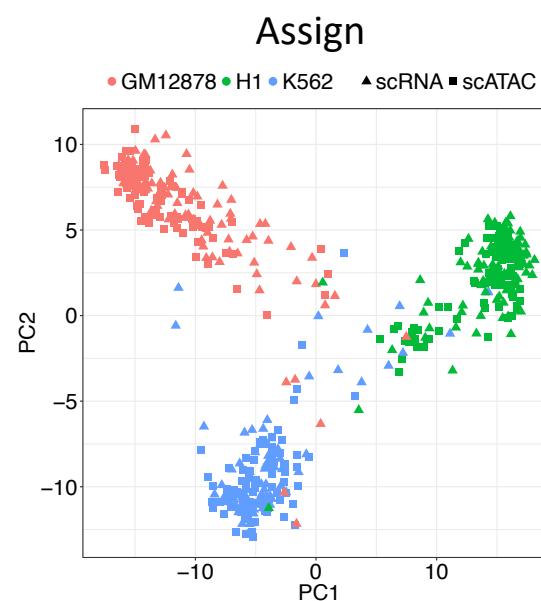
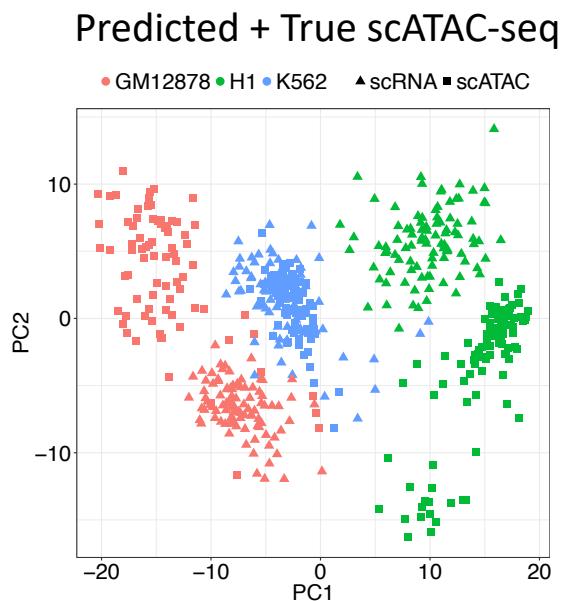
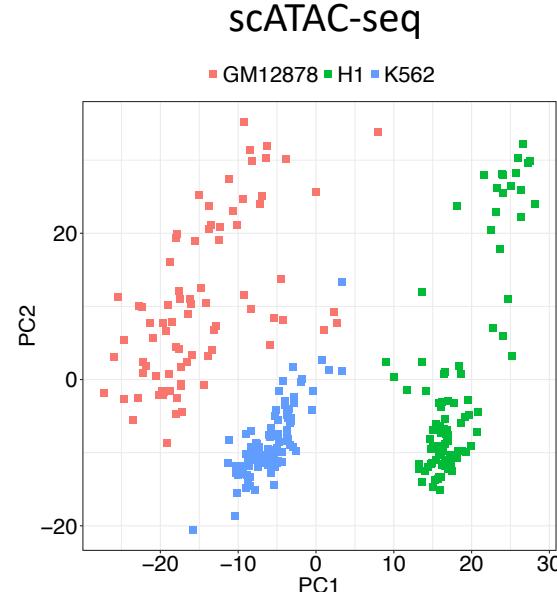
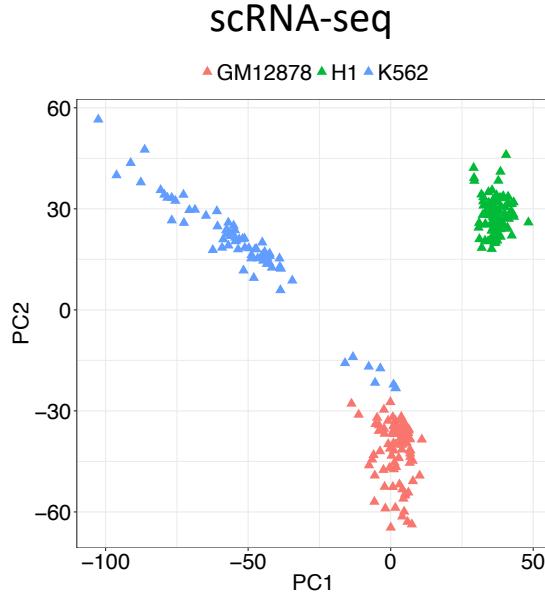


Fisher's extract test:

$$N_{P\_obs} / N_{T\_obs} = N_P / N_T$$

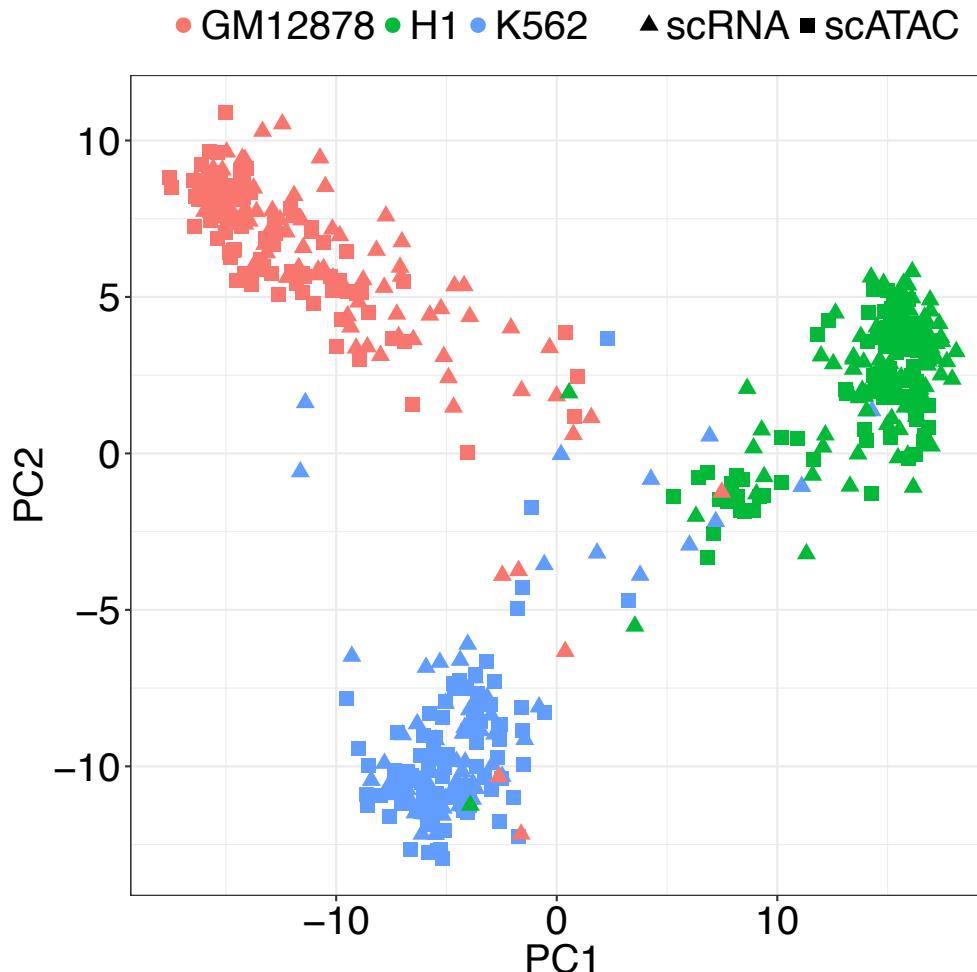


# Application to a test dataset



# Comparison with other methods

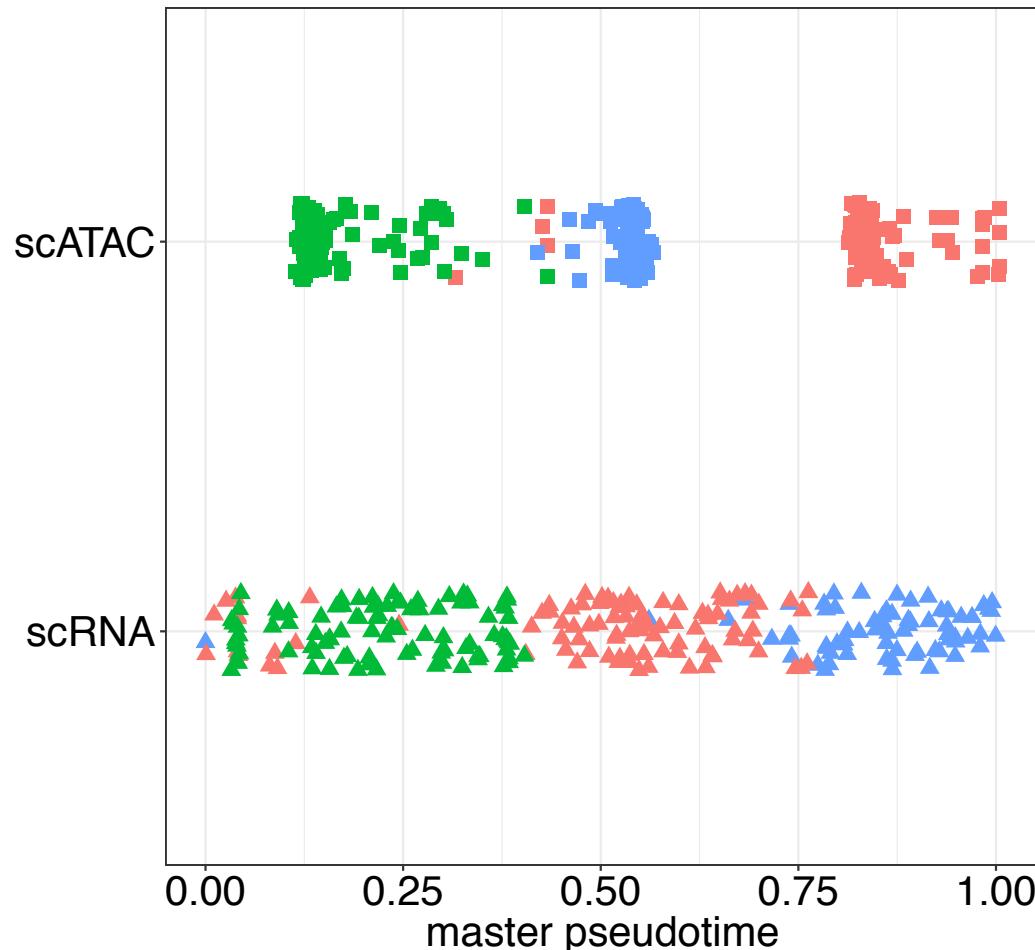
ASSIGN: 94% correct pairs



# Comparison with other methods

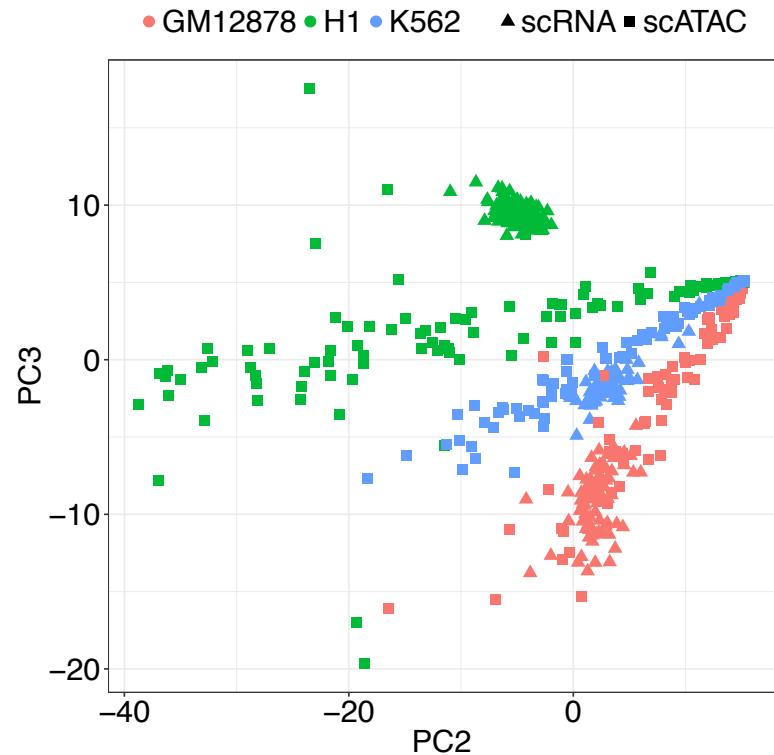
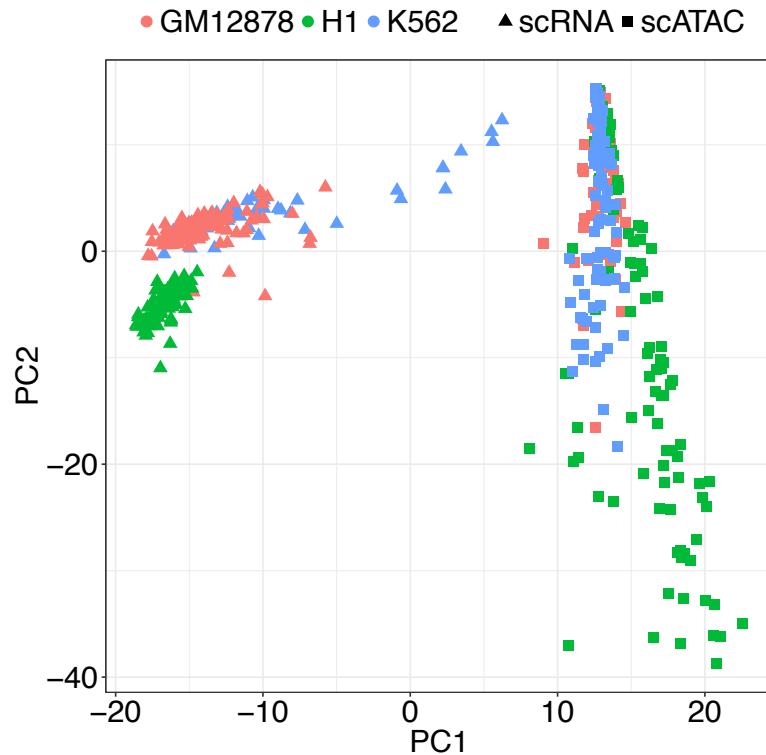
MATCHER: 43% correct pairs

• GM12878 • H1 • K562 ▲ scRNA ■ scATAC



# Comparison with other methods

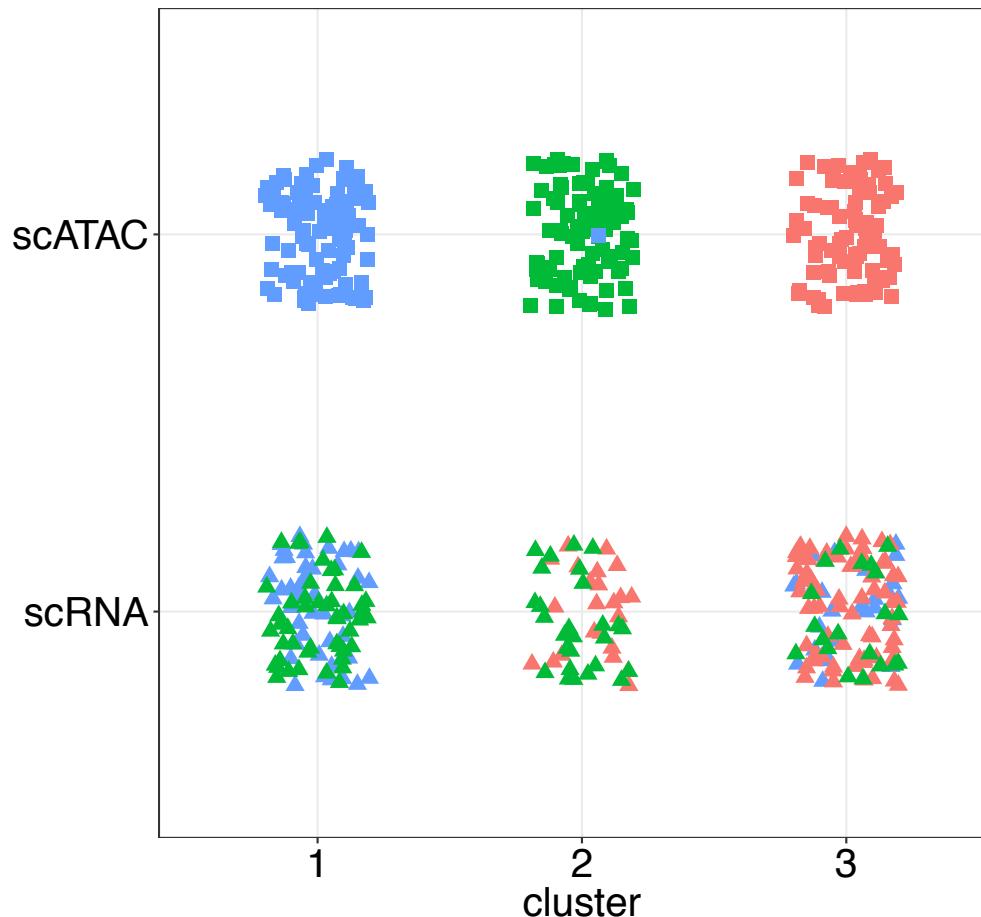
Cicero: 79% correct pairs



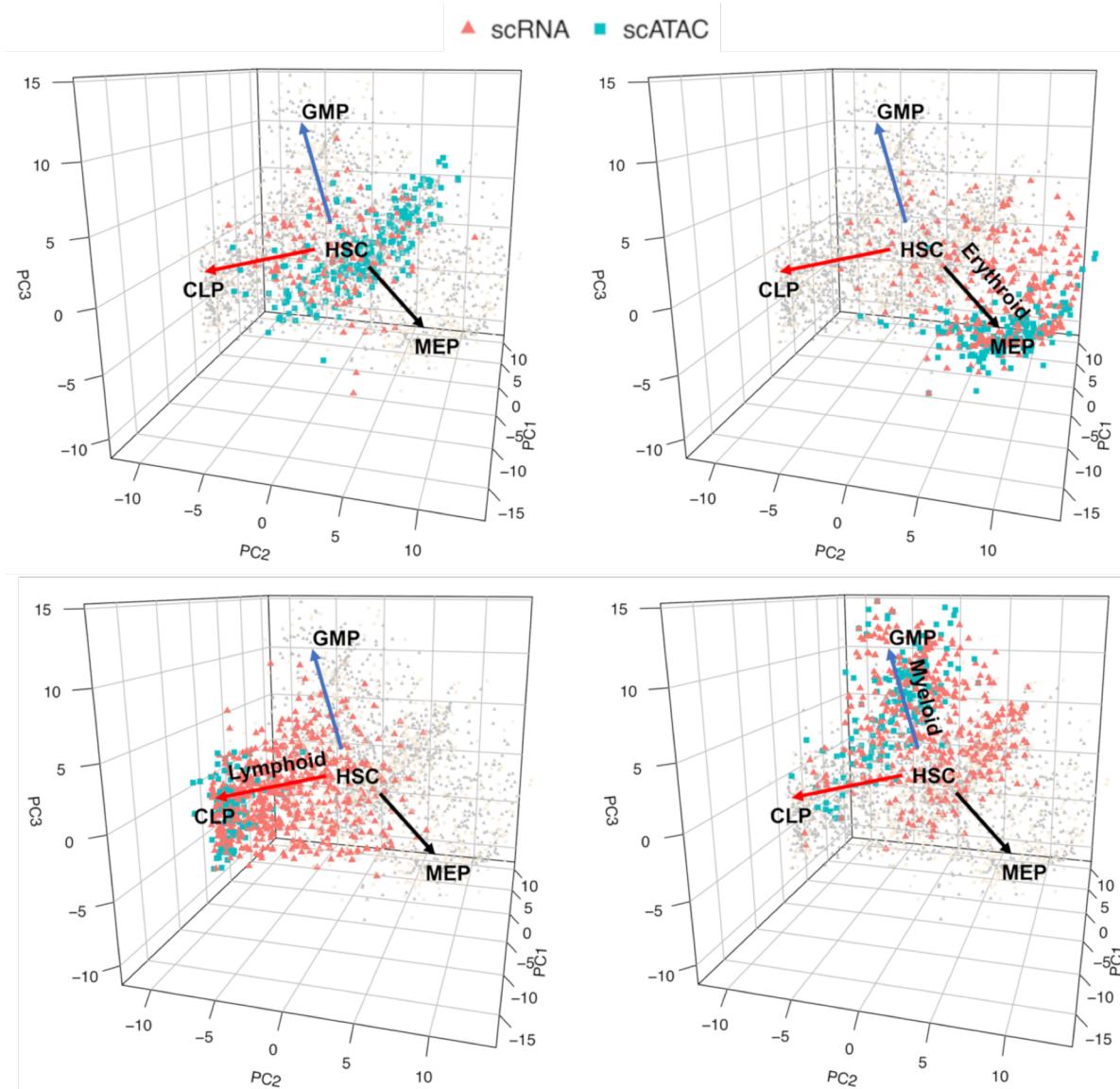
# Comparison with other methods

couple-NMF: 56% correct pairs

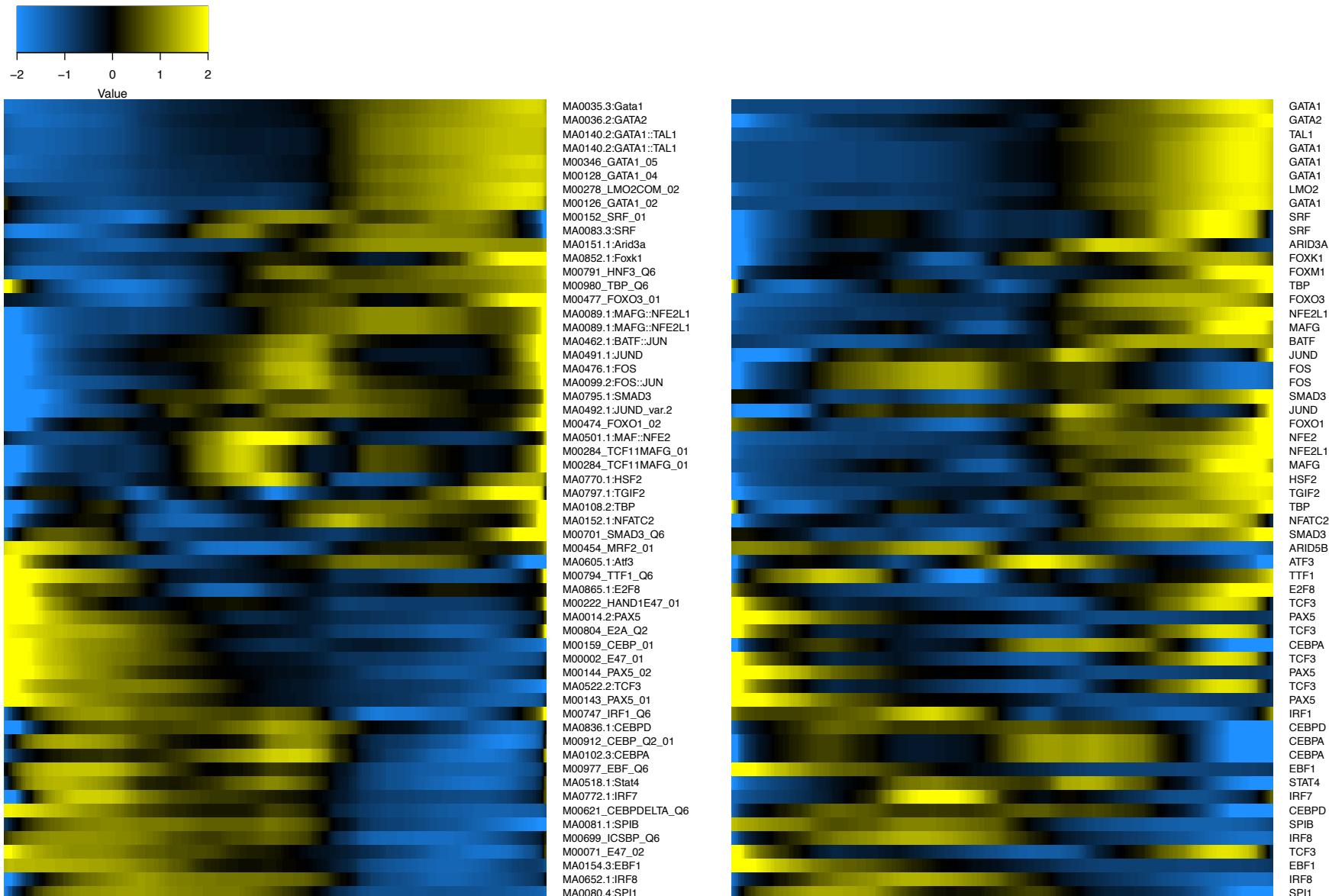
● GM12878 ● H1 ● K562 ▲ scRNA ■ scATAC



# Application to HCA data

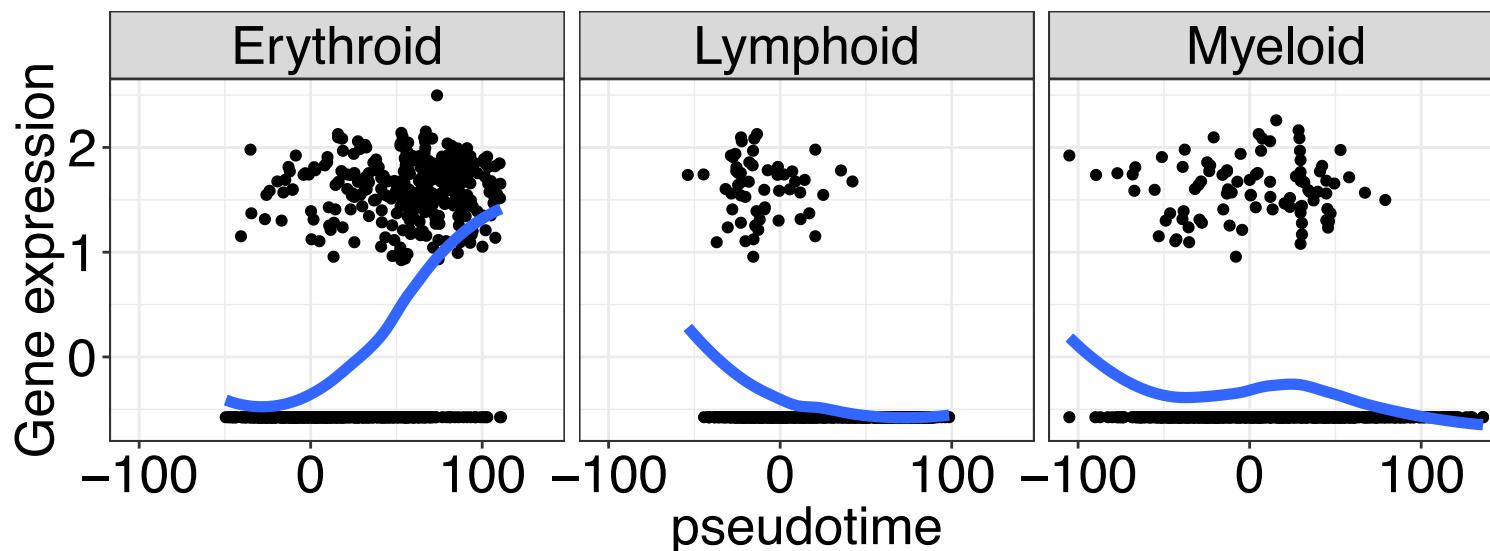
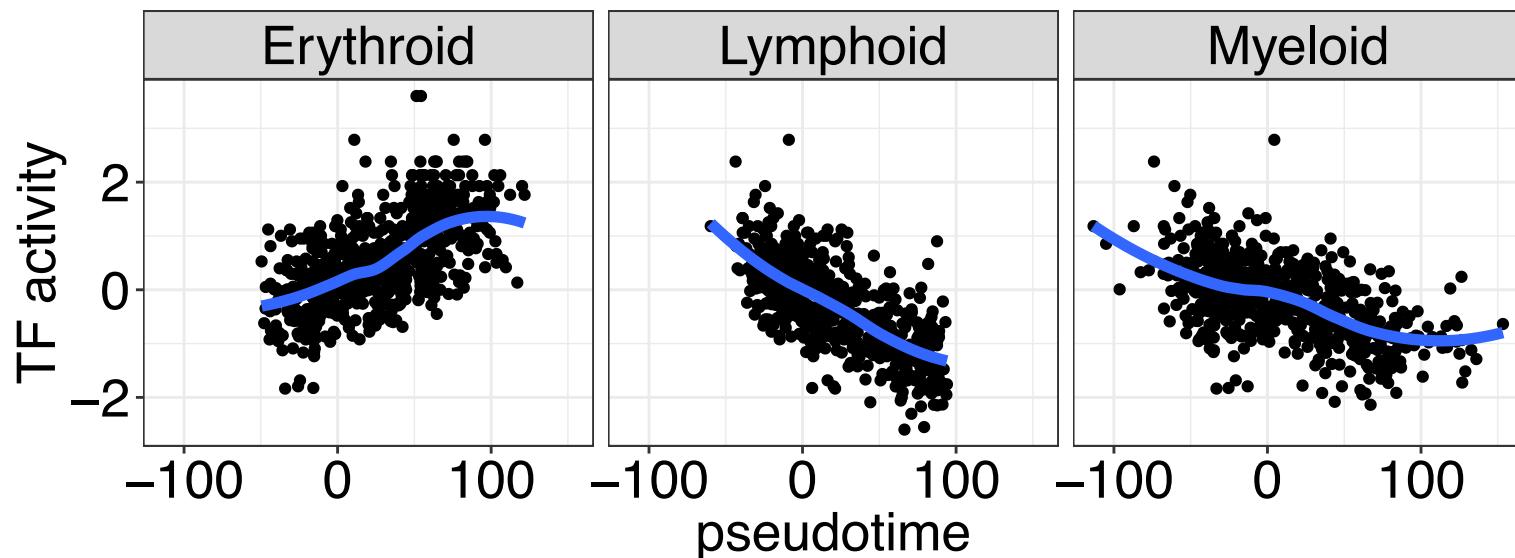


# Application to HCA data



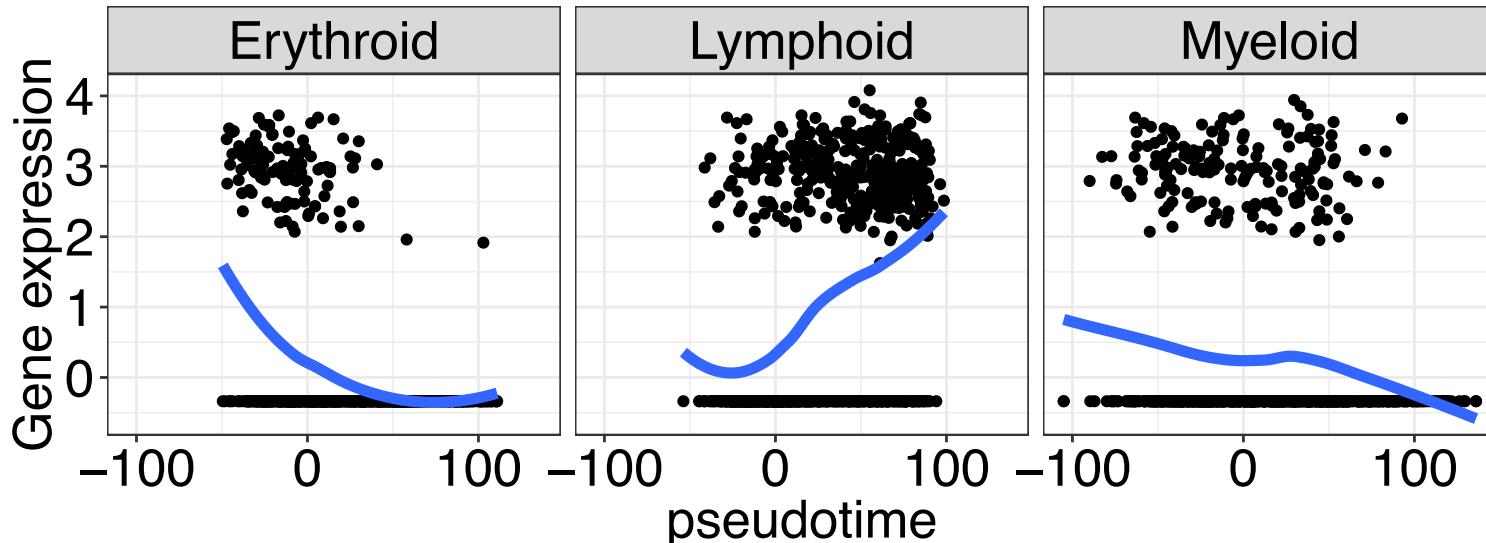
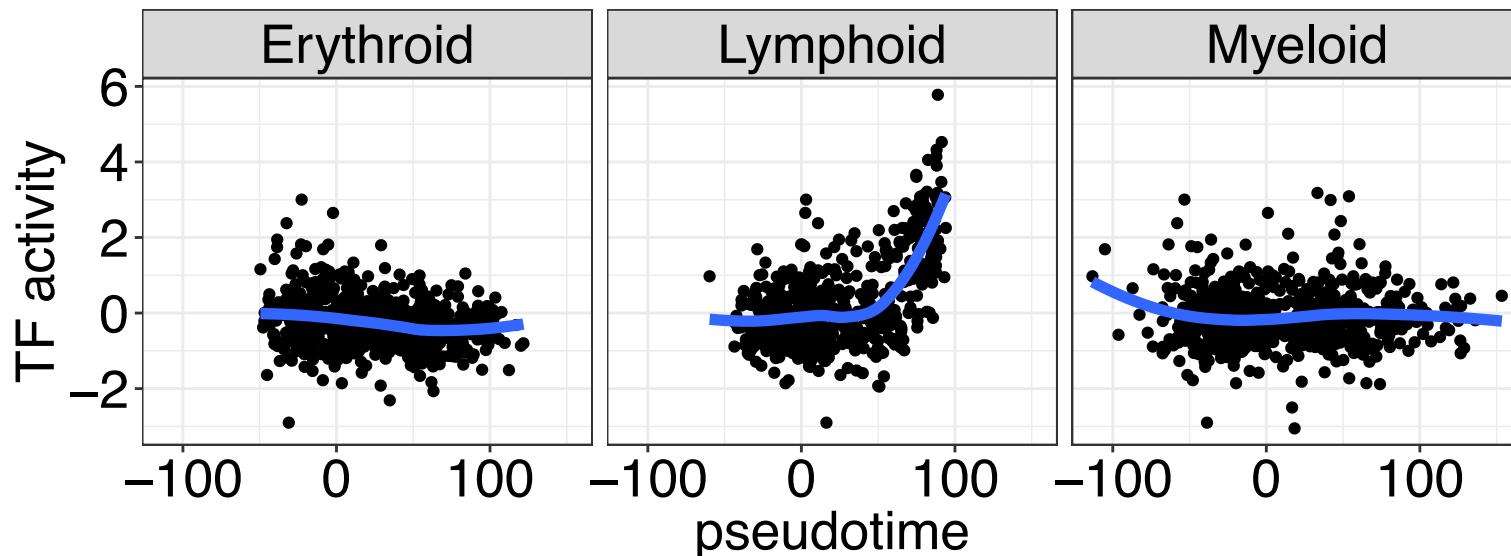
# Application to HCA data

GATA1



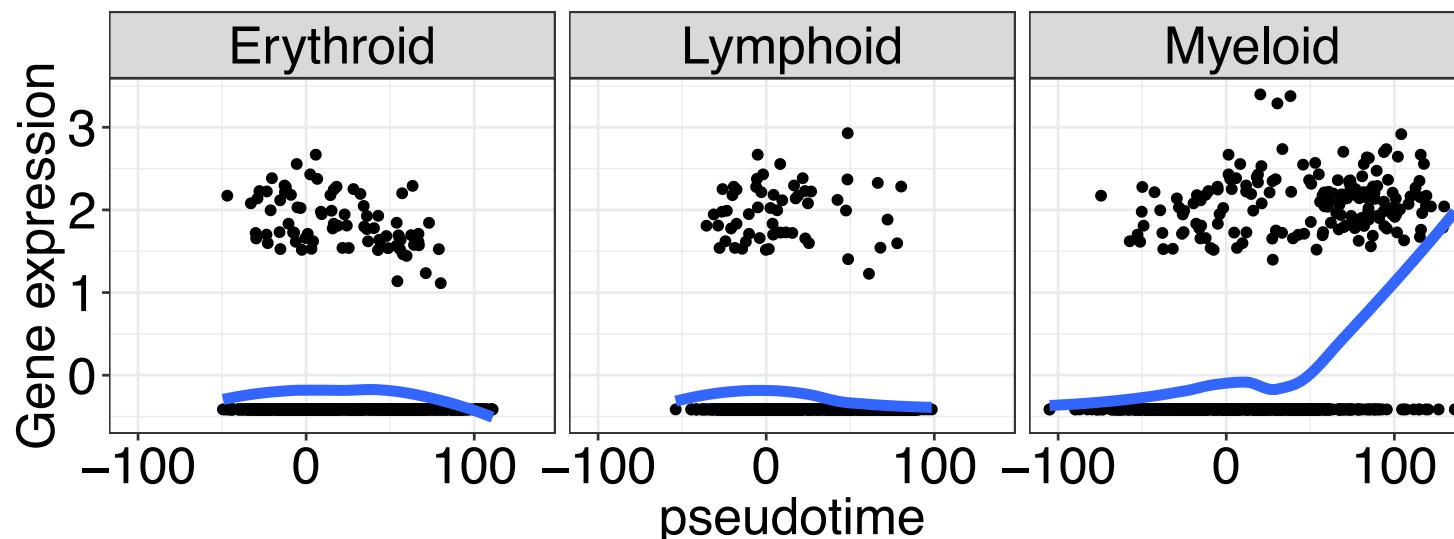
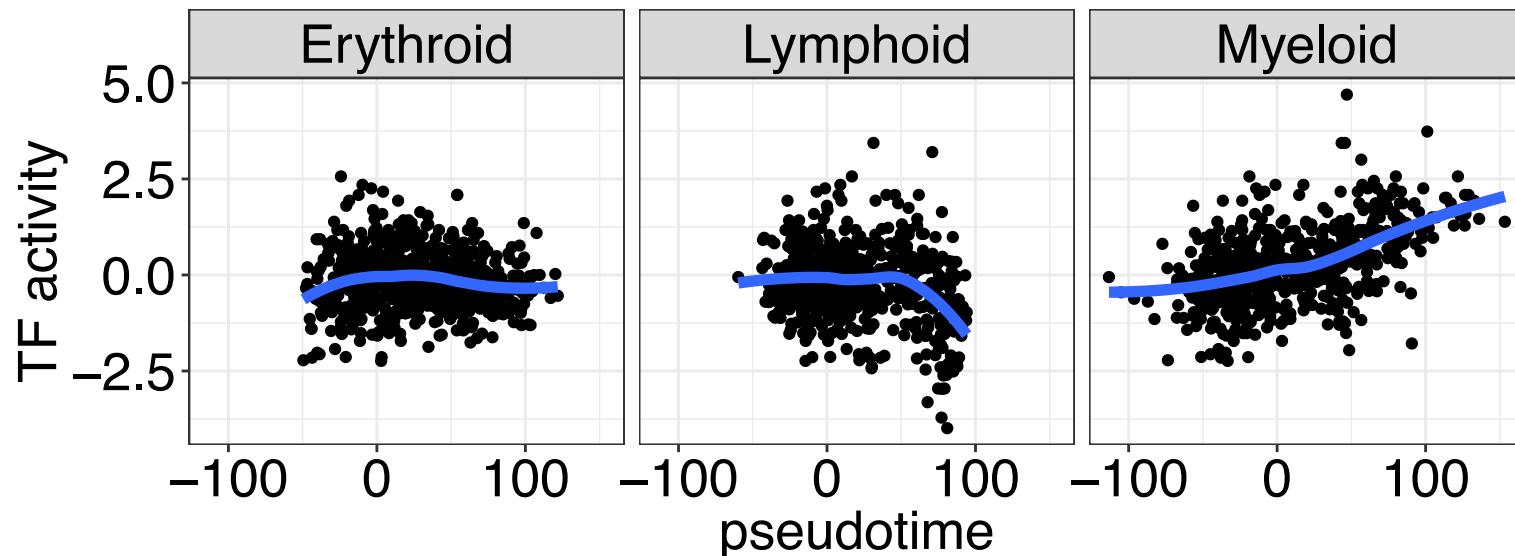
# Application to HCA data

EBF1



# Application to HCA data

CEBPD



# Acknowledgement



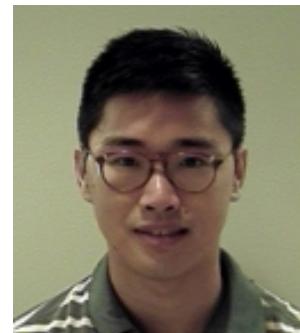
**Hongkai Ji**

Professor,  
Department of Biostatistics,  
Johns Hopkins Bloomberg  
School of Public Health



**Zhicheng Ji**

PhD Candidate,  
Department of Biostatistics,  
Johns Hopkins Bloomberg  
School of Public Health



**Weixiang Fang**

PhD Candidate,  
Department of Biostatistics,  
Johns Hopkins Bloomberg  
School of Public Health

