

# Statistical Methods for High-dimensional Biology



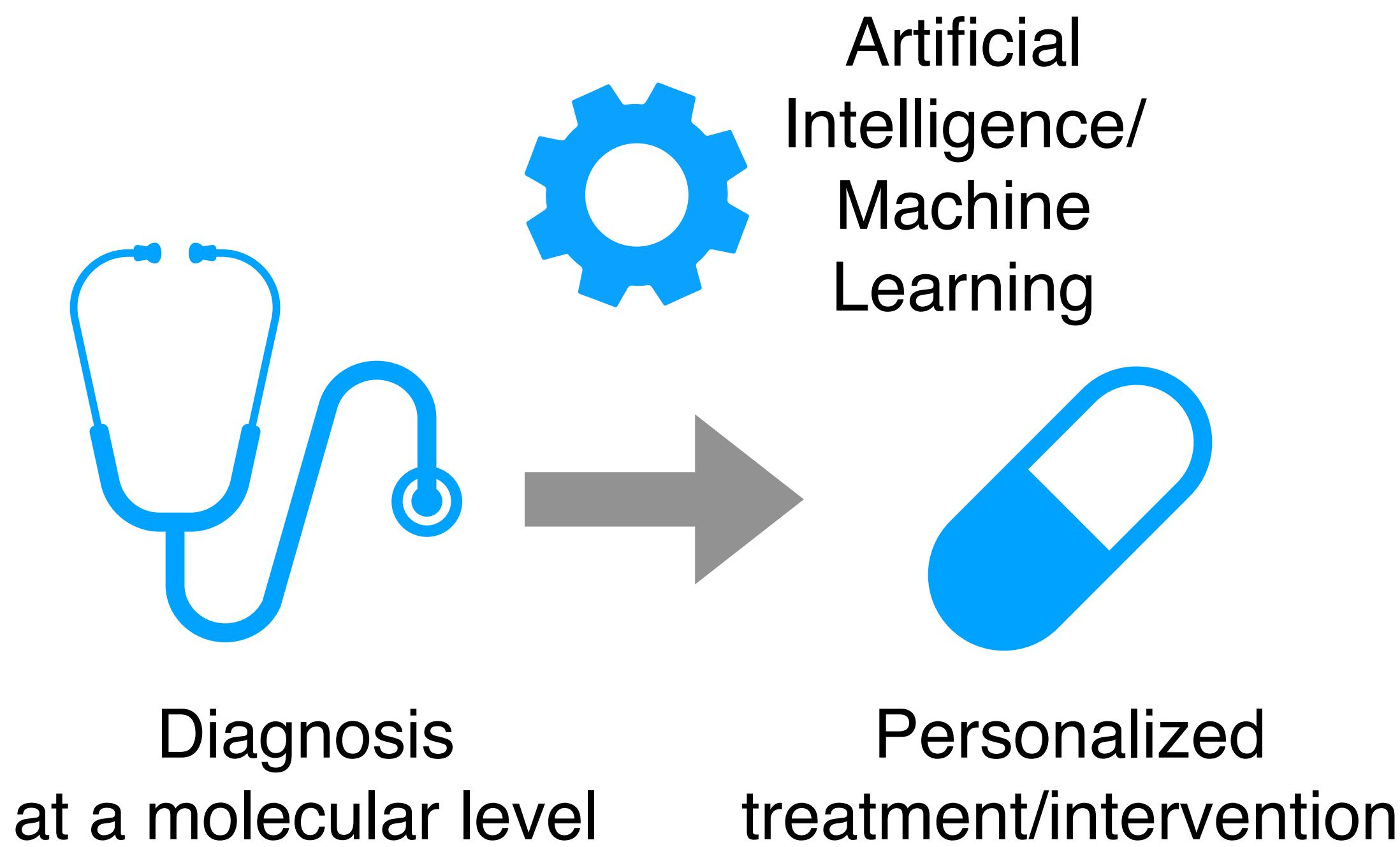
## Multomics data integration

Yongjin Park, UBC Path&Lab, STAT, BC Cancer

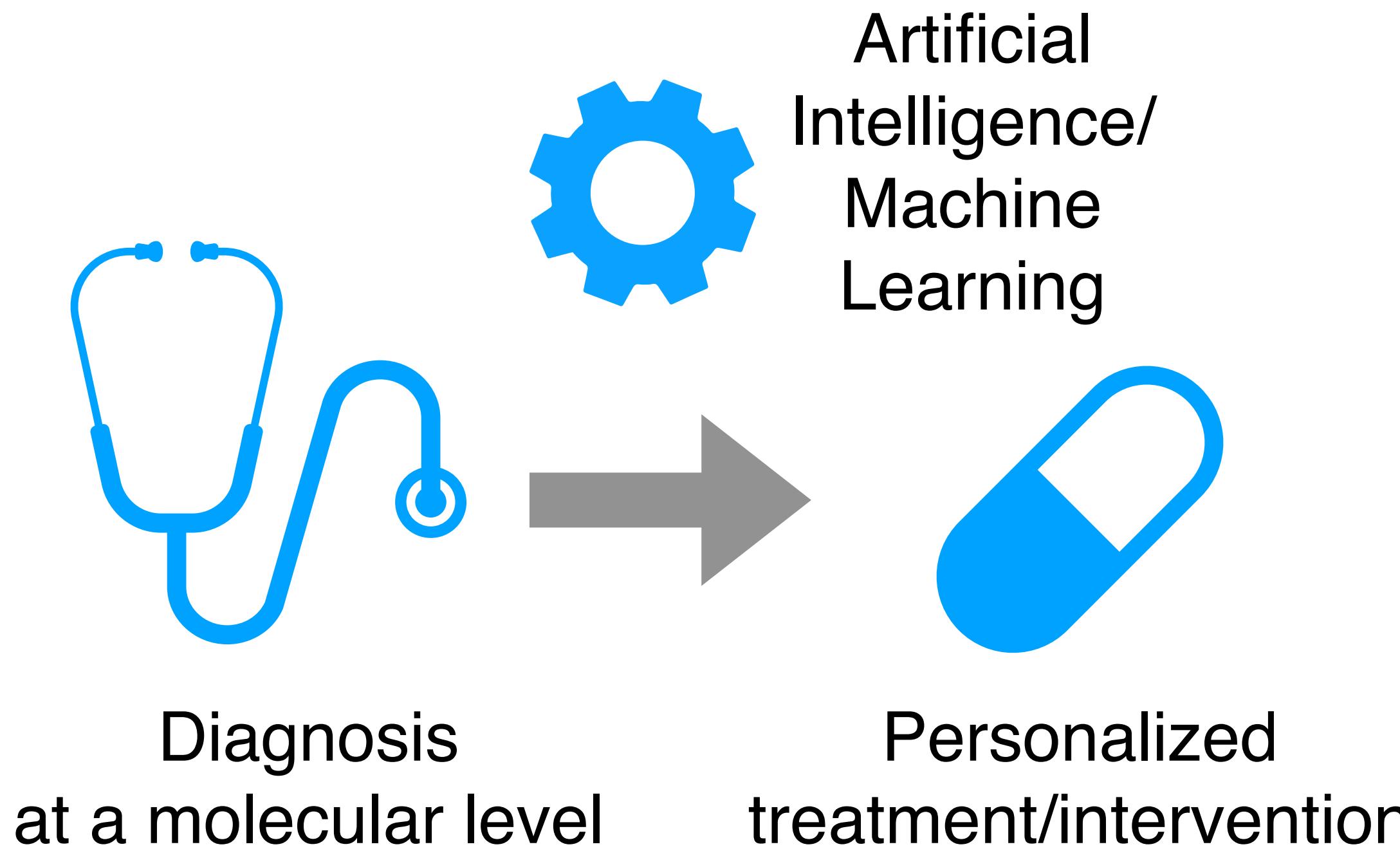
# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms

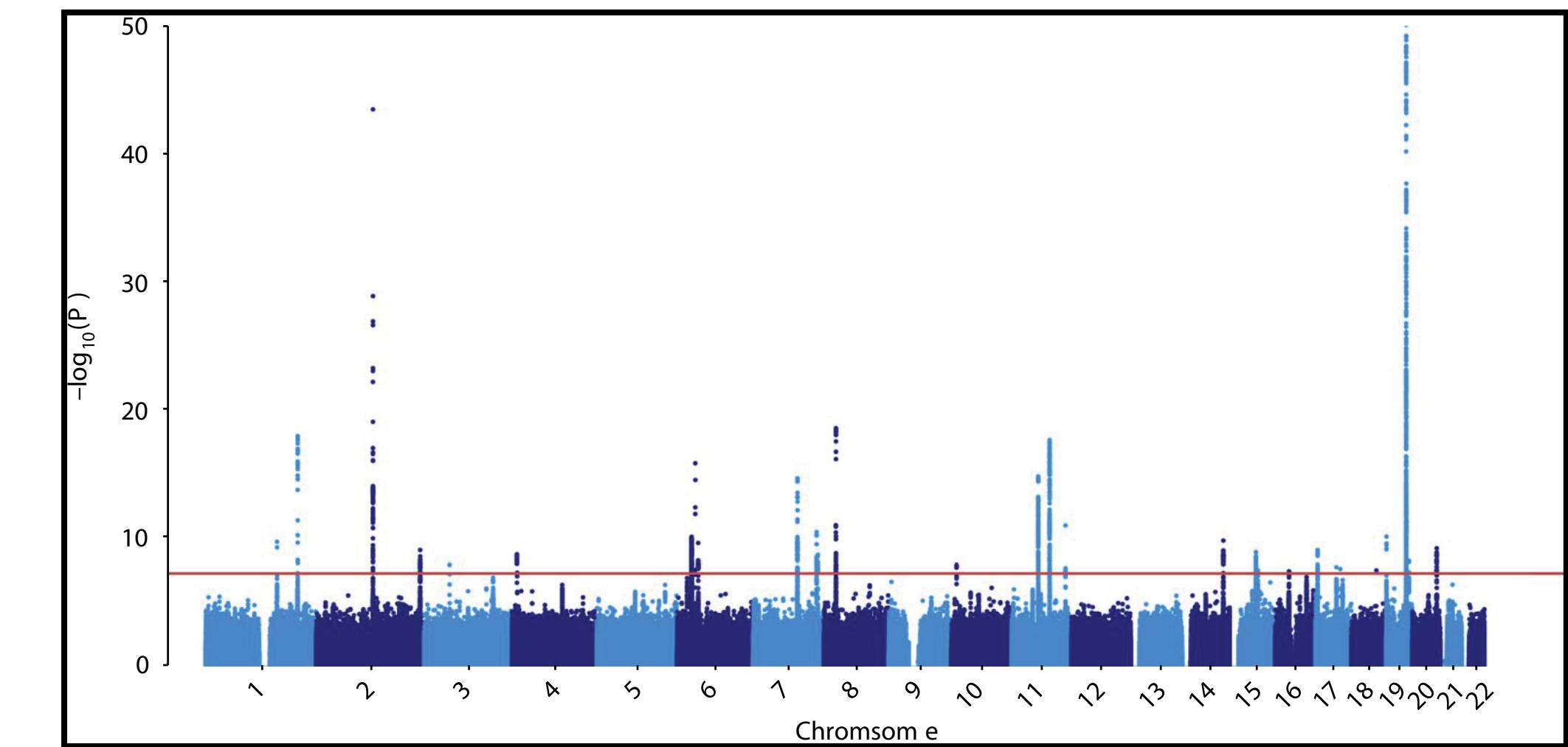
# Genomic Medicine: Promises and Challenges



# Genomic Medicine: Promises and Challenges

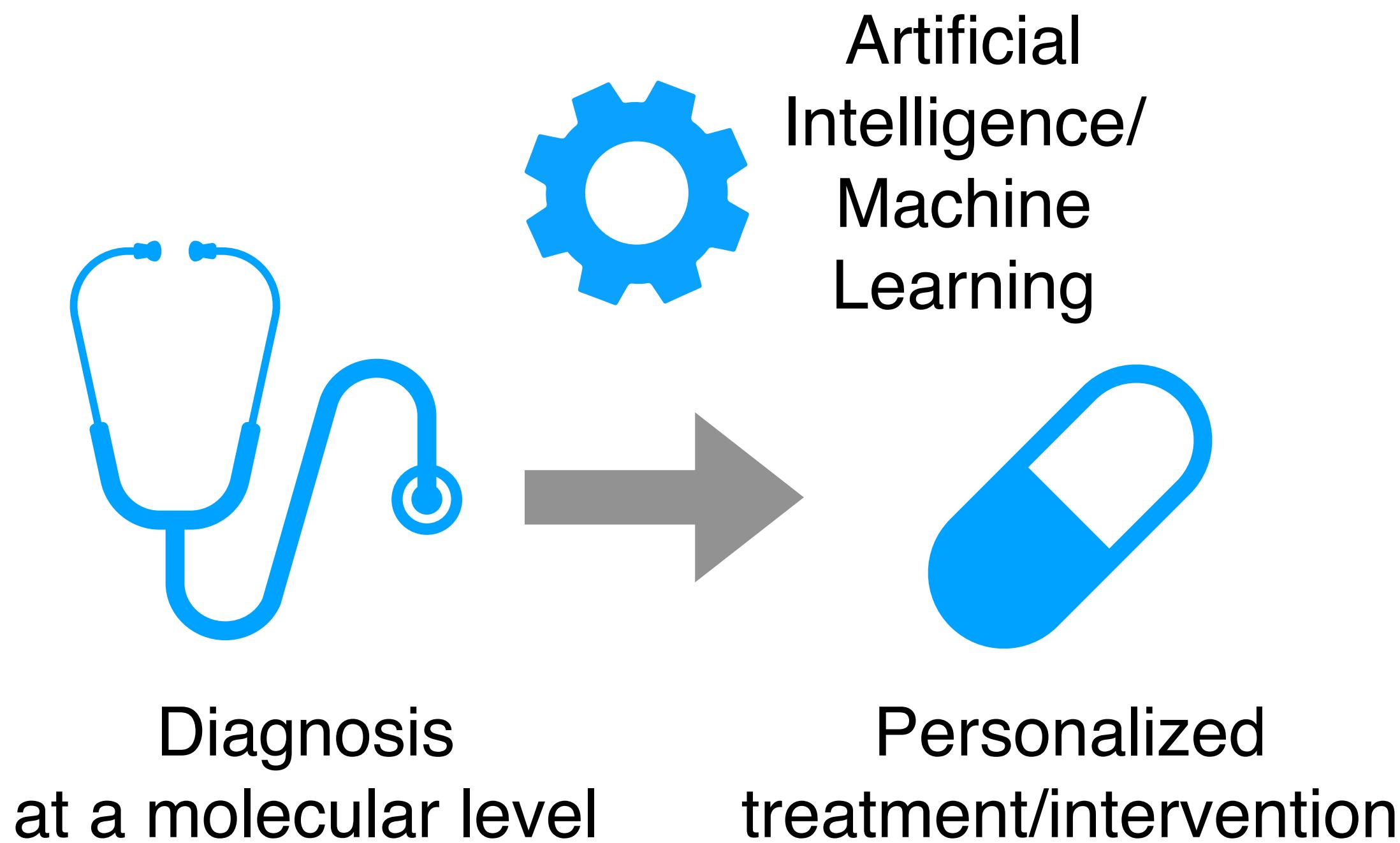


**Genome-wide association studies  
of complex disease traits**

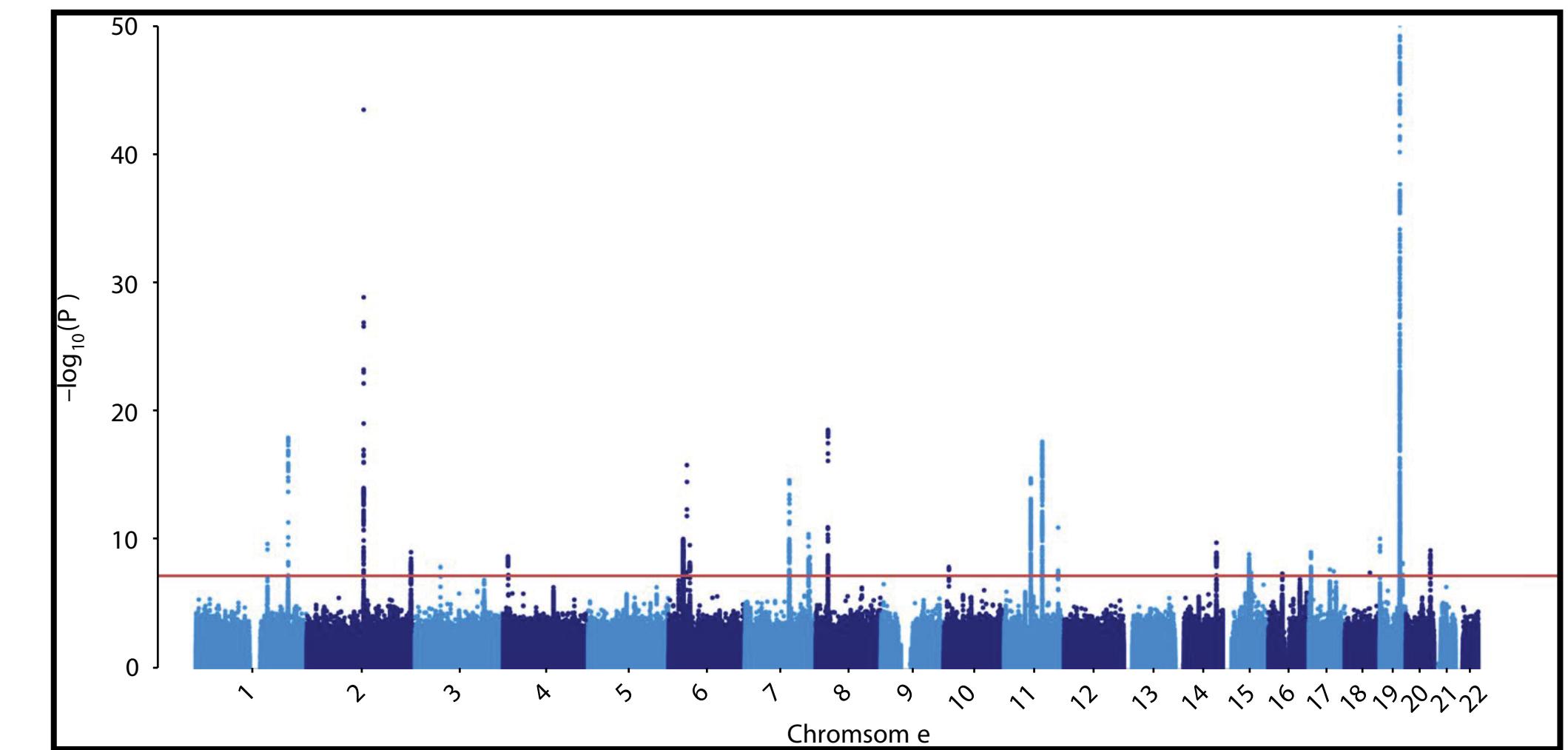


**Lectures 13, 14, 15**

# Genomic Medicine: Promises and Challenges



- 90% of disease hits are difficult to interpret.
- Causal variants & target genes, pathways are unknown.
- Target cell types are unknown.

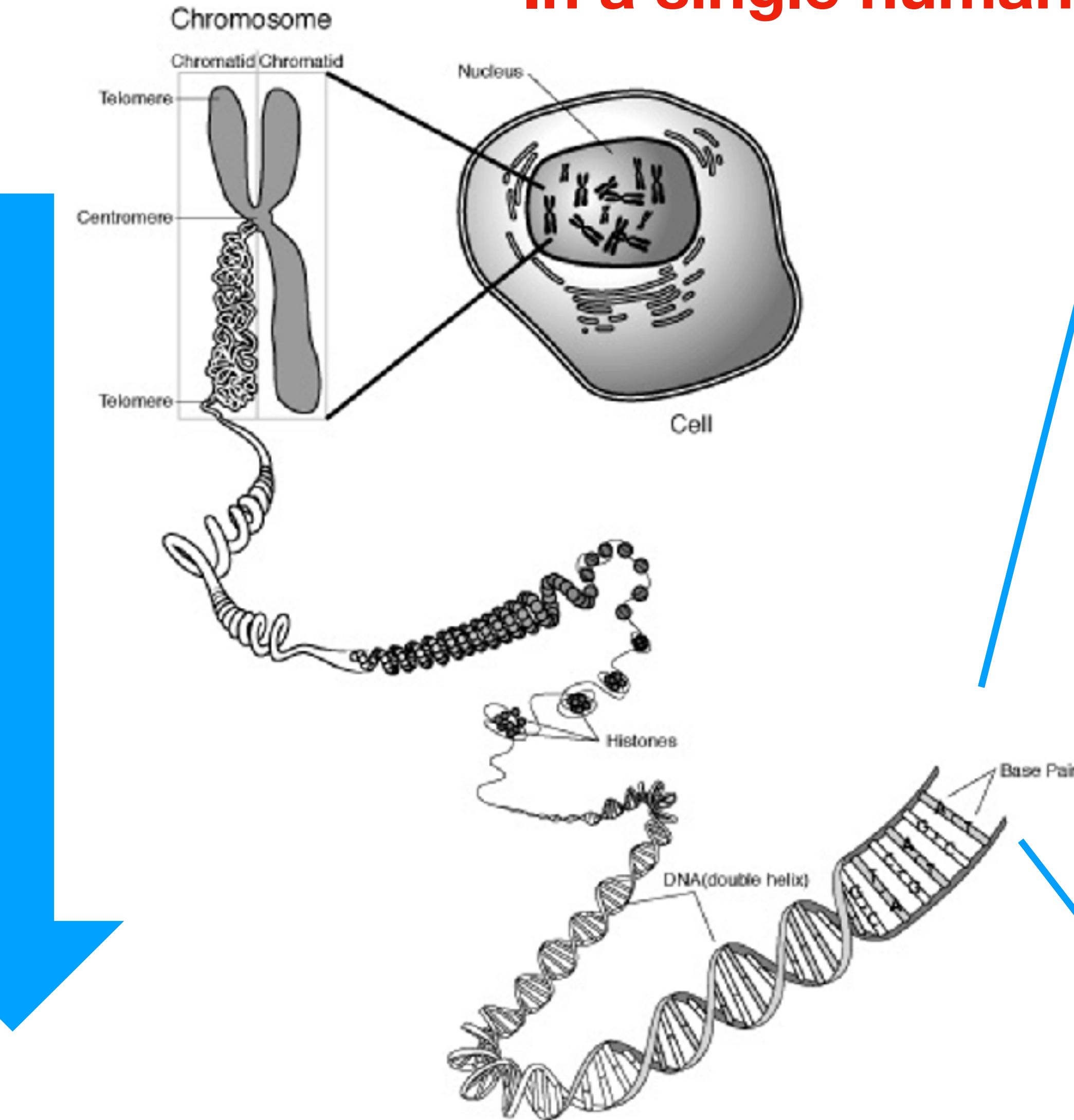
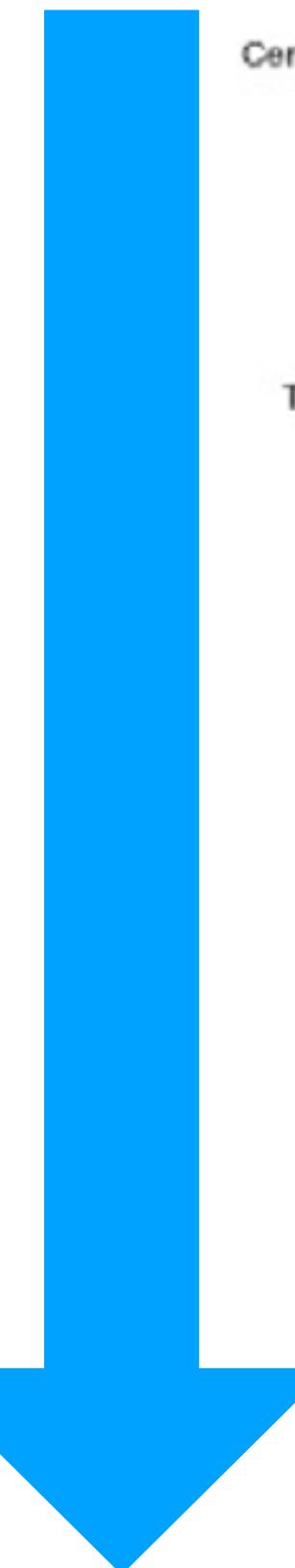


- Not all disease mechanisms are observable.
- “Causal” interventions are often too invasive.
- Networks are highly inter-connected.
- Heterogeneity across individuals.

# Where is the source of complexity?

In a single human cell ( $n=1$ )

zoom-in  
down  
to  
the DNA  
base  
pair  
level



- **3.2B** base pairs (ATGC characters)
- **30k** mRNA/genes (protein coding)  $\pm$  5k to 100k non-coding
- **100+** Epigenetic modifications per gene regulatory element
- **4+** isoforms due to post-transcriptional modifications
- Post-translational modifications

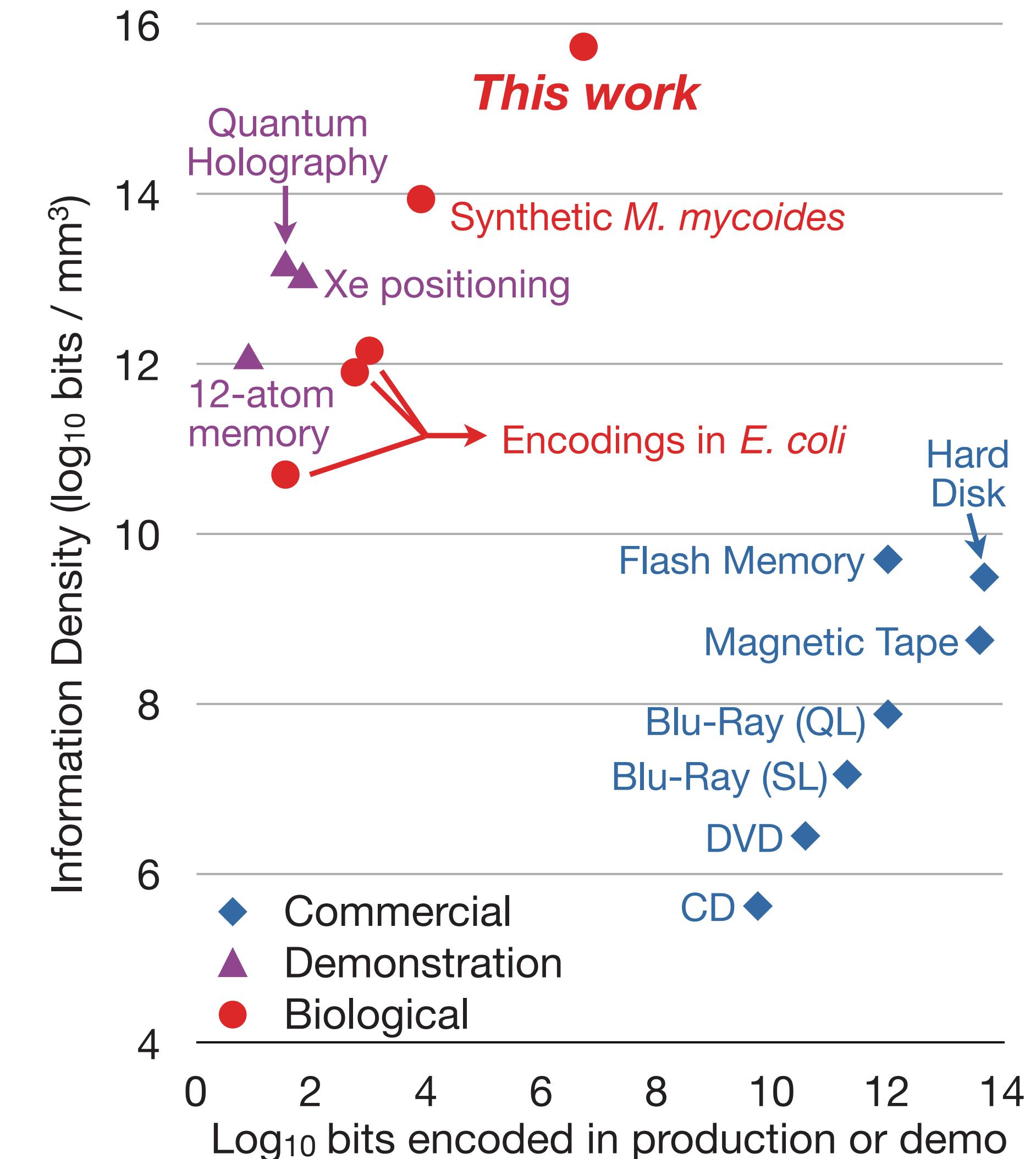
# DNA can store a lot of information



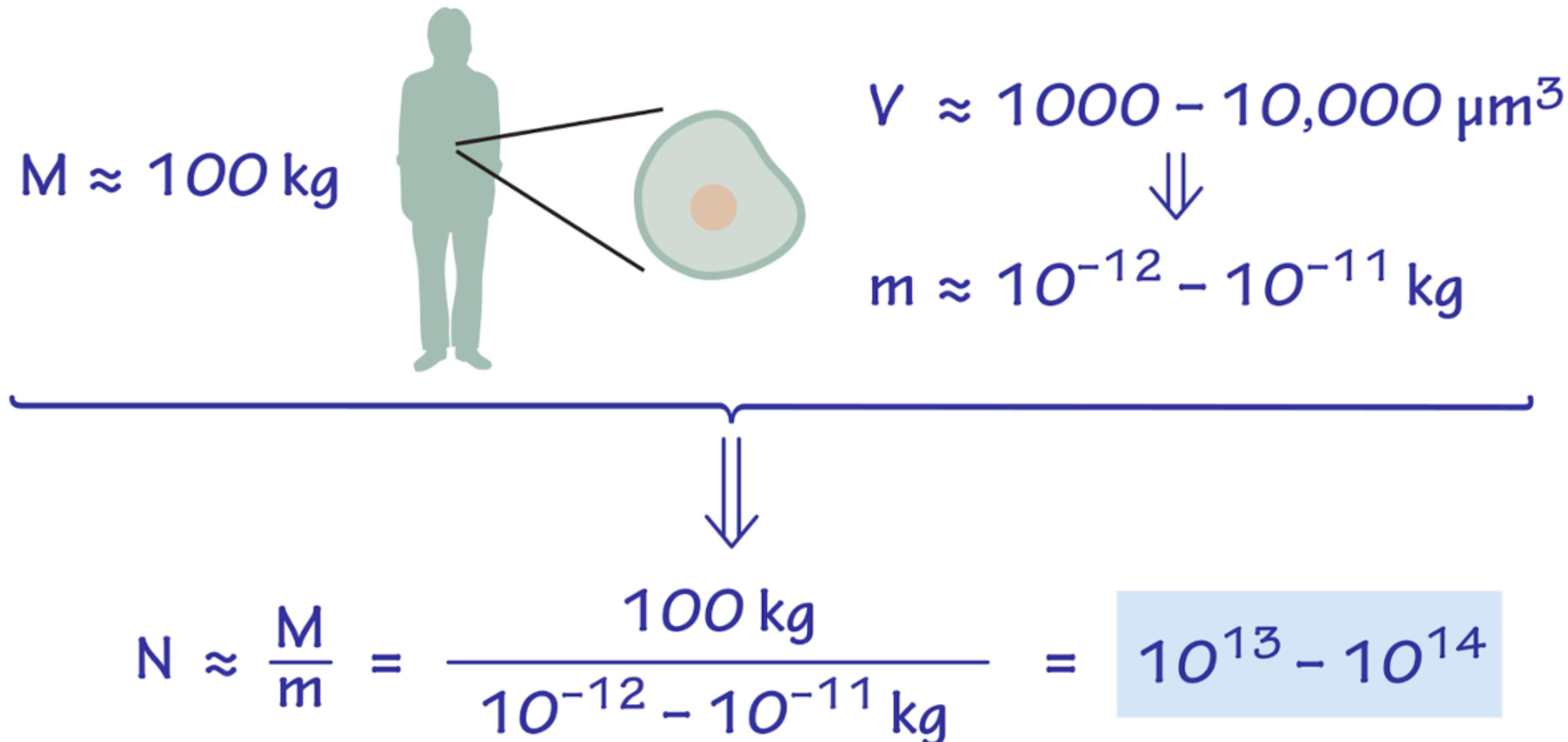
## Next-Generation Digital Information Storage in DNA

George M. Church,<sup>1,2</sup> Yuan Gao,<sup>3</sup> Sriram Kosuri<sup>1,2\*</sup>

- Hard to write but easy to read for a cell
- An efficient way to copy and edit
- Highly parallelized operations
- A long history of evolution



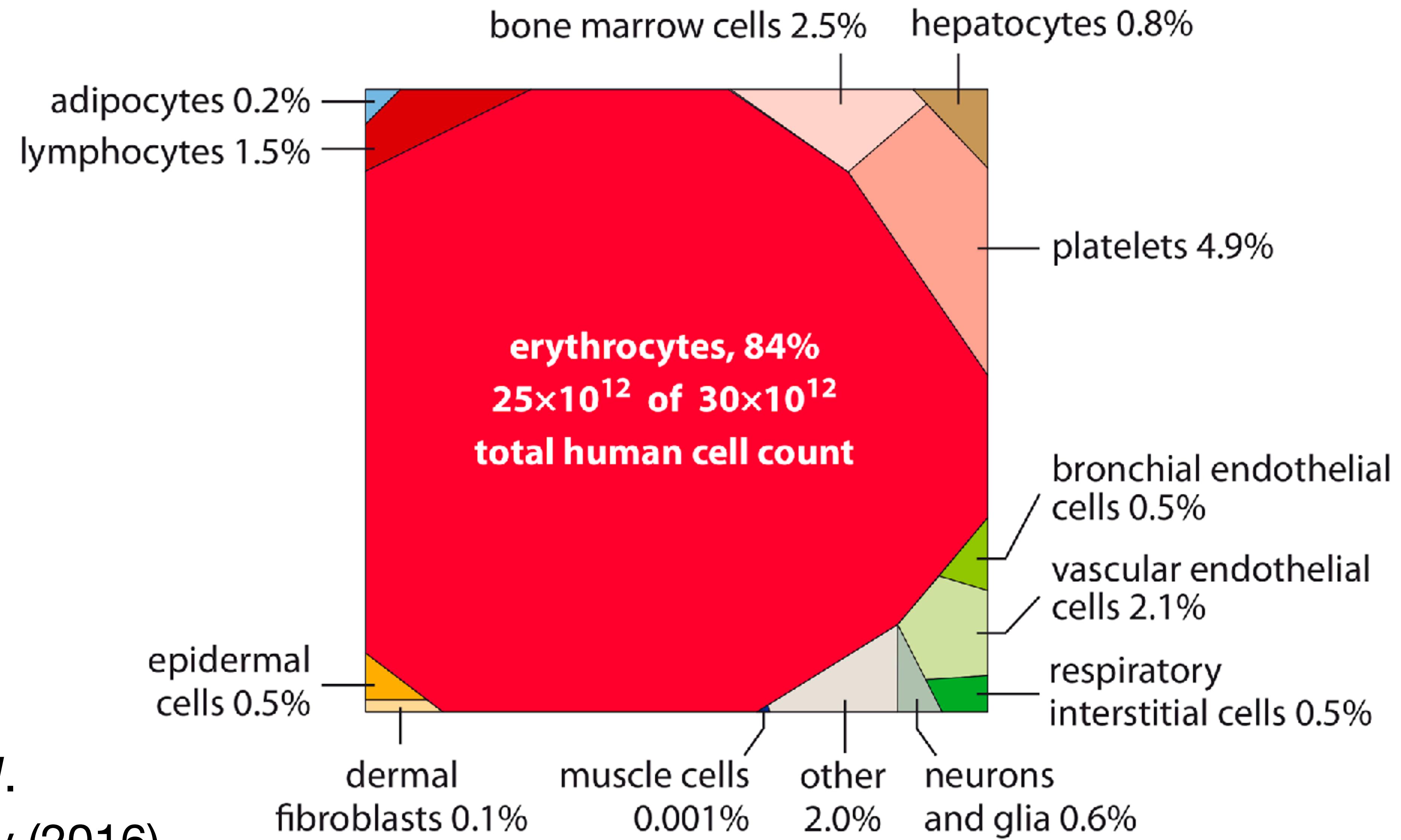
# Trillions of cells



Sender *et al.*

Plos Biology (2016)

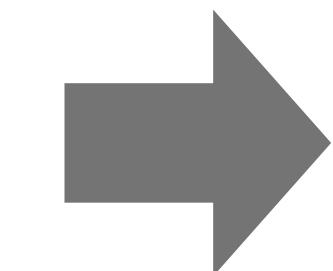
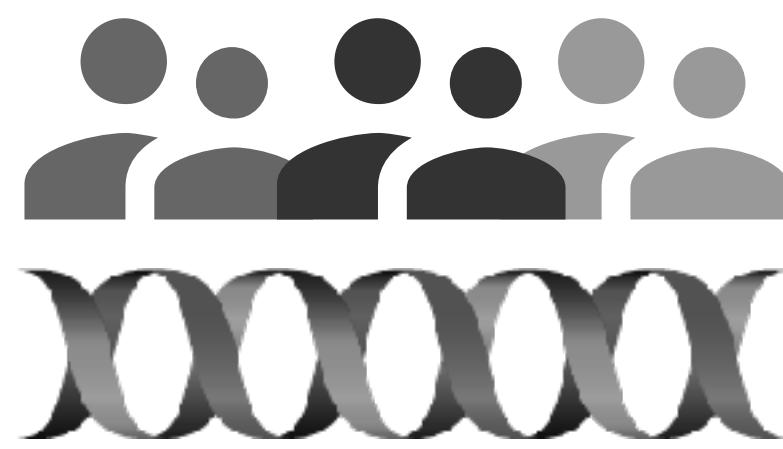
# > 30 trillions of such complexity



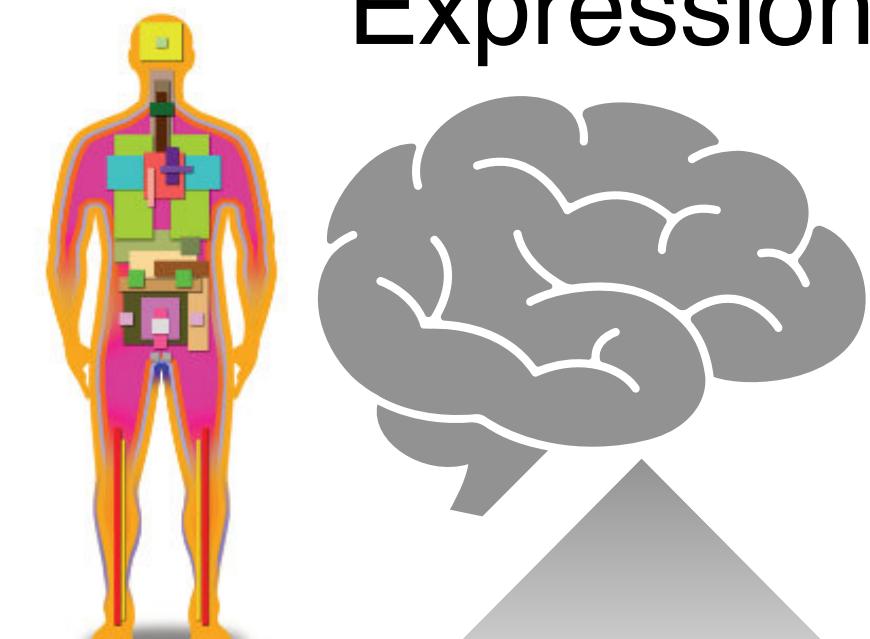
Where is  
**CCTCTGTG**  
**TCGGA**

# We need to understand how it works

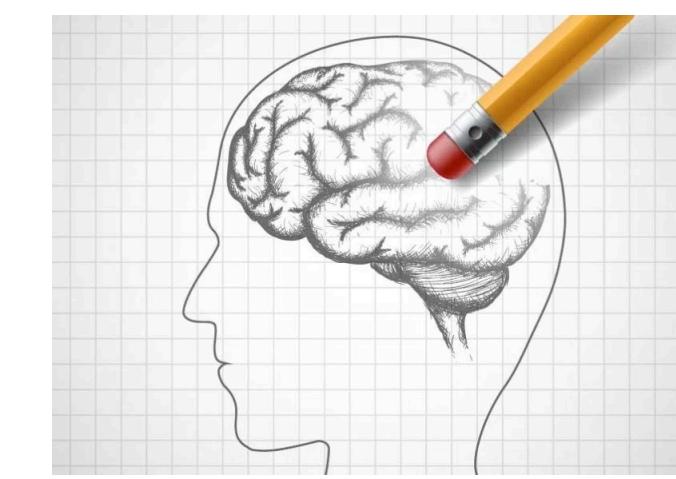
Genotype  
information (X)



Expression (M)



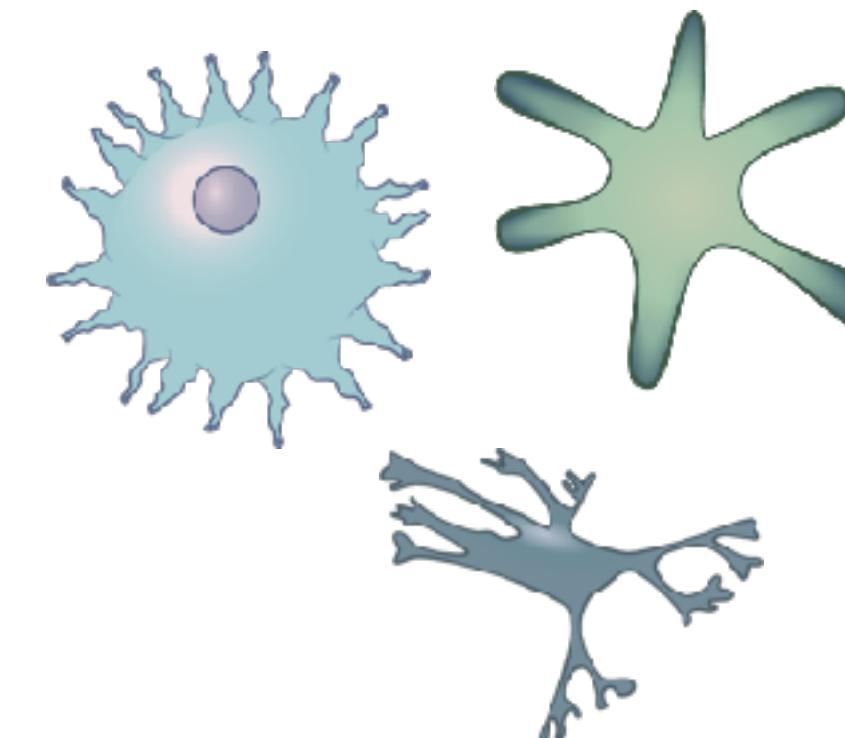
Phenotype (y)



interplay w/  
epigenetic  
mediations

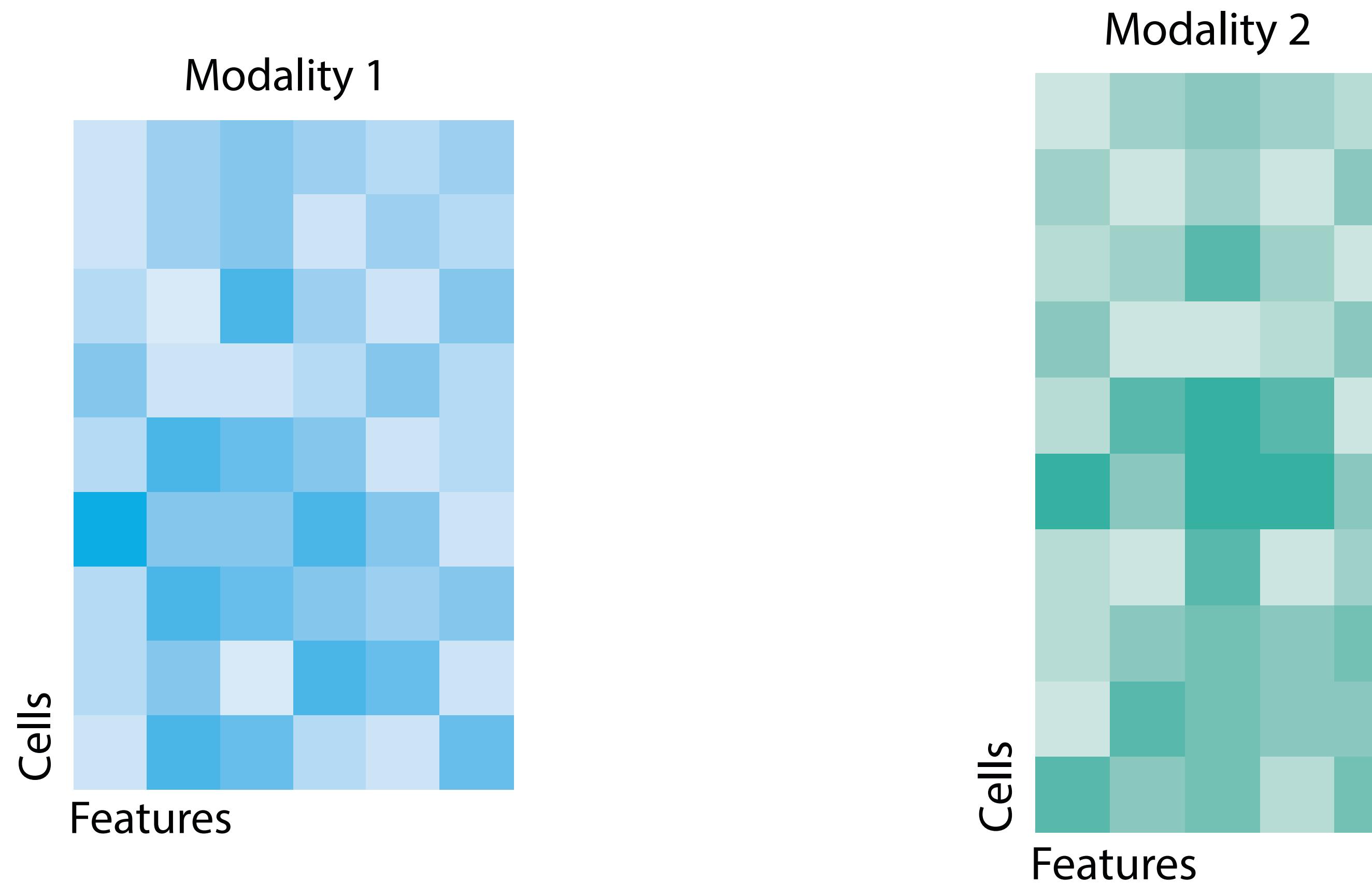


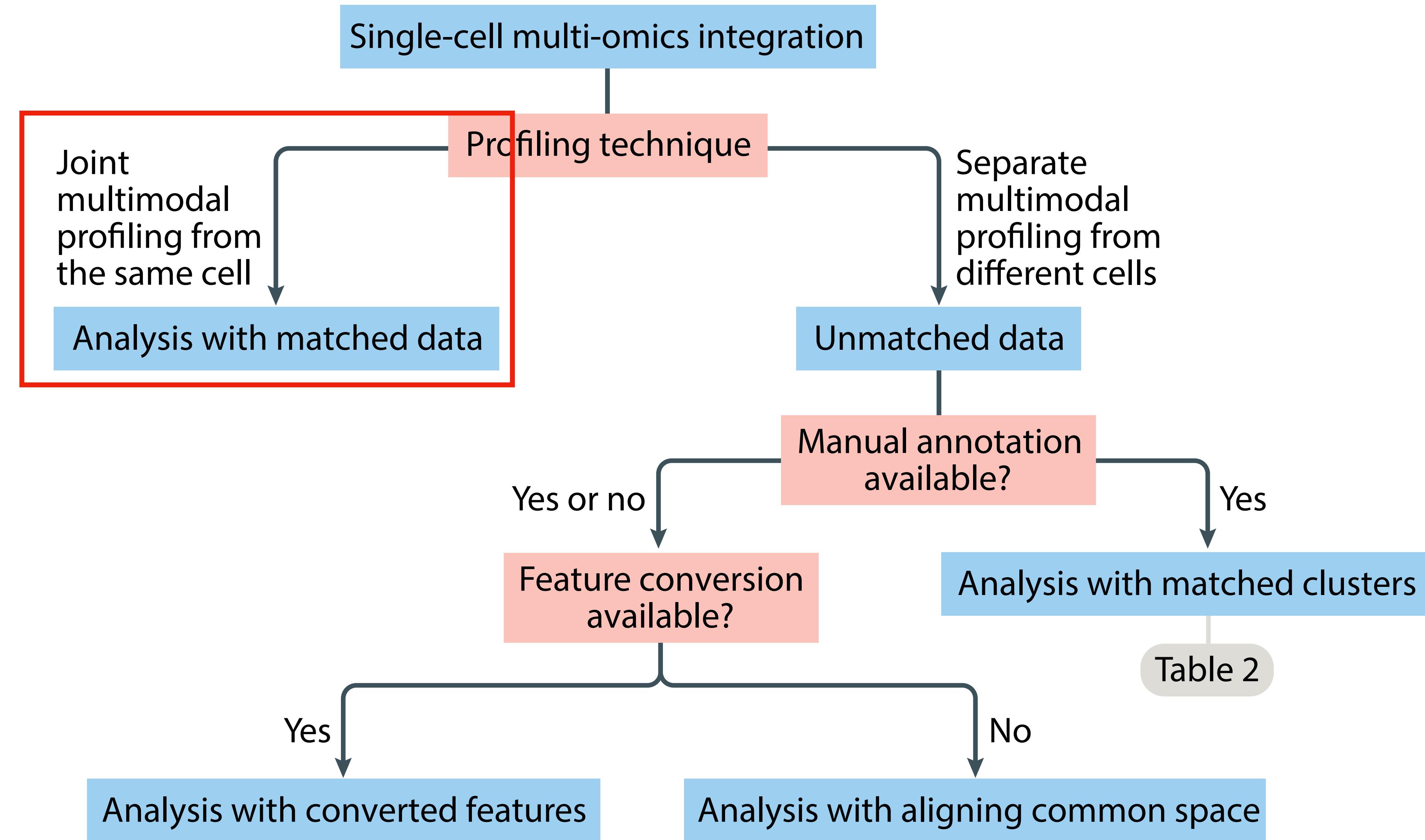
causal cell types



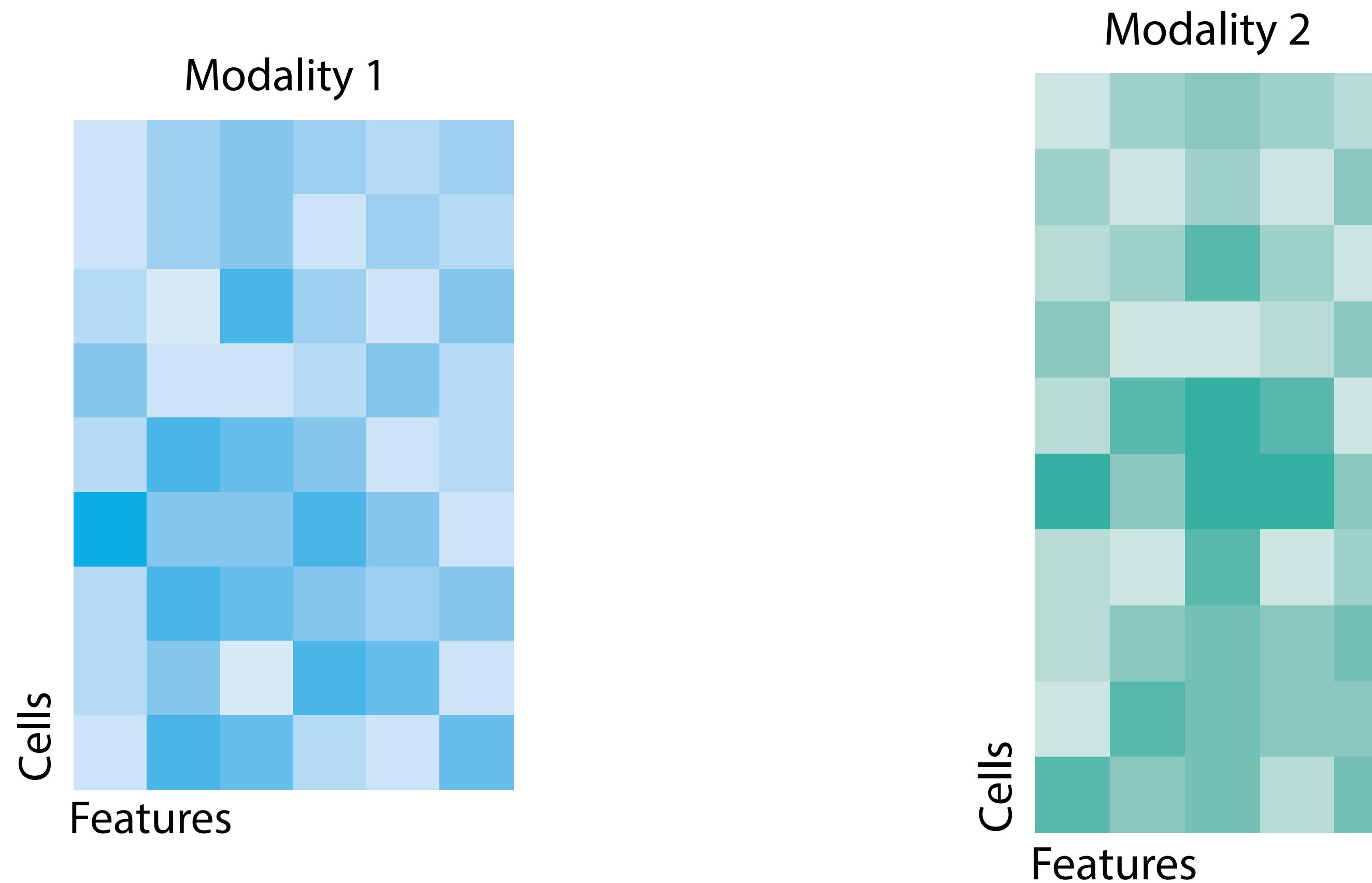
gene-gene interaction,  
gene-environment,  
trans-effect mediated  
by *cis*-effect  
(...)

# Multimodal integration can go in several directions

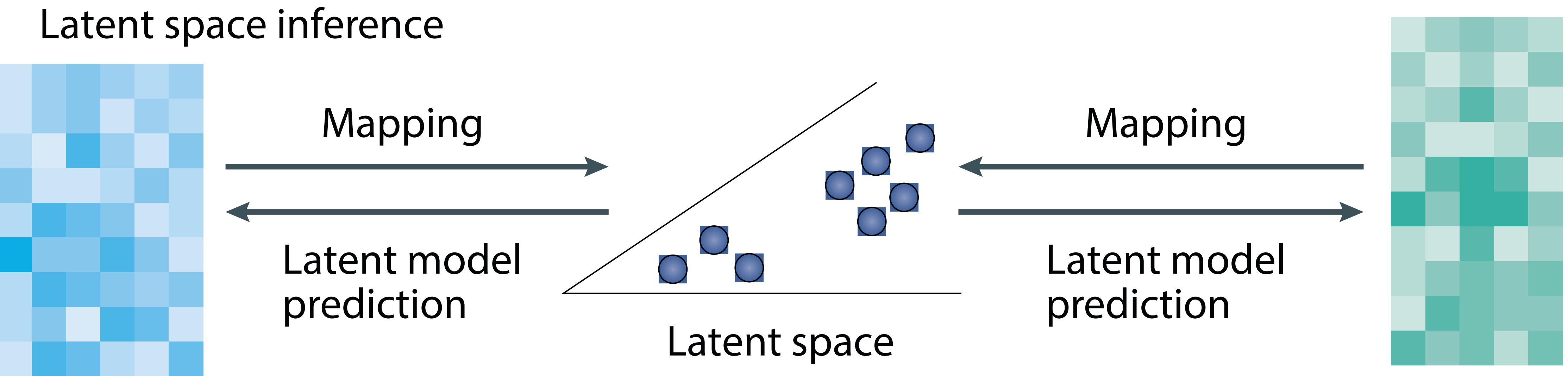




# Multimodal integration can go in several directions

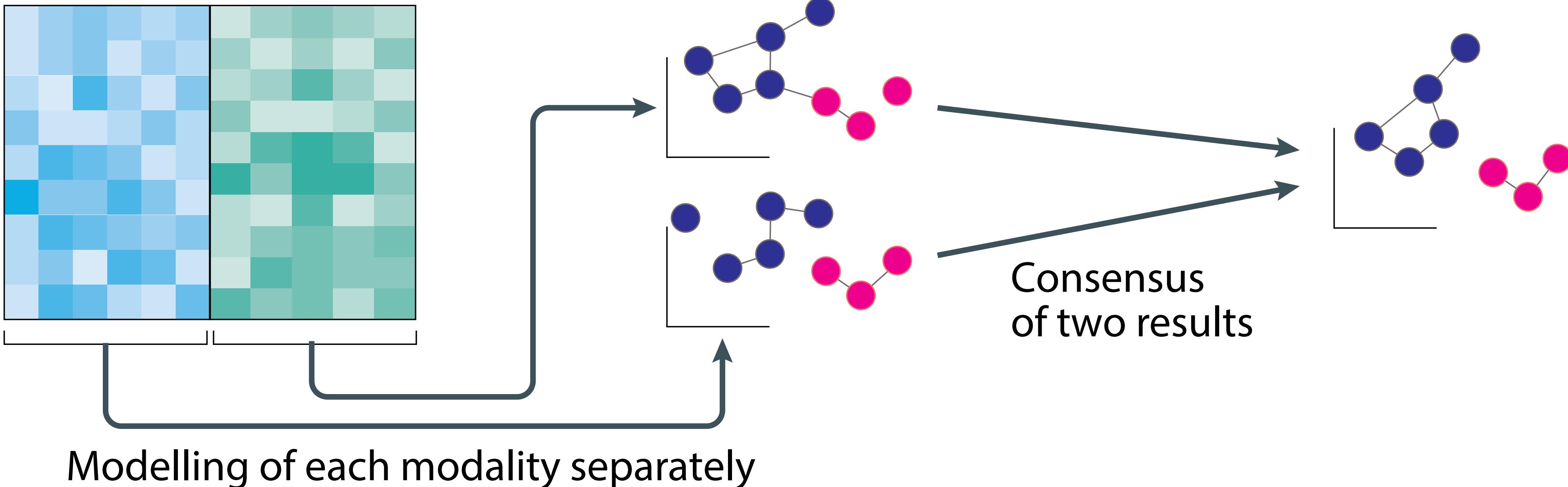


# Different types may implicate different aspects of data generation process



If we have multiple types of data for each cell...  
we can identify consensus groups/modules of cells

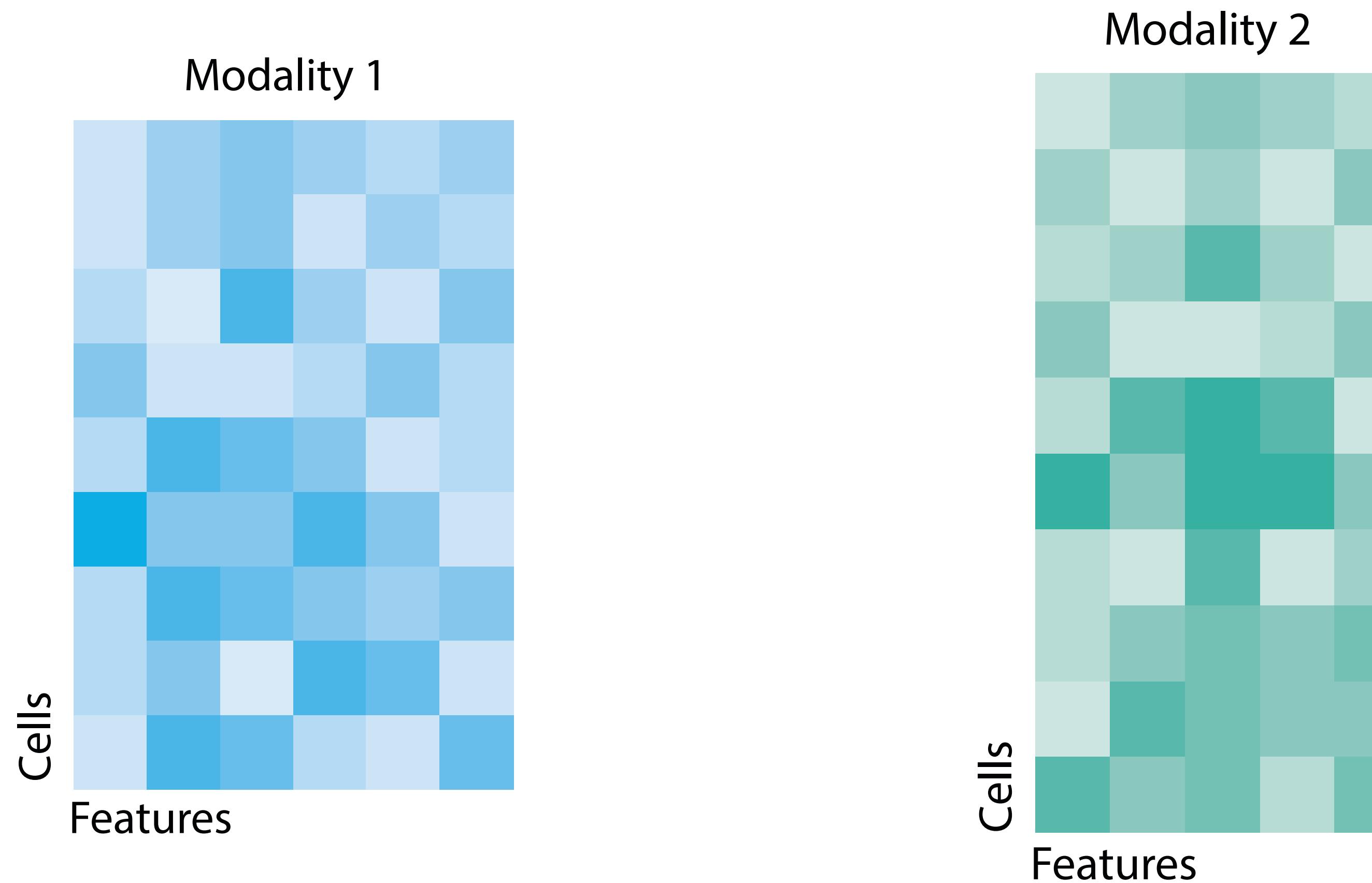
Consensus of individual inferences (late integration)



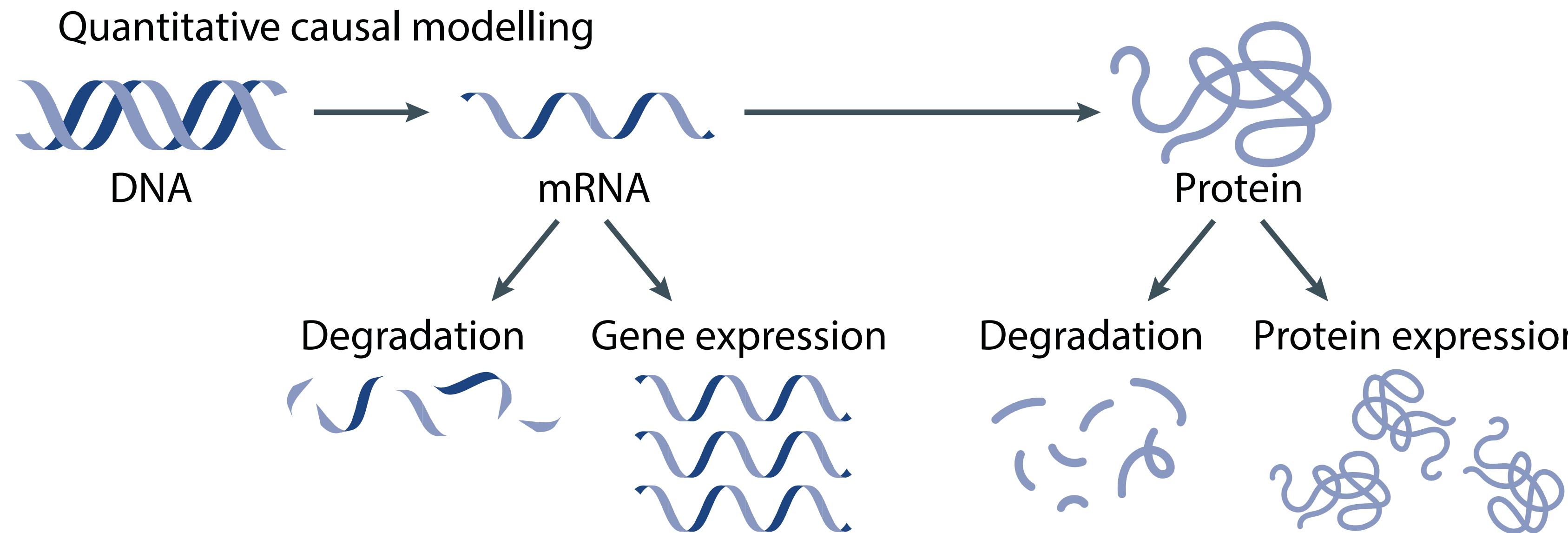
# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms

# Multimodal integration can go in several directions

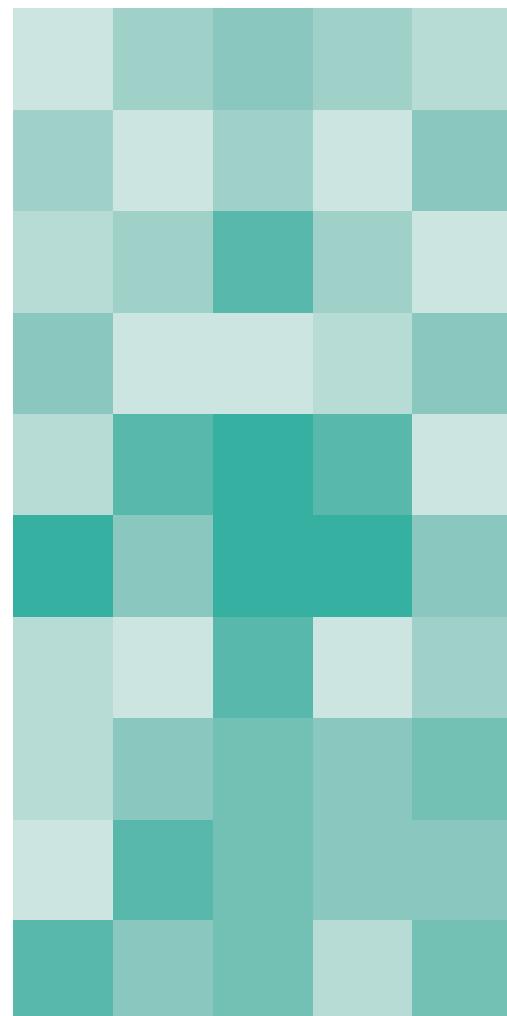


# Ideally, we want to model causal relationships across multiple regulatory mechanisms

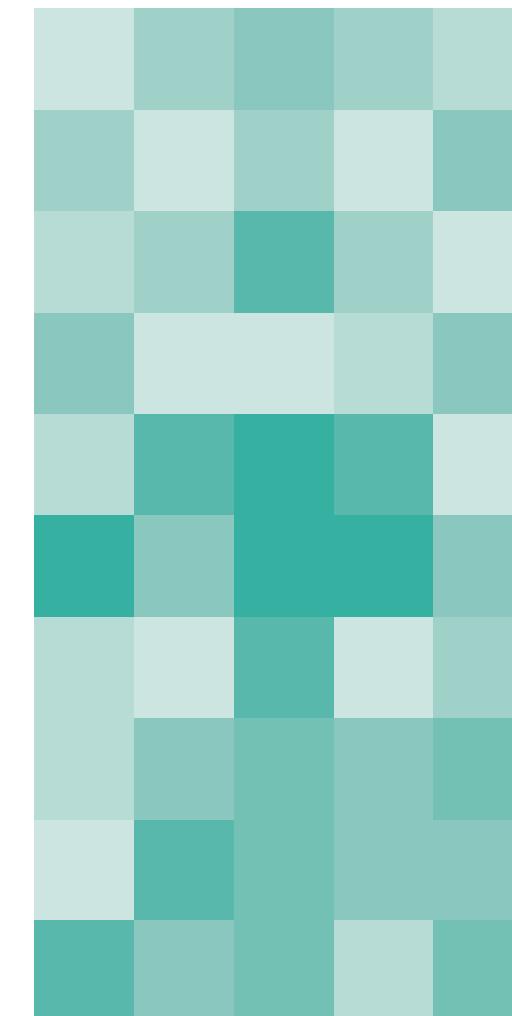


# Can we simply match one feature to the other features?

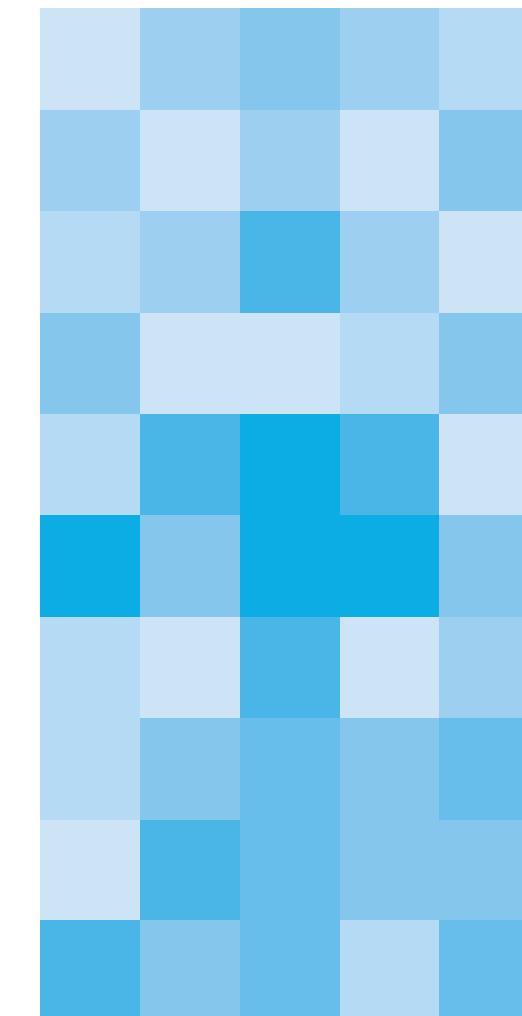
Statistical modelling between features



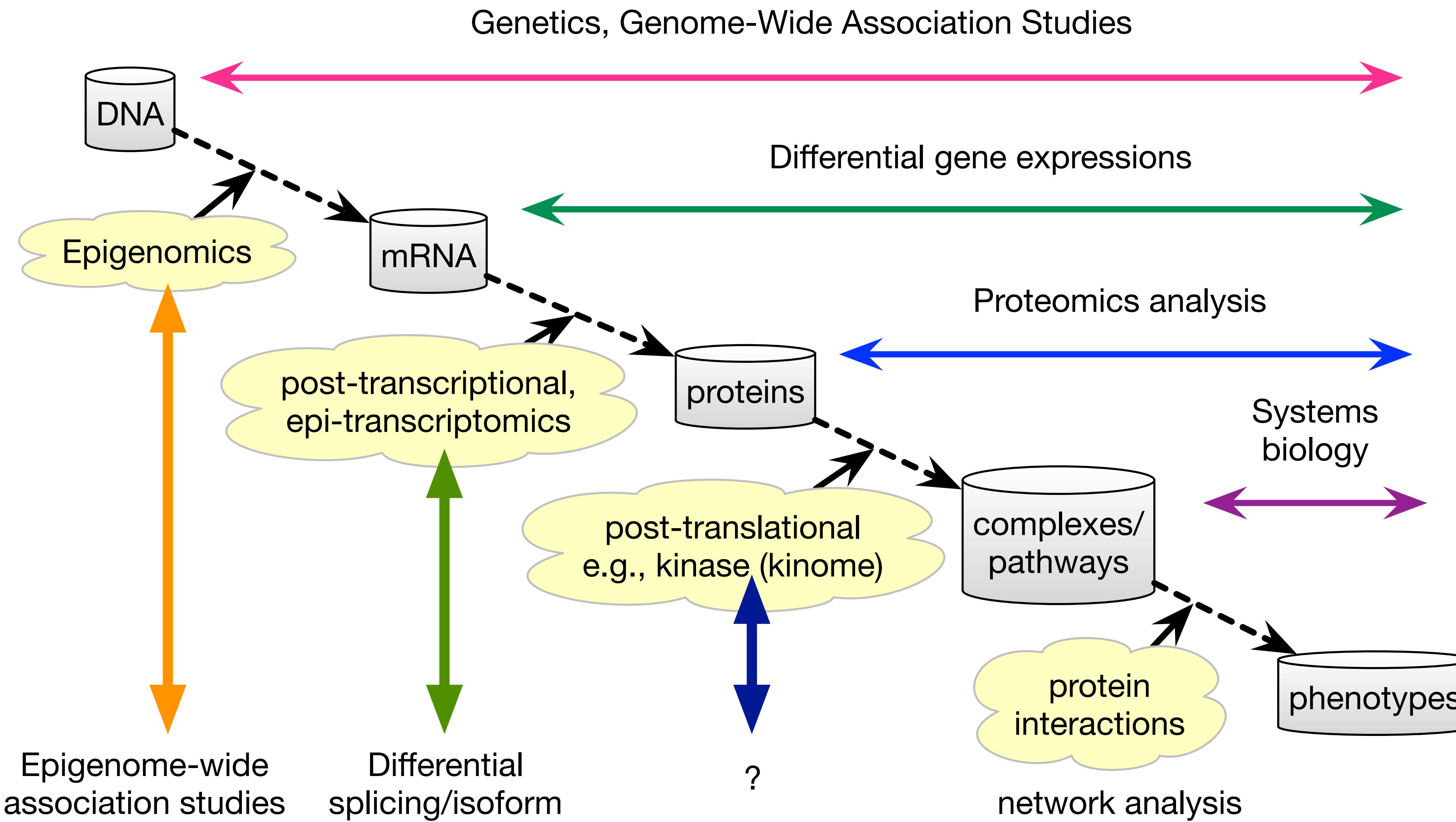
Calibration

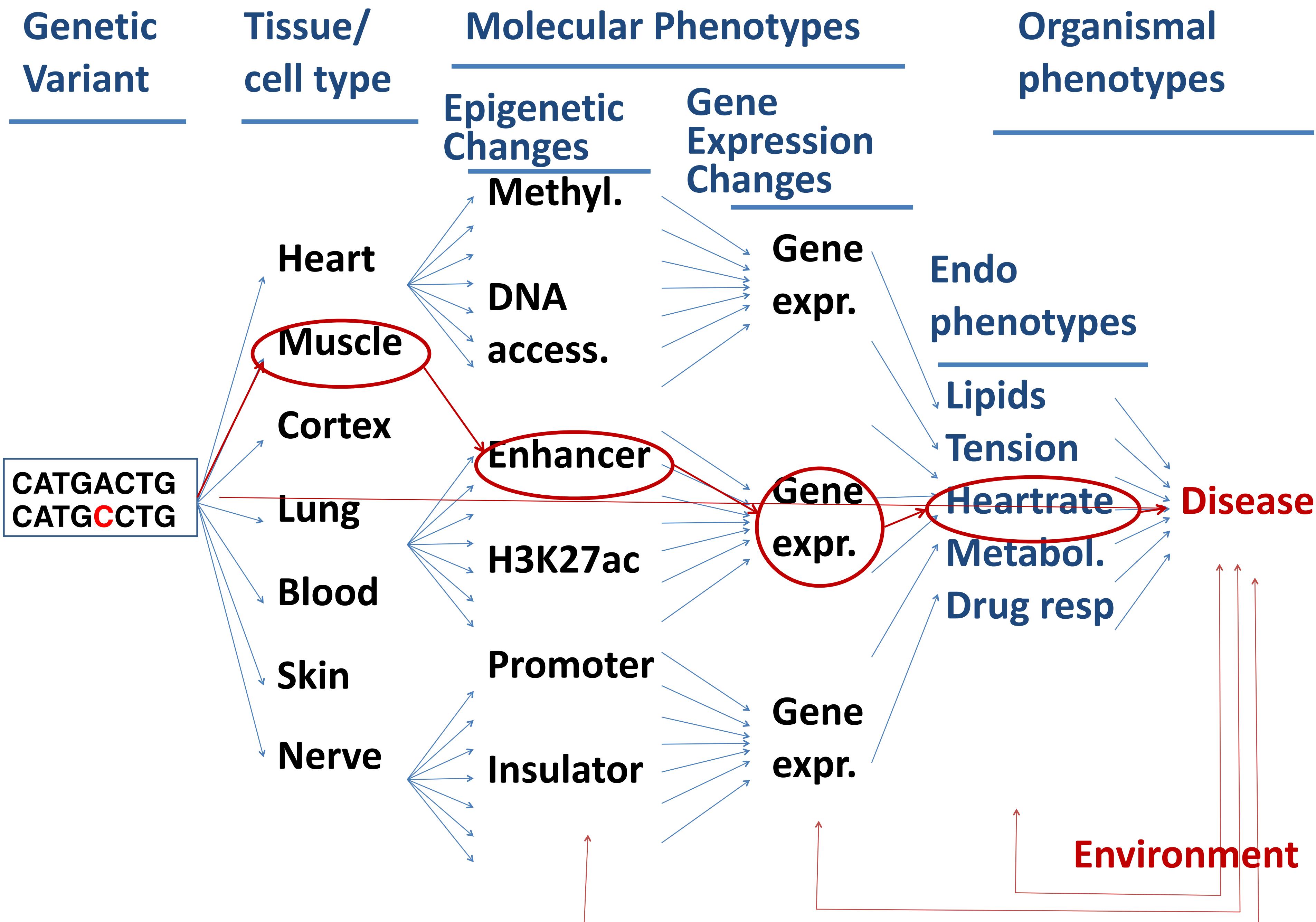


Features can be matched between modalities by modelling a statistical relationship



# Information flow and network across multiple layers of regulations





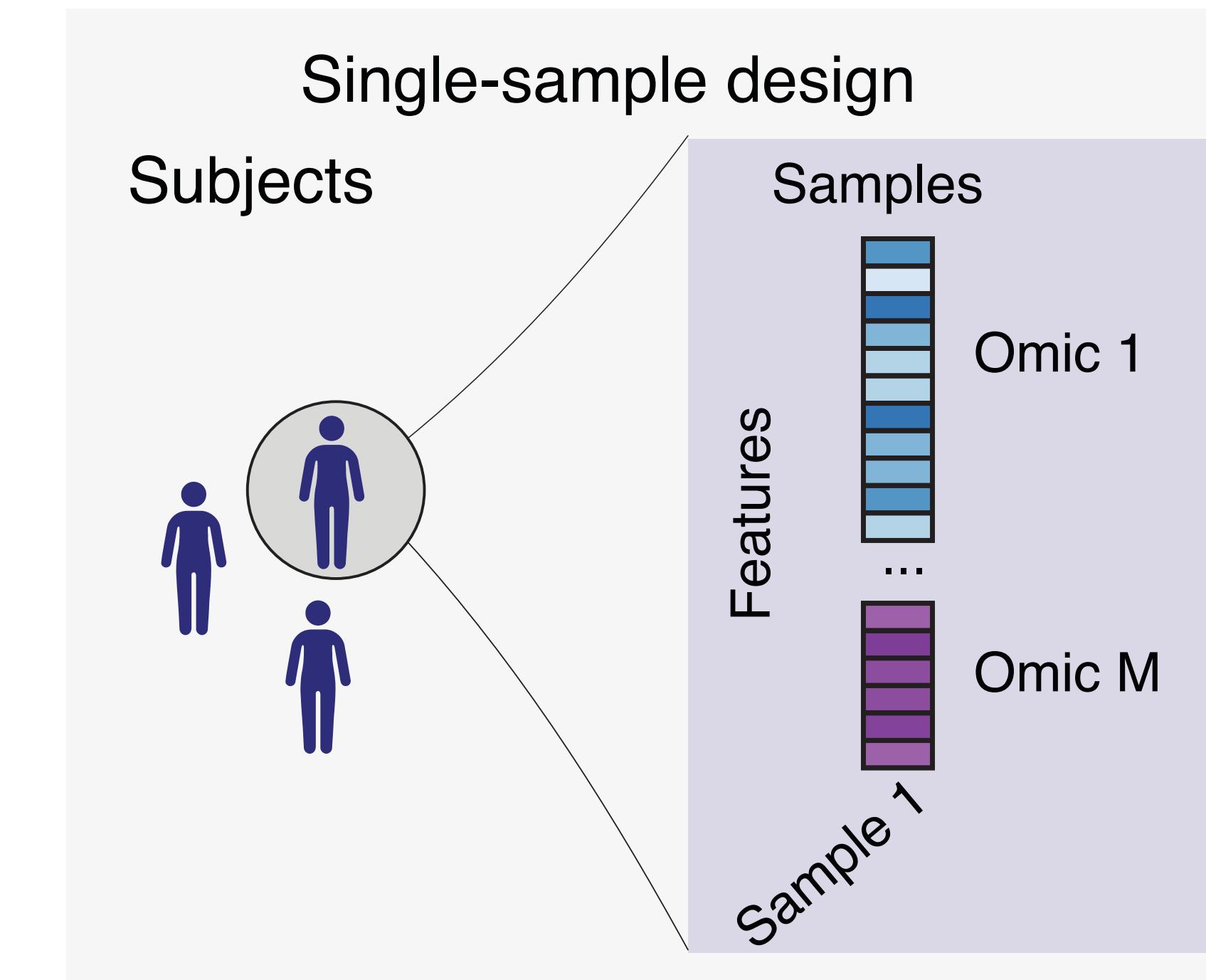
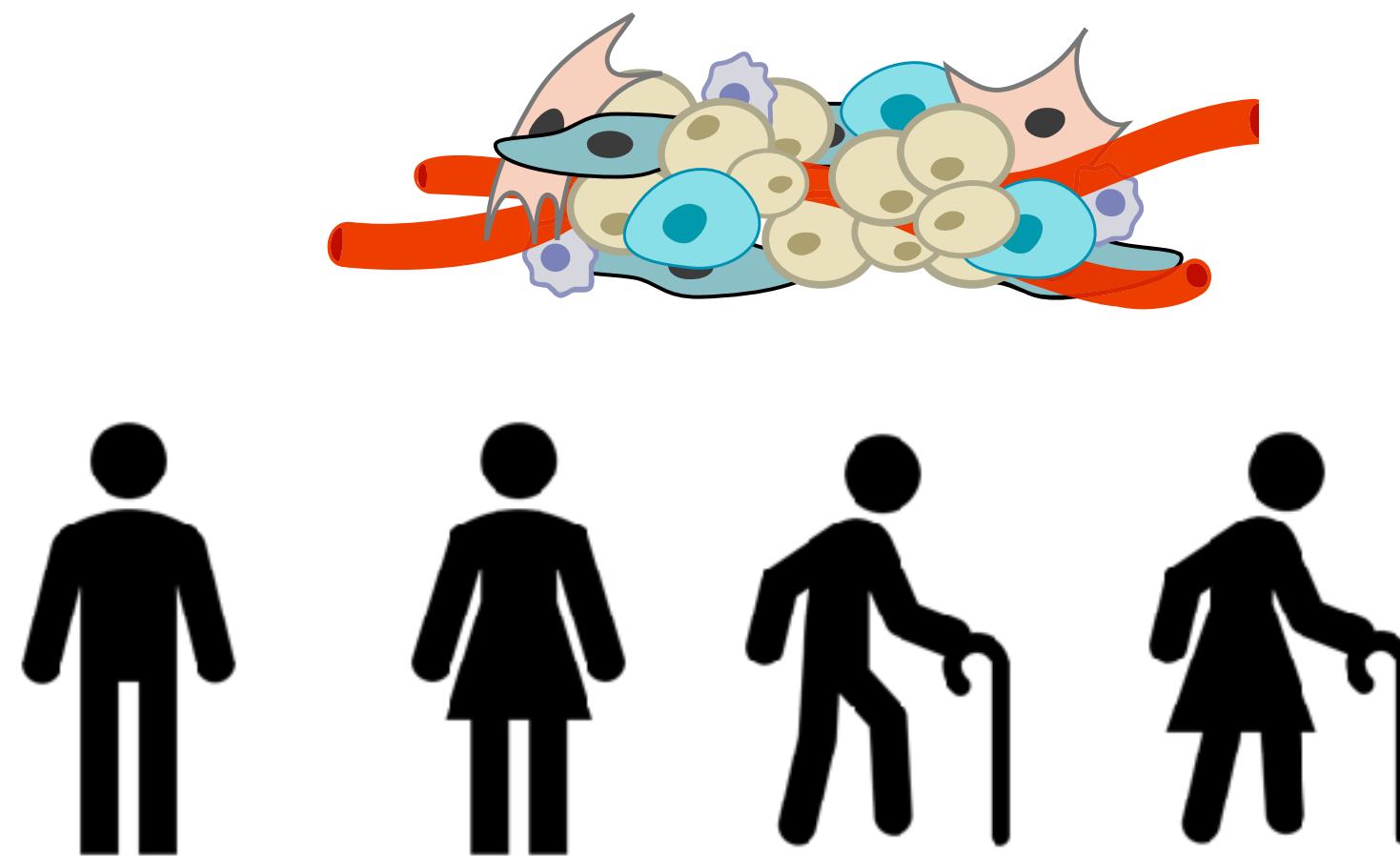
Slide credit: Manolis Kellis

**Feedback from environment / disease state**

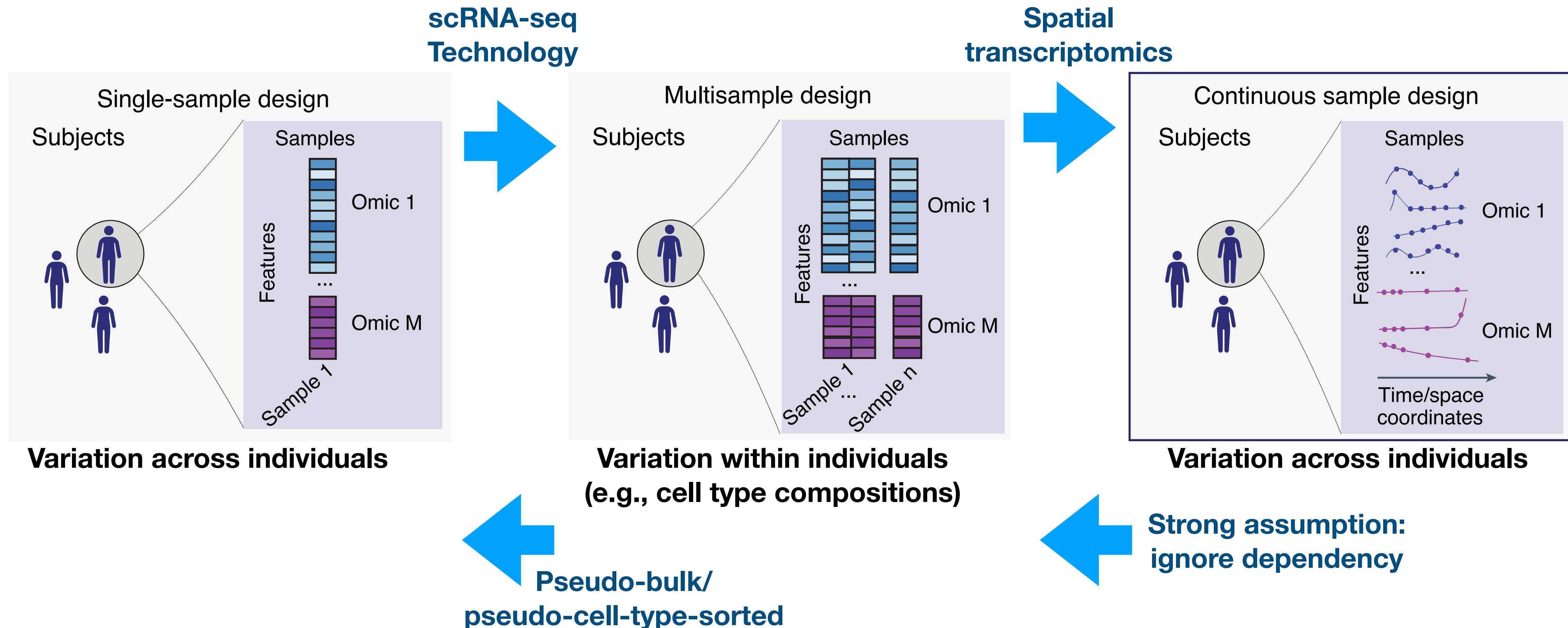
# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms

# A traditional design for Omics studies: one sample per individual



# A new experimental design: many samples per individual



*Method*

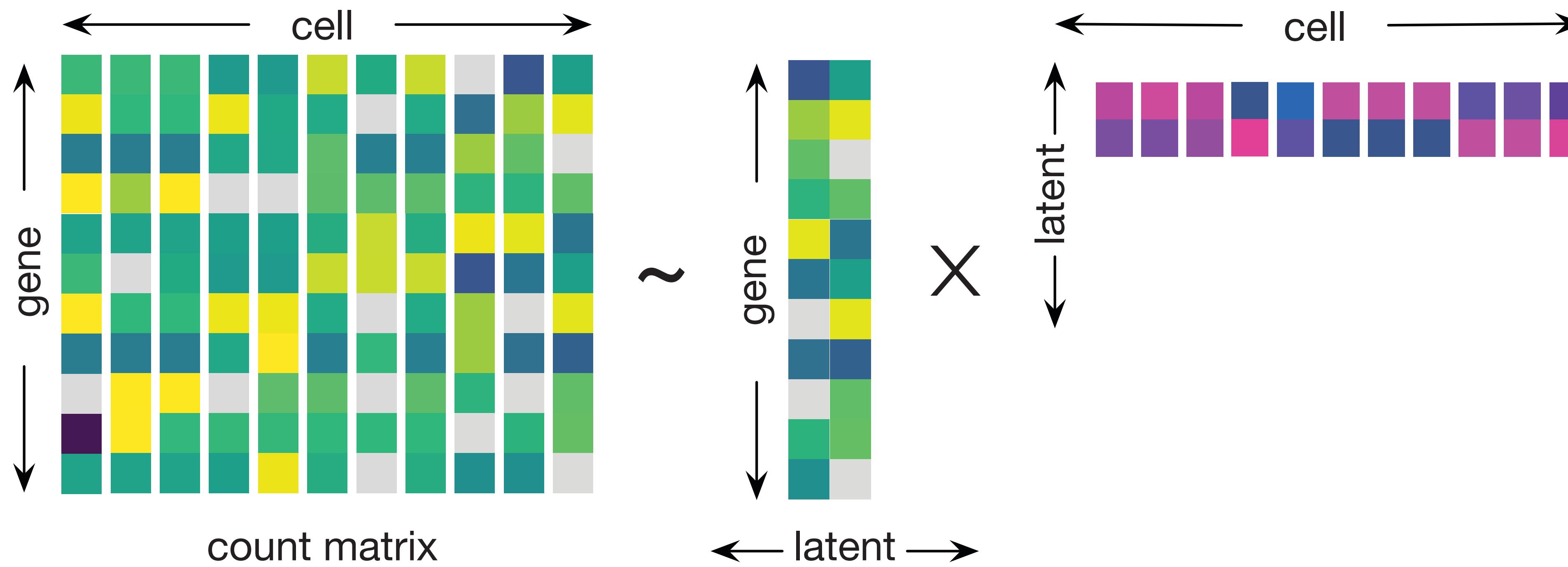


molecular  
systems  
biology

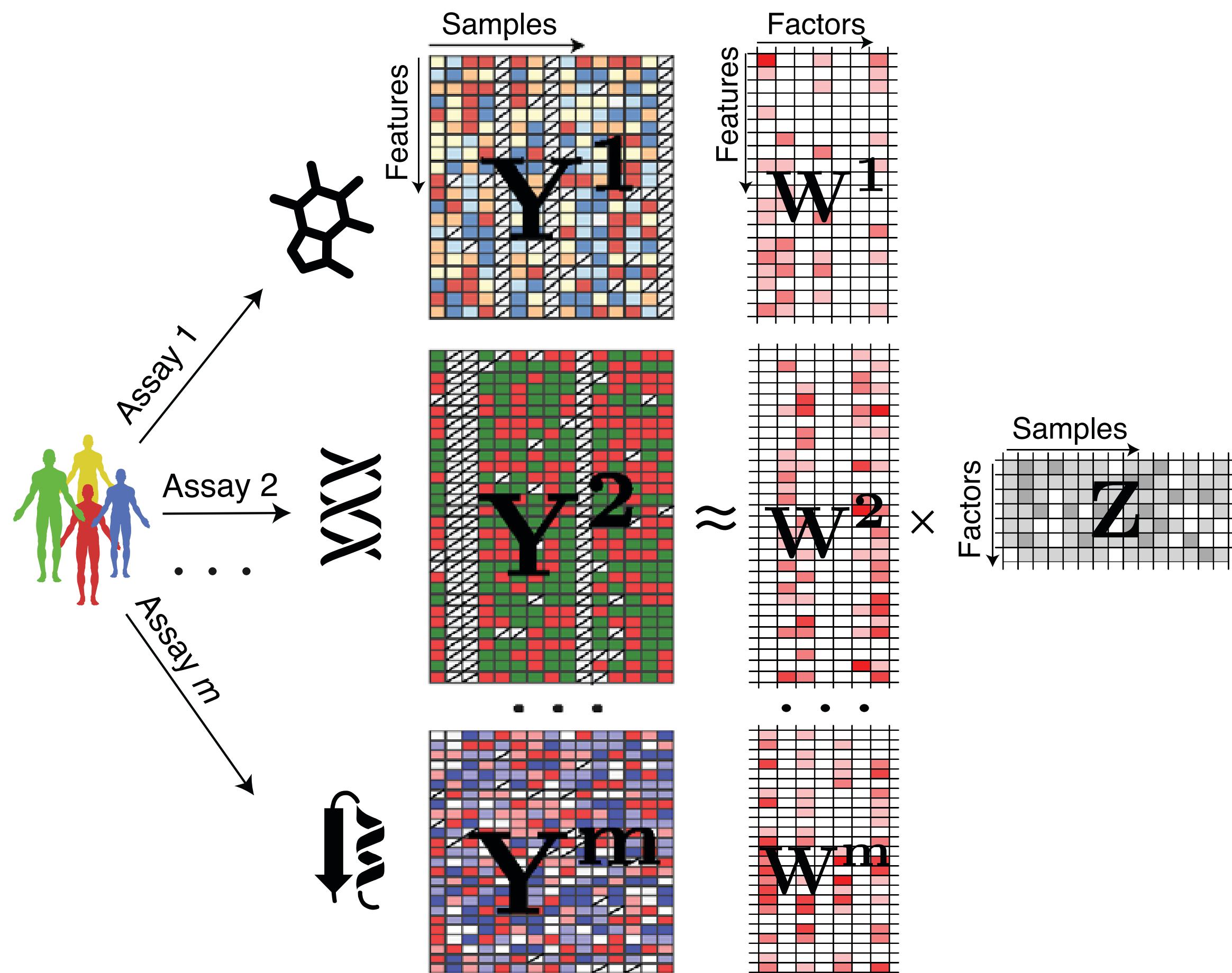
# Multi-Omics Factor Analysis—a framework for unsupervised integration of multi-omics data sets

Ricard Argelaguet<sup>1,†</sup> , Britta Velten<sup>2,†</sup> , Damien Arnol<sup>1</sup> , Sascha Dietrich<sup>3</sup> , Thorsten Zenz<sup>3,4,5</sup> , John C Marioni<sup>1,6,7</sup> , Florian Buettner<sup>1,8,\*</sup> , Wolfgang Huber<sup>2,\*\*</sup>  & Oliver Stegle<sup>1,2,\*\*\*</sup> 

# Factorization of a single data matrix



# MOFA: Multi-Omics Factor Analysis



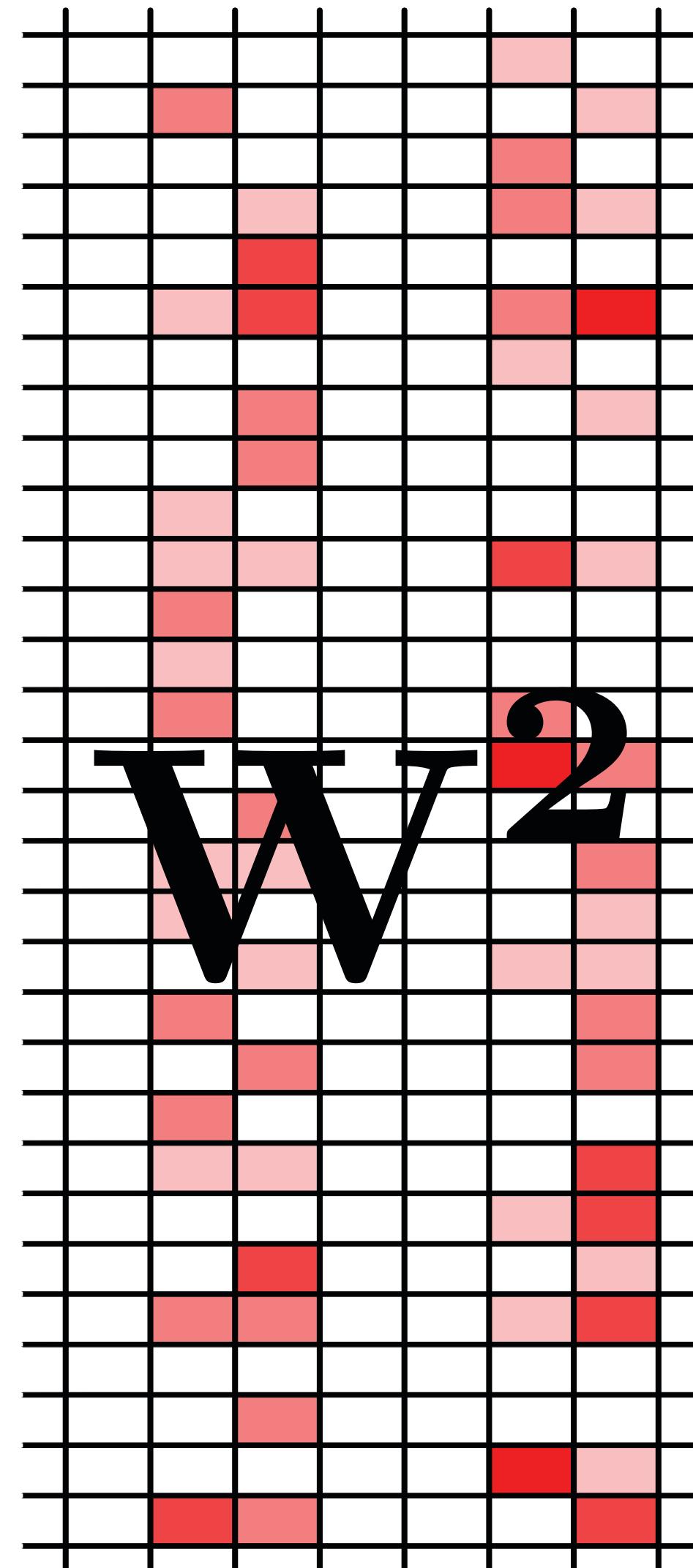
For each data modality  $m$

$$\mathbf{Y}^{(m)} = \mathbf{Z}\mathbf{W}^{(m)^\top} + \boldsymbol{\epsilon}^{(m)}$$

$$W = S \hat{W}$$

column-wise  
spike (0/1 for each rank)

# MOFA: Multi-Omics Factor Analysis



For each data modality  $m$

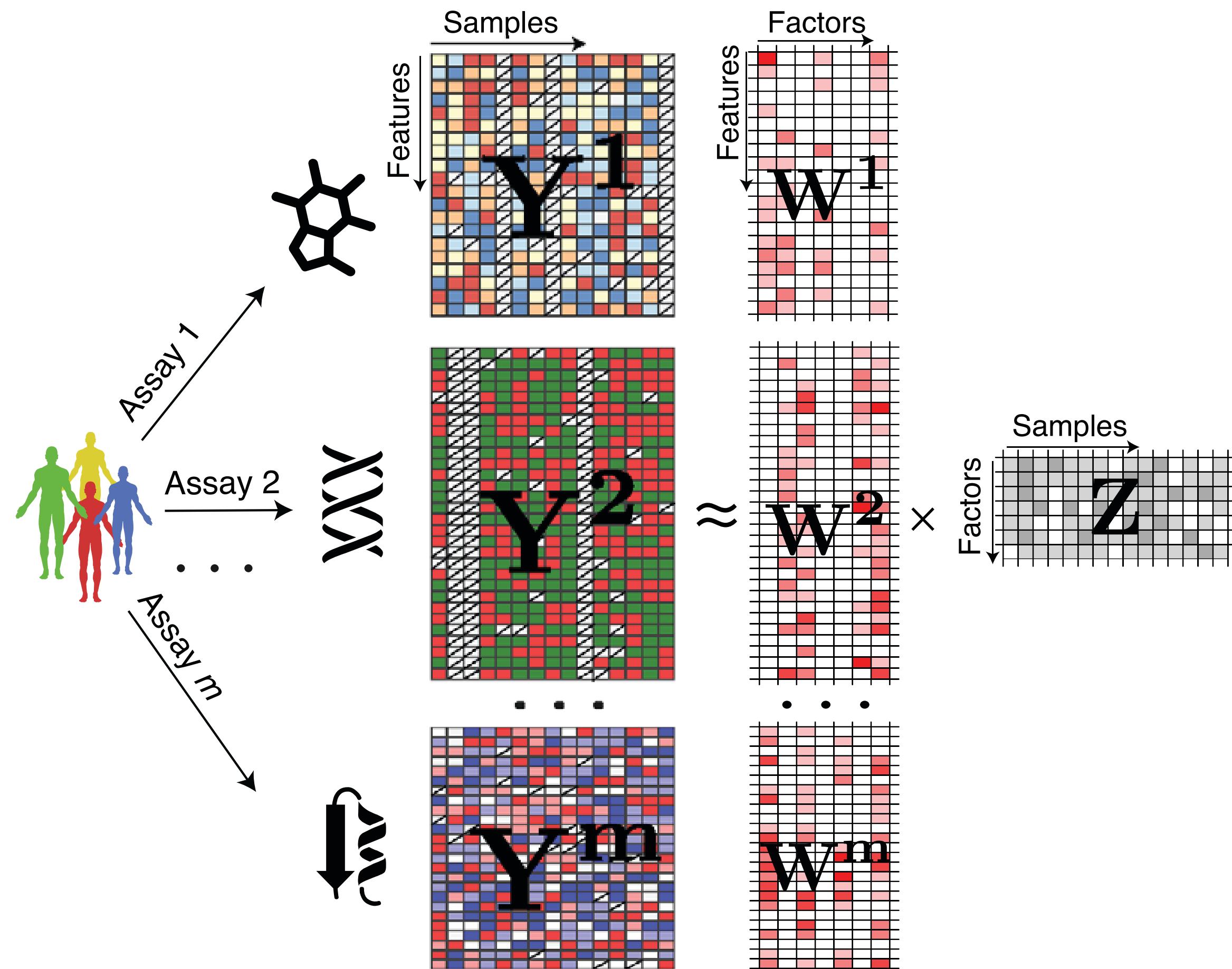
$$\mathbf{Y}^{(m)} = \mathbf{Z}\mathbf{W}^{(m)\top} + \boldsymbol{\epsilon}^{(m)}$$

$$\mathbf{W} = \mathbf{S} \hat{\mathbf{W}}$$

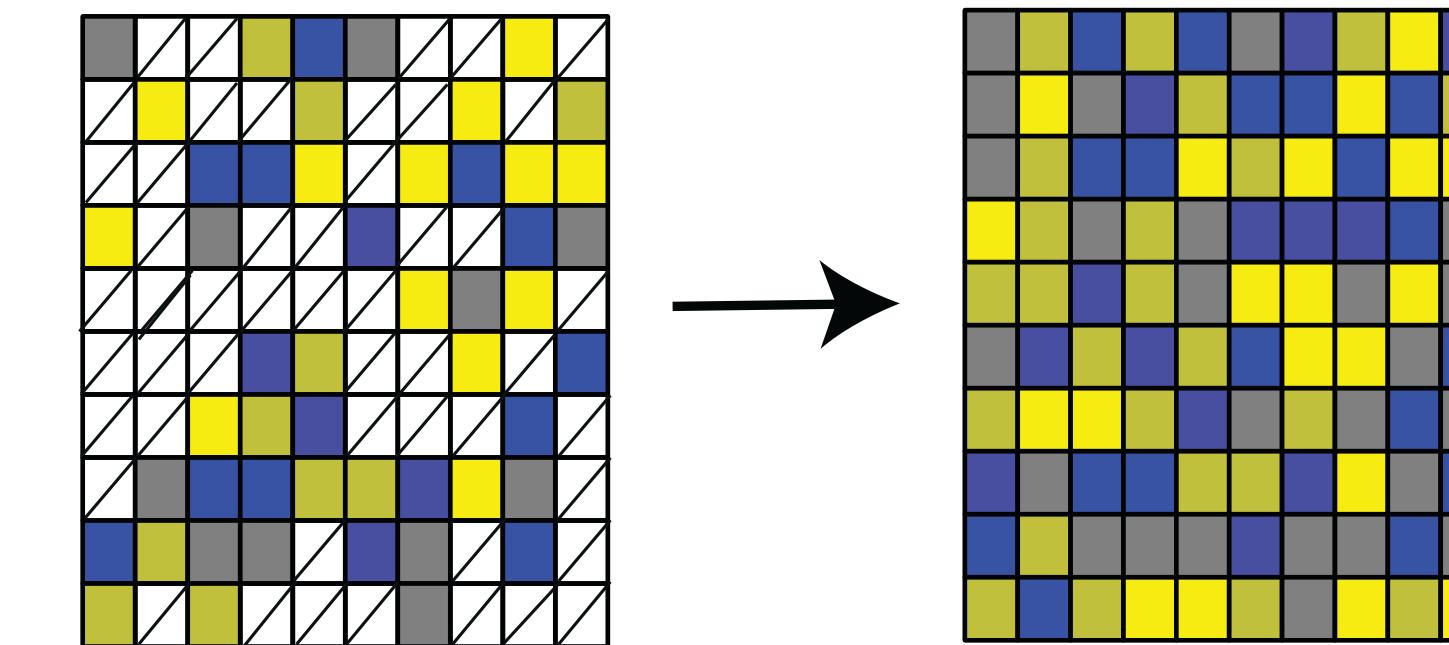


column-wise  
spike (0/1 for each rank)

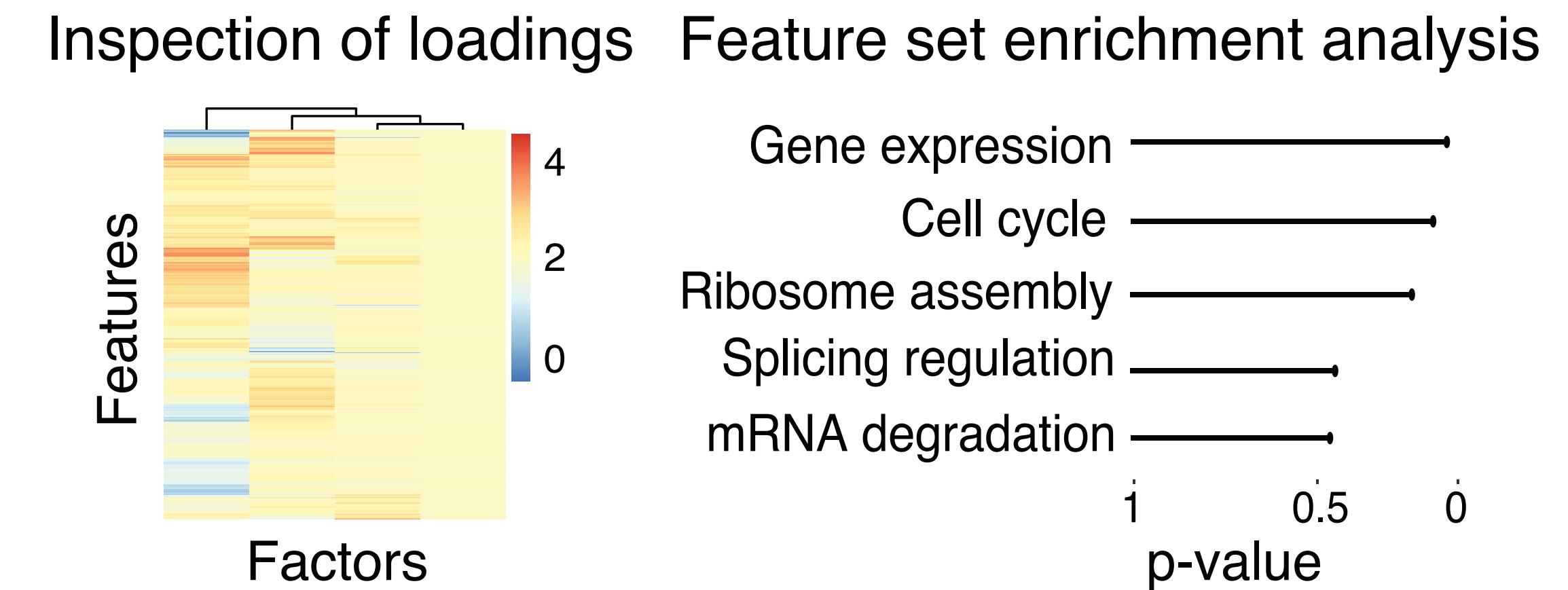
# When is MOFA useful?



## Imputation of missing values

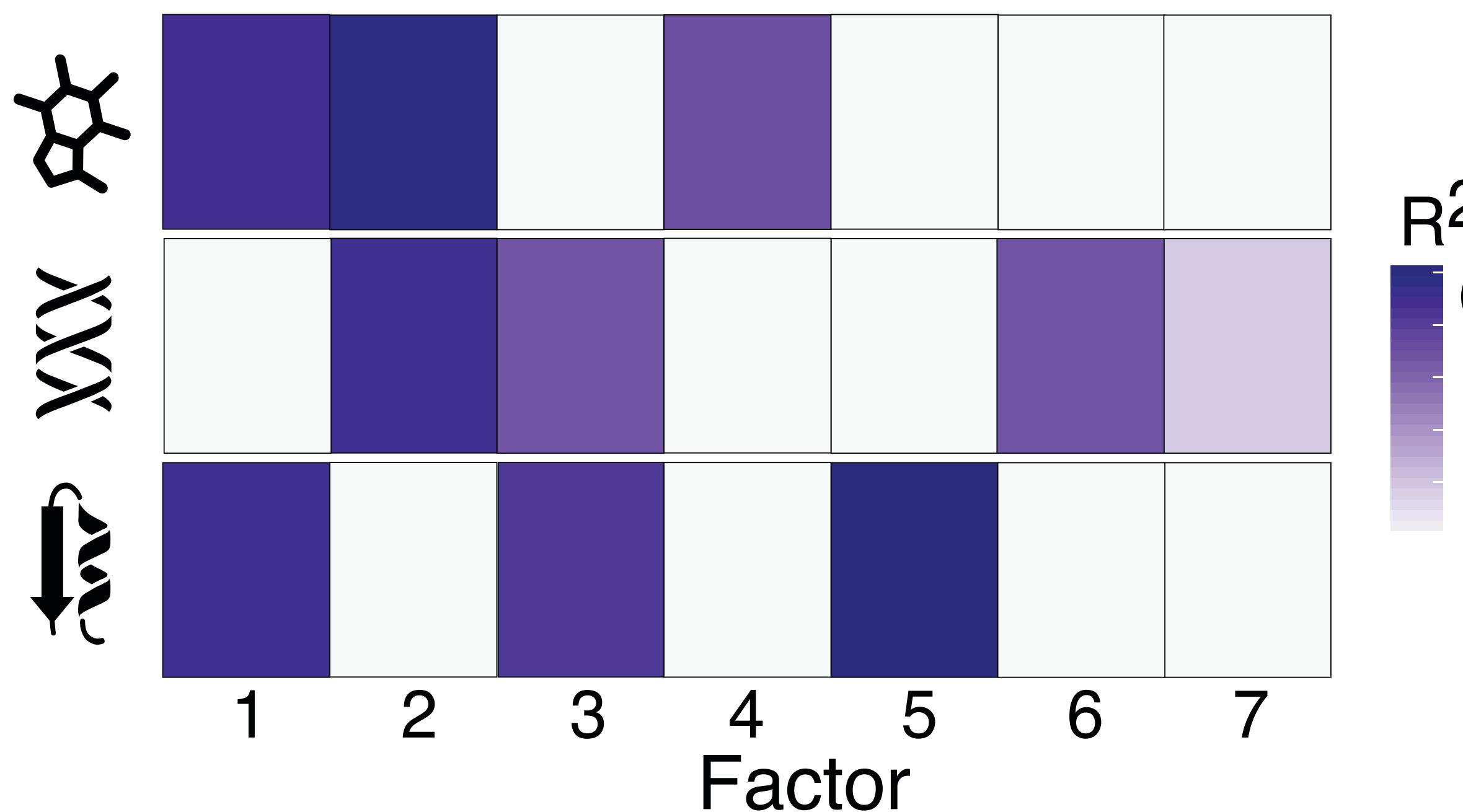


## Annotation of factors

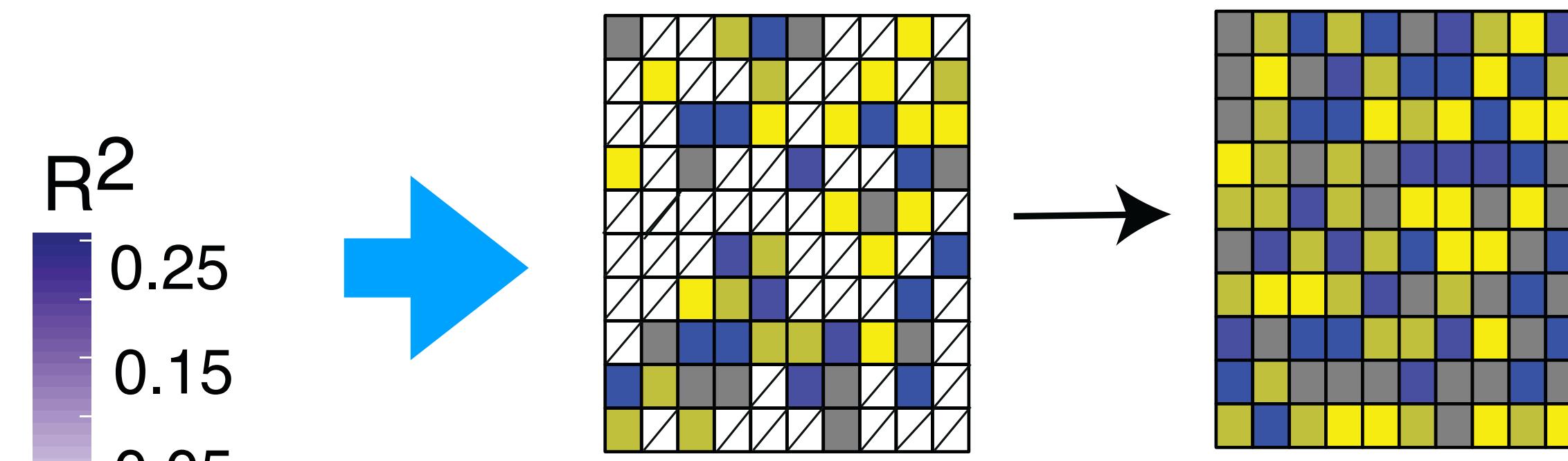


# When is MOFA useful?

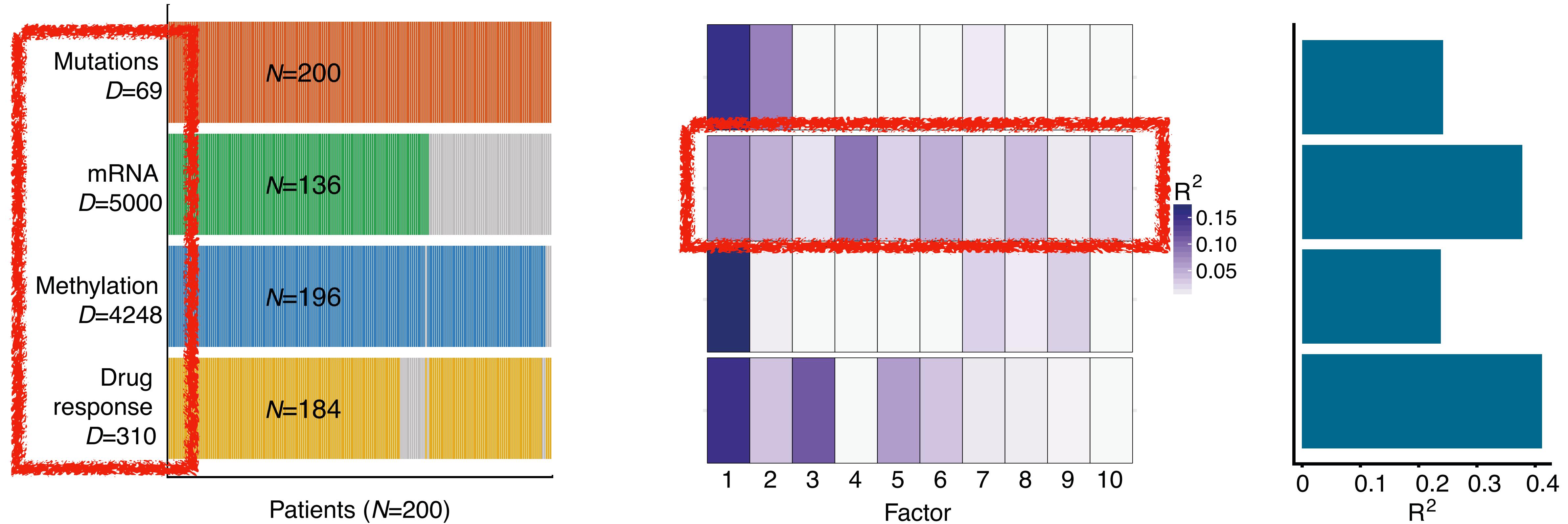
Variance decomposition by factor



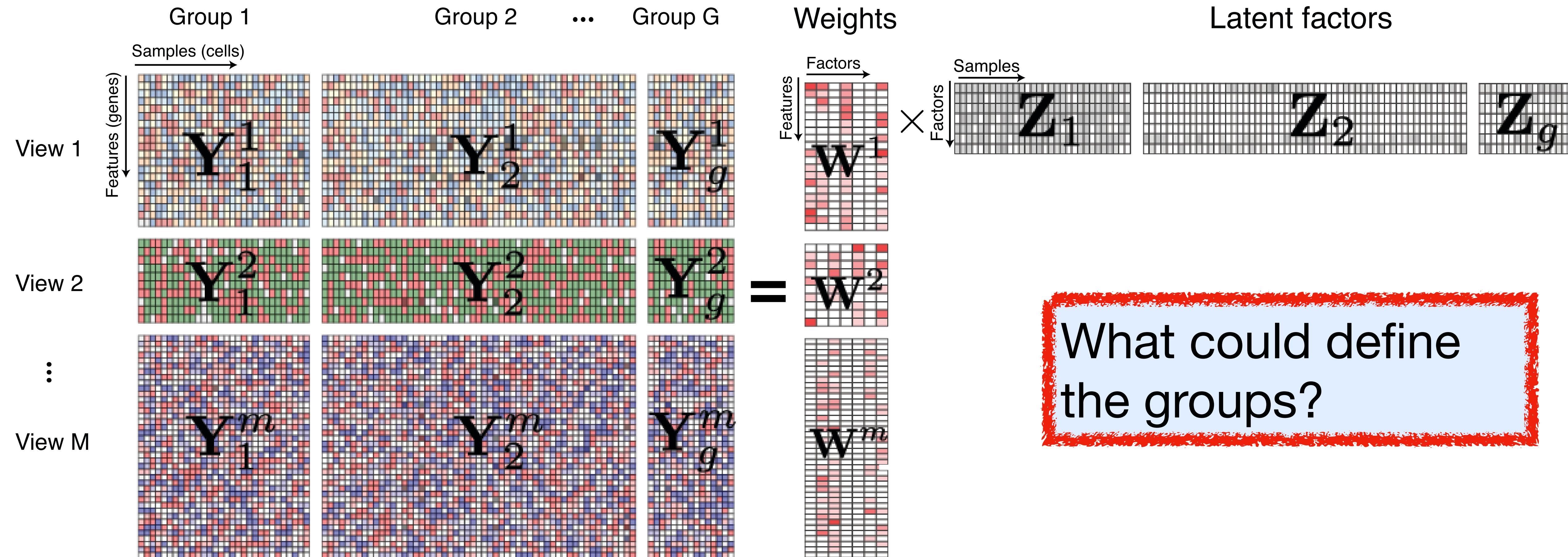
Imputation of missing values



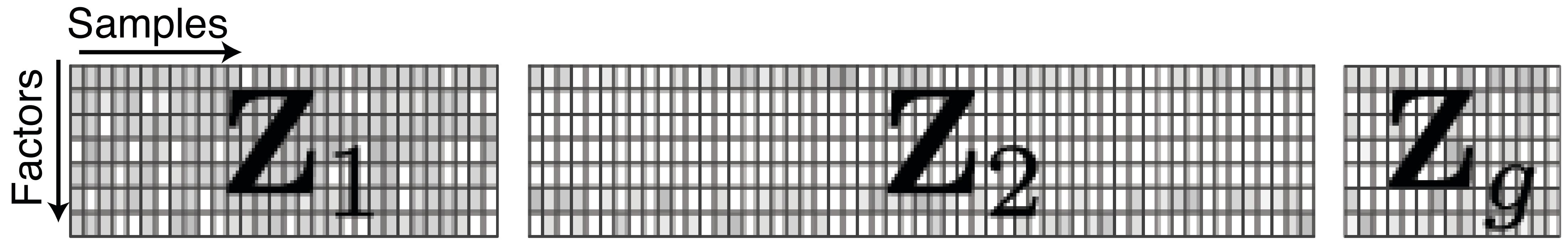
# When is MOFA *not* really useful?



# MOFA+ is multi-group and modality modelling

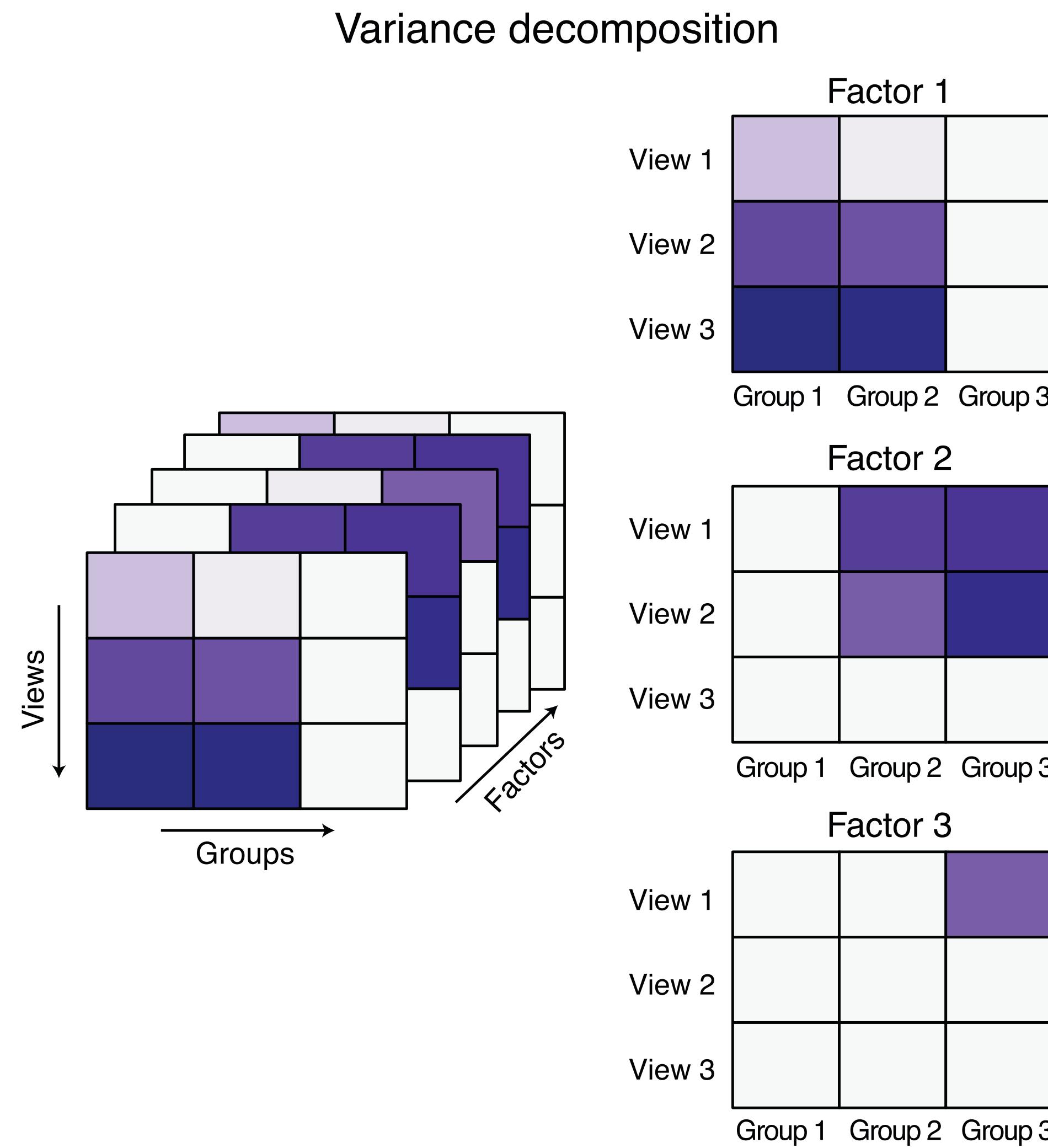


# MOFA+ can model heterogeneity across groups

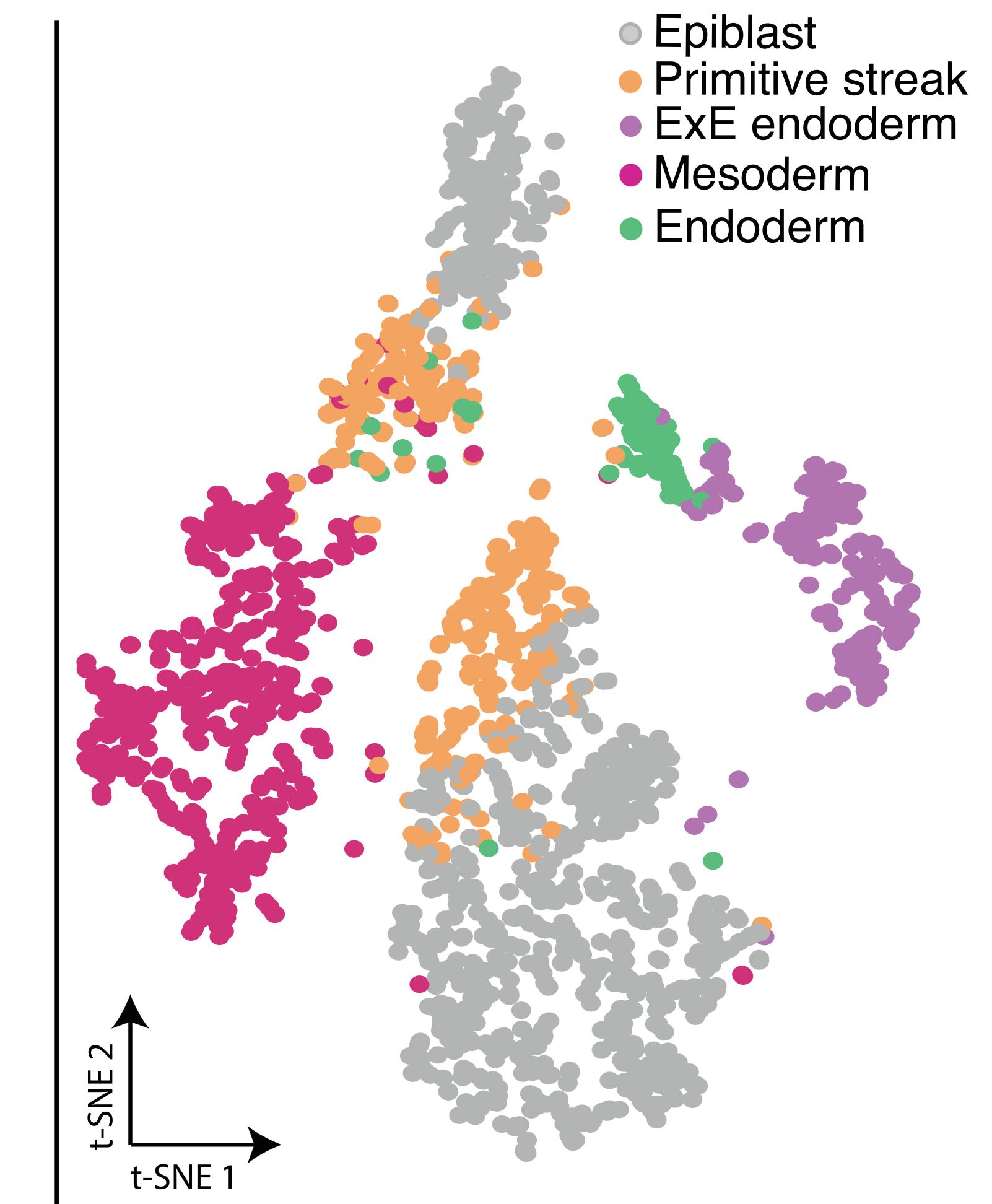
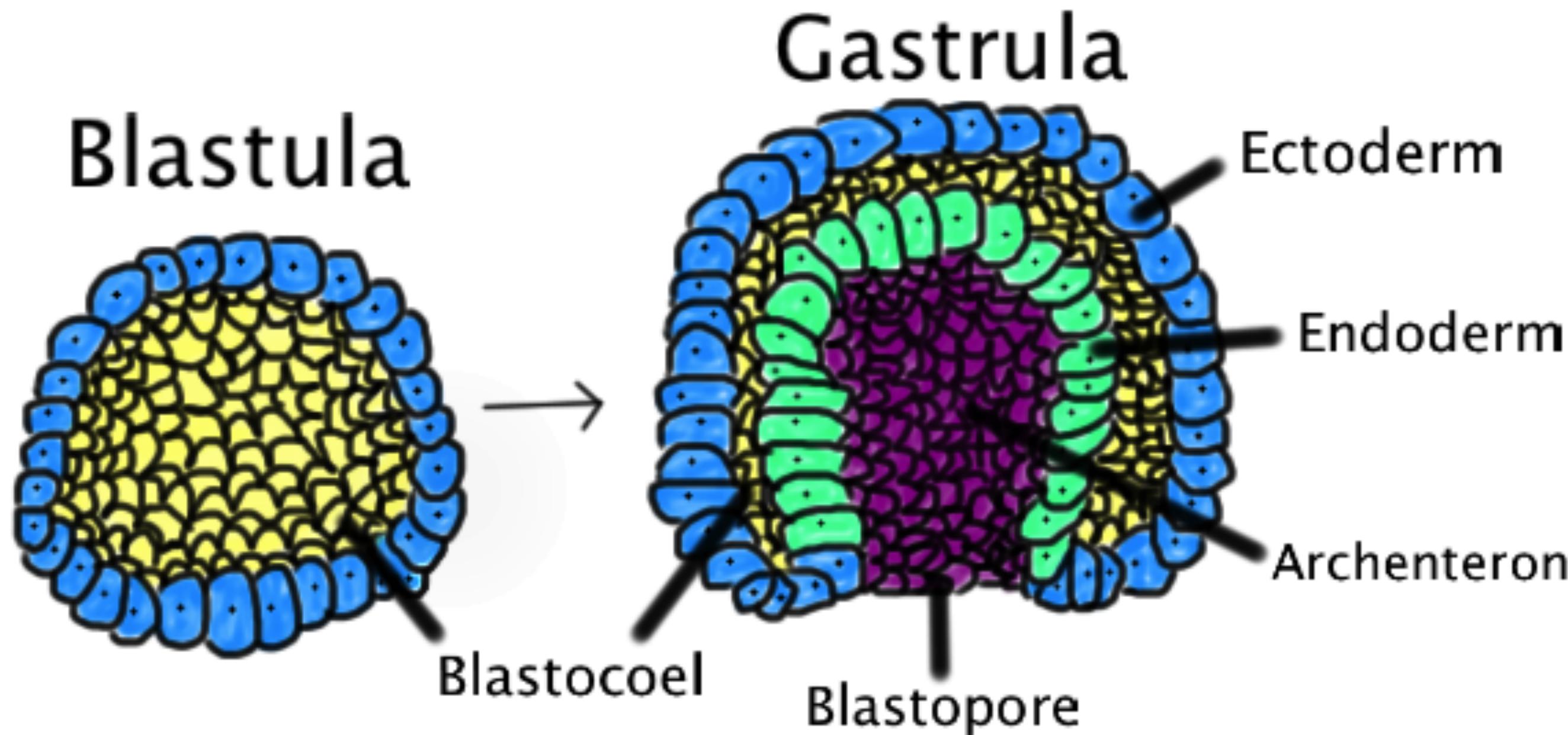


$$\mathbf{Y}^{(g,m)} = \mathbf{Z}^{(g)} \mathbf{W}^{(m)\top} + \boldsymbol{\epsilon}^{(g,m)}$$

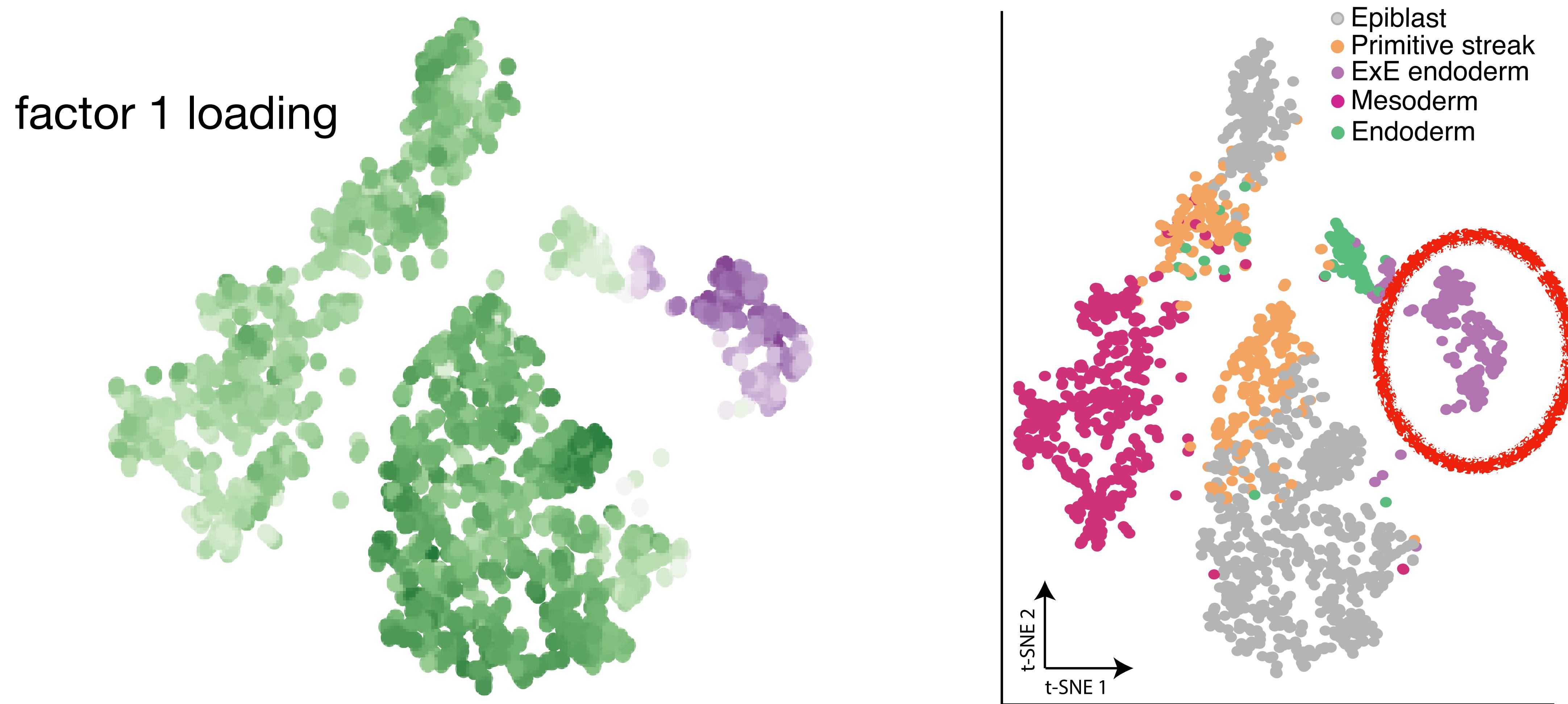
# MOFA+ can factorize group and modality variance



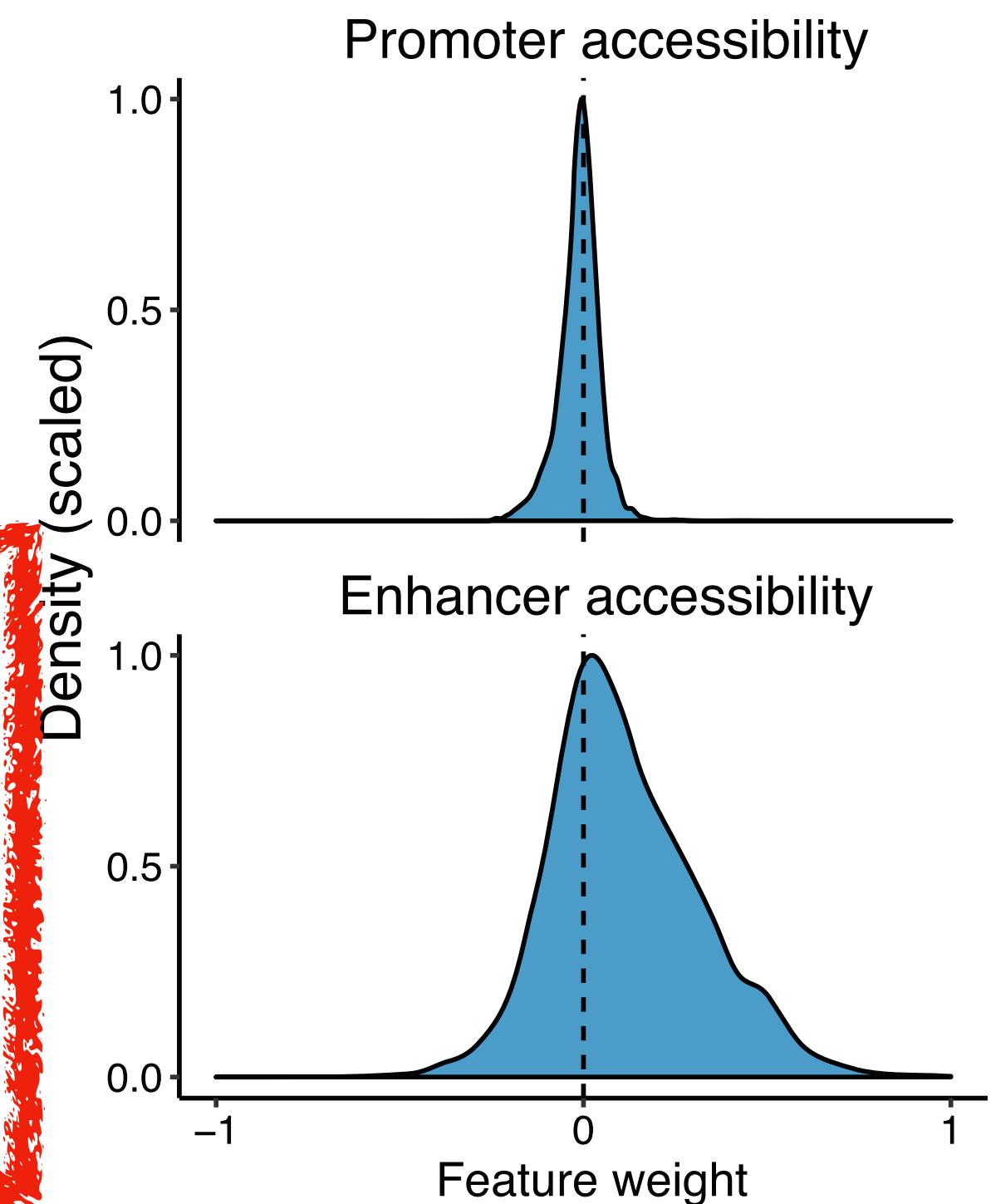
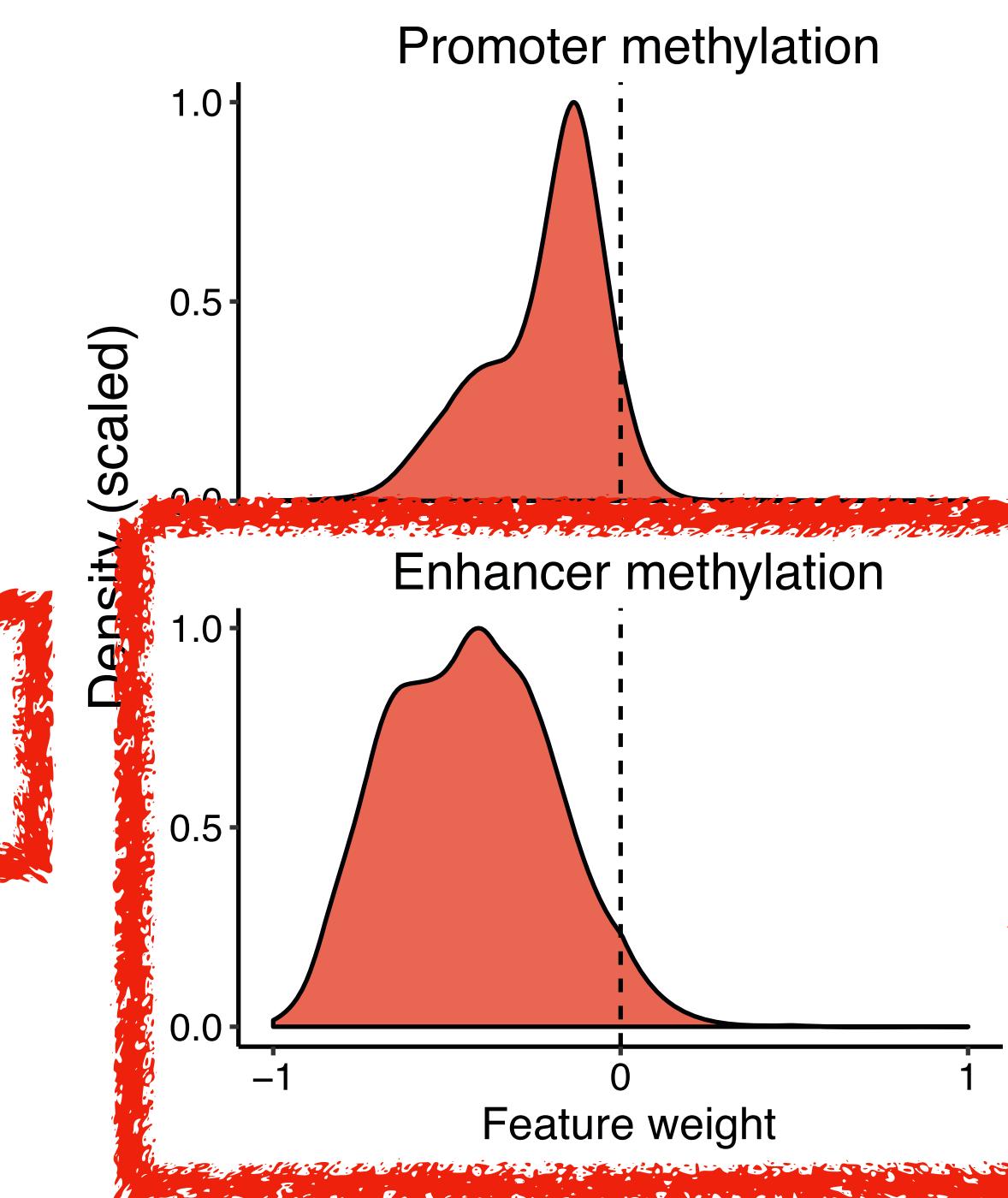
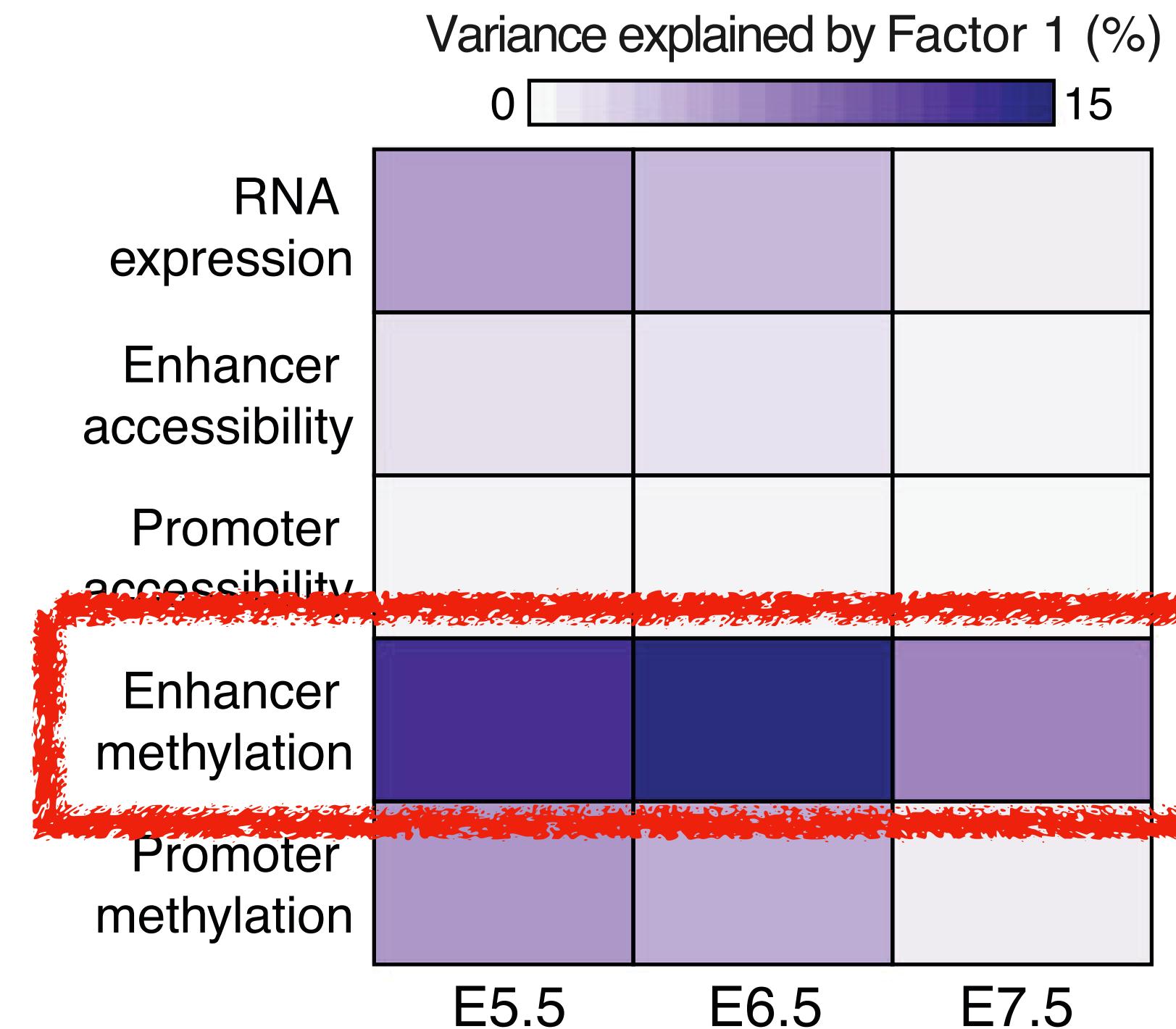
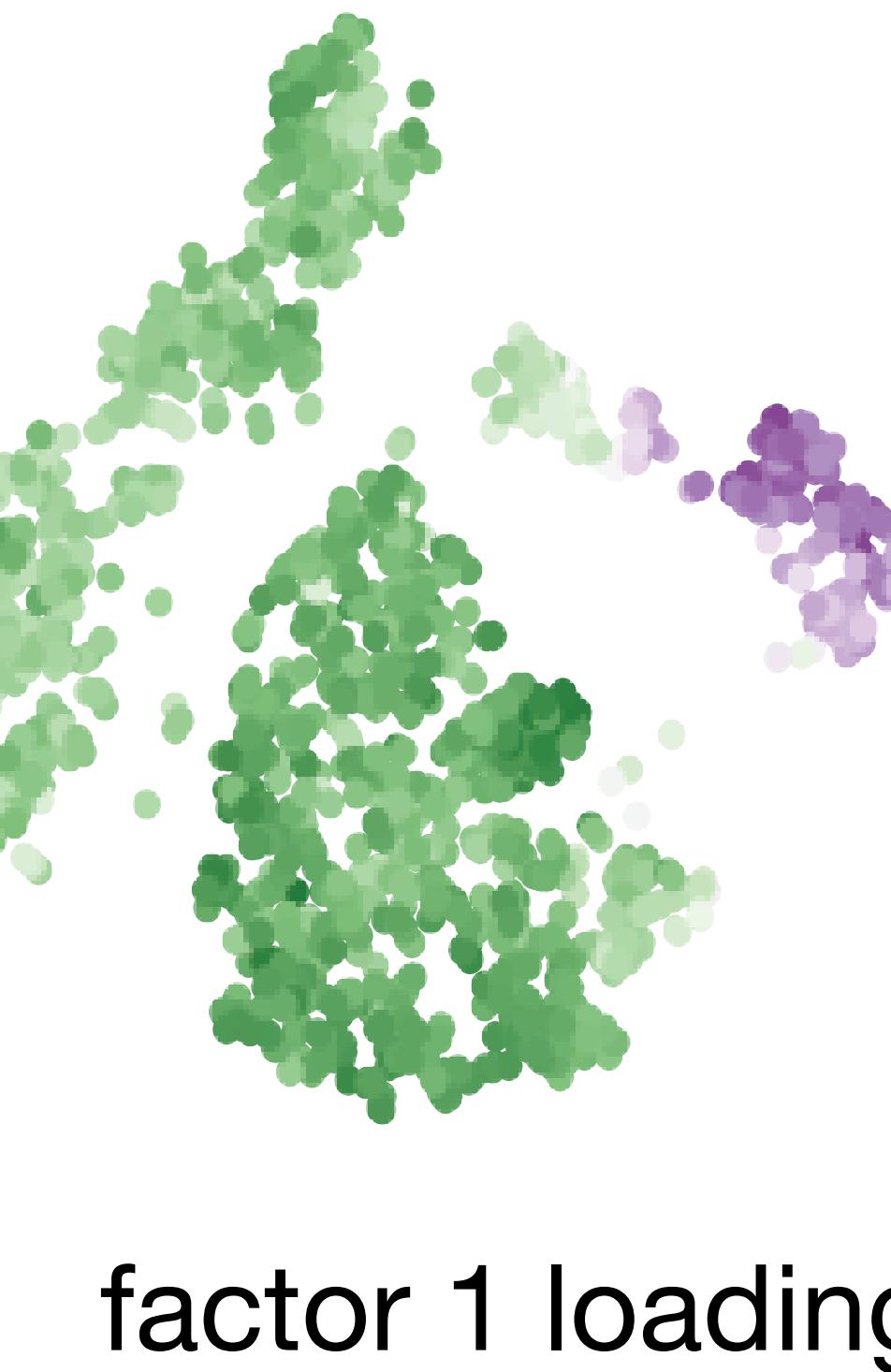
# An example of mouse gastrulation



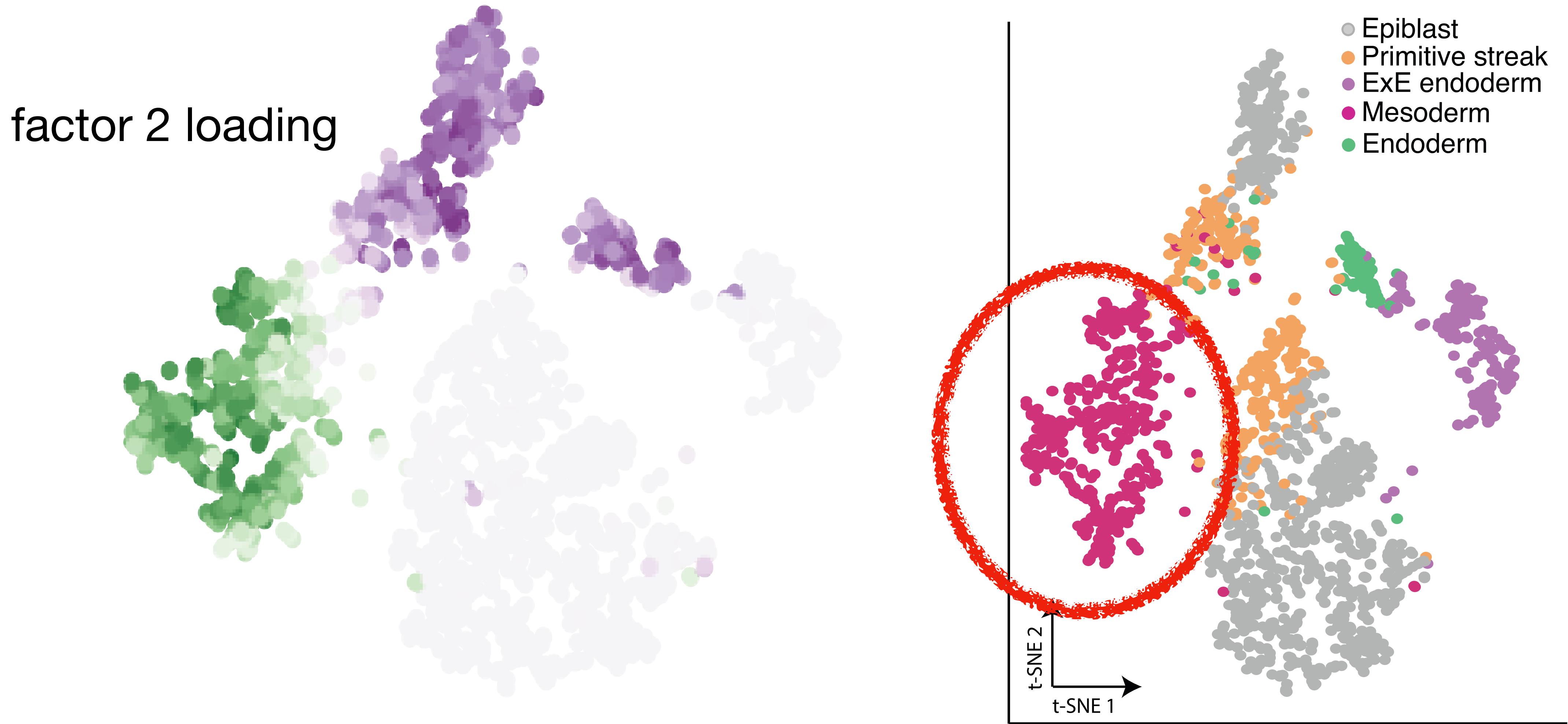
# Factor 1 $\approx$ ExE endoderm



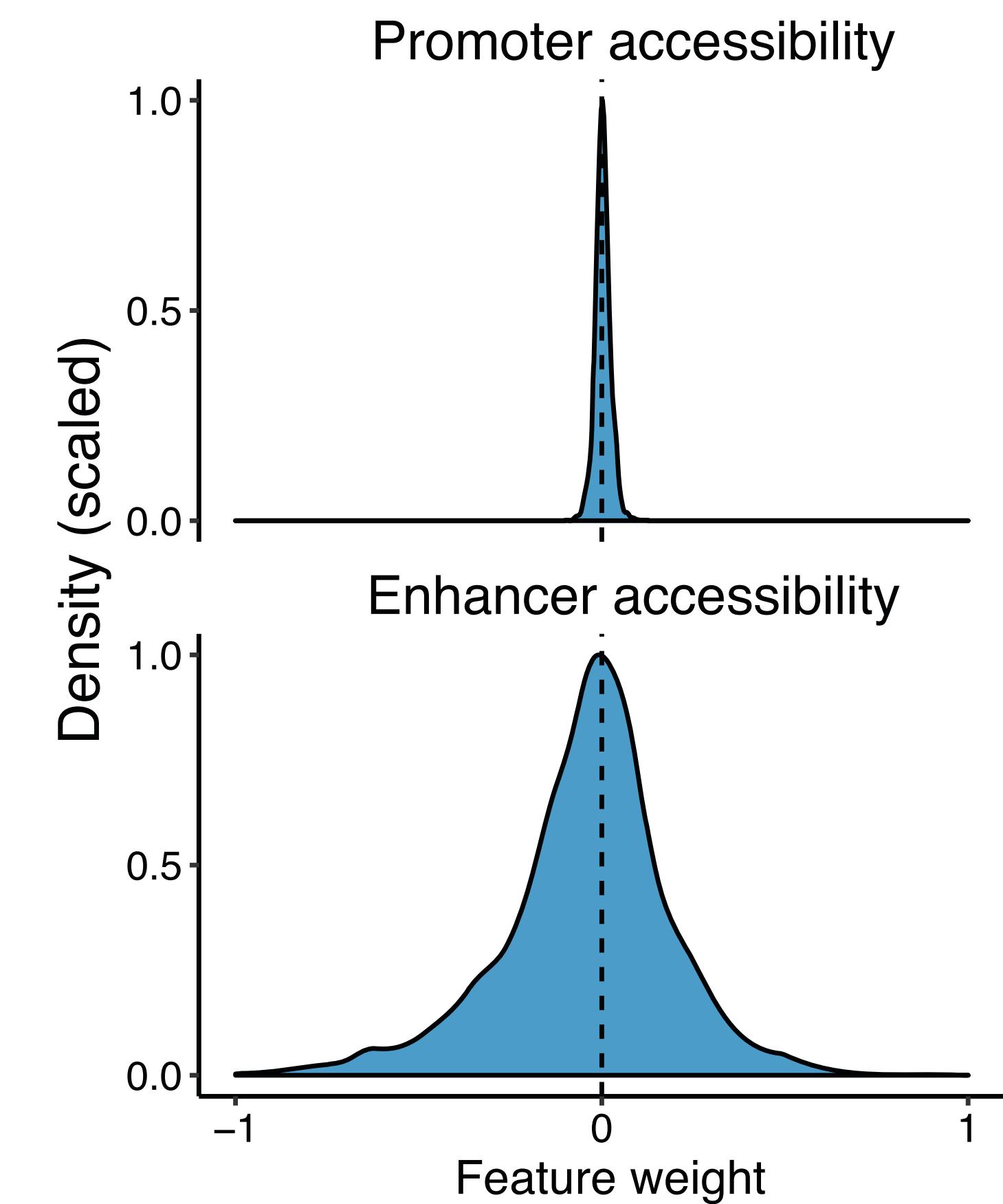
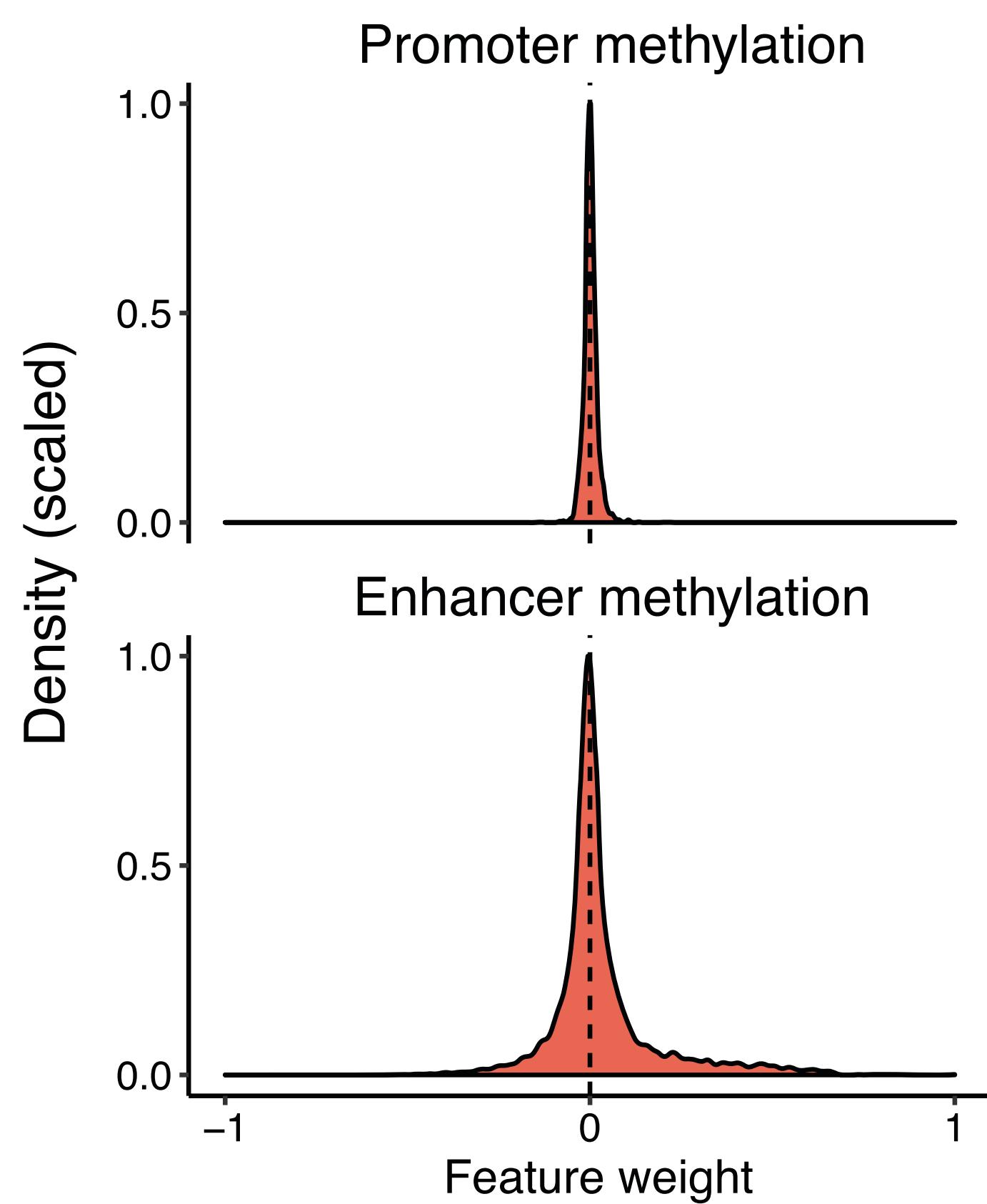
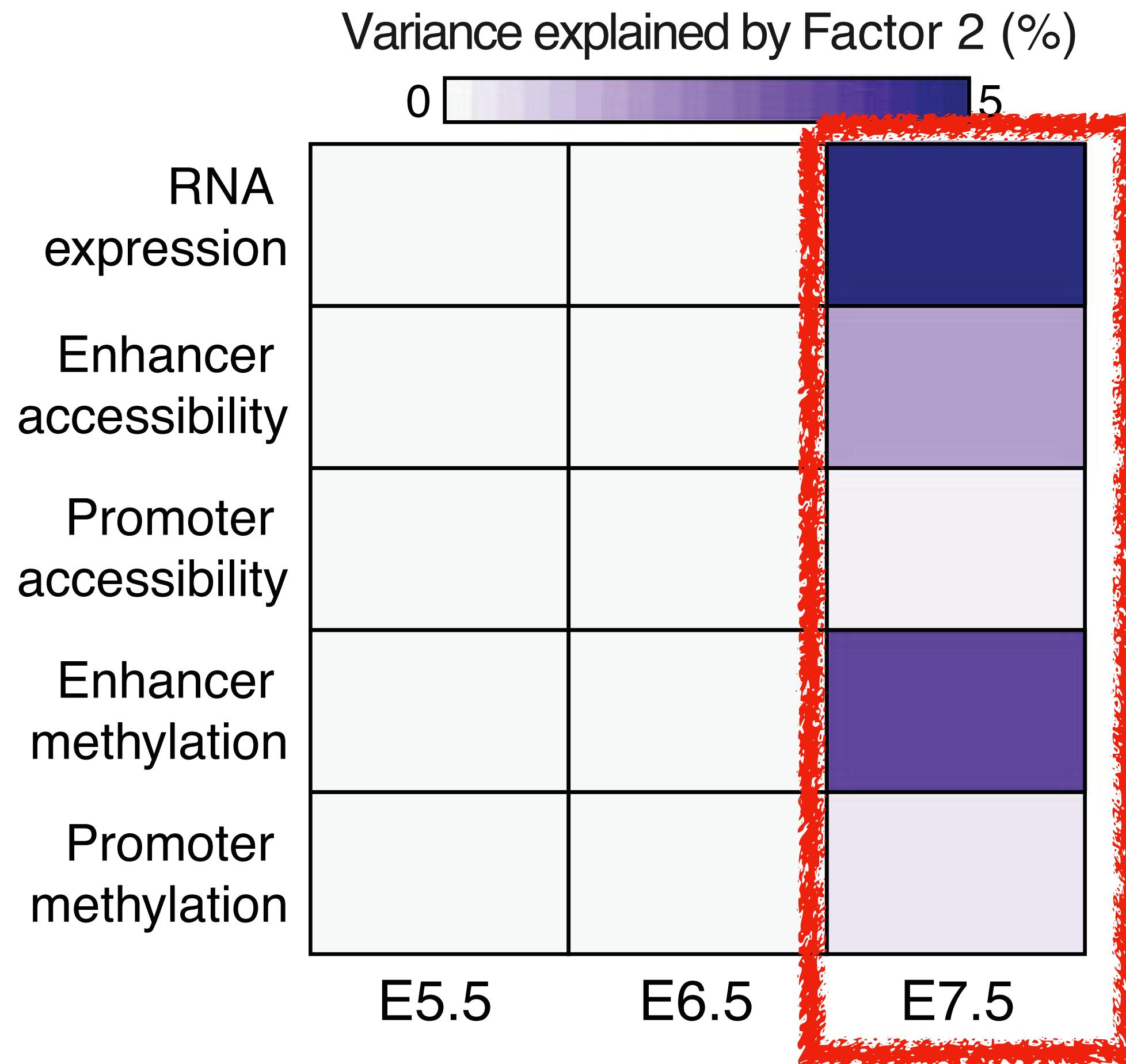
# Factor 1 $\approx$ ExE endoderm + Enh methylation



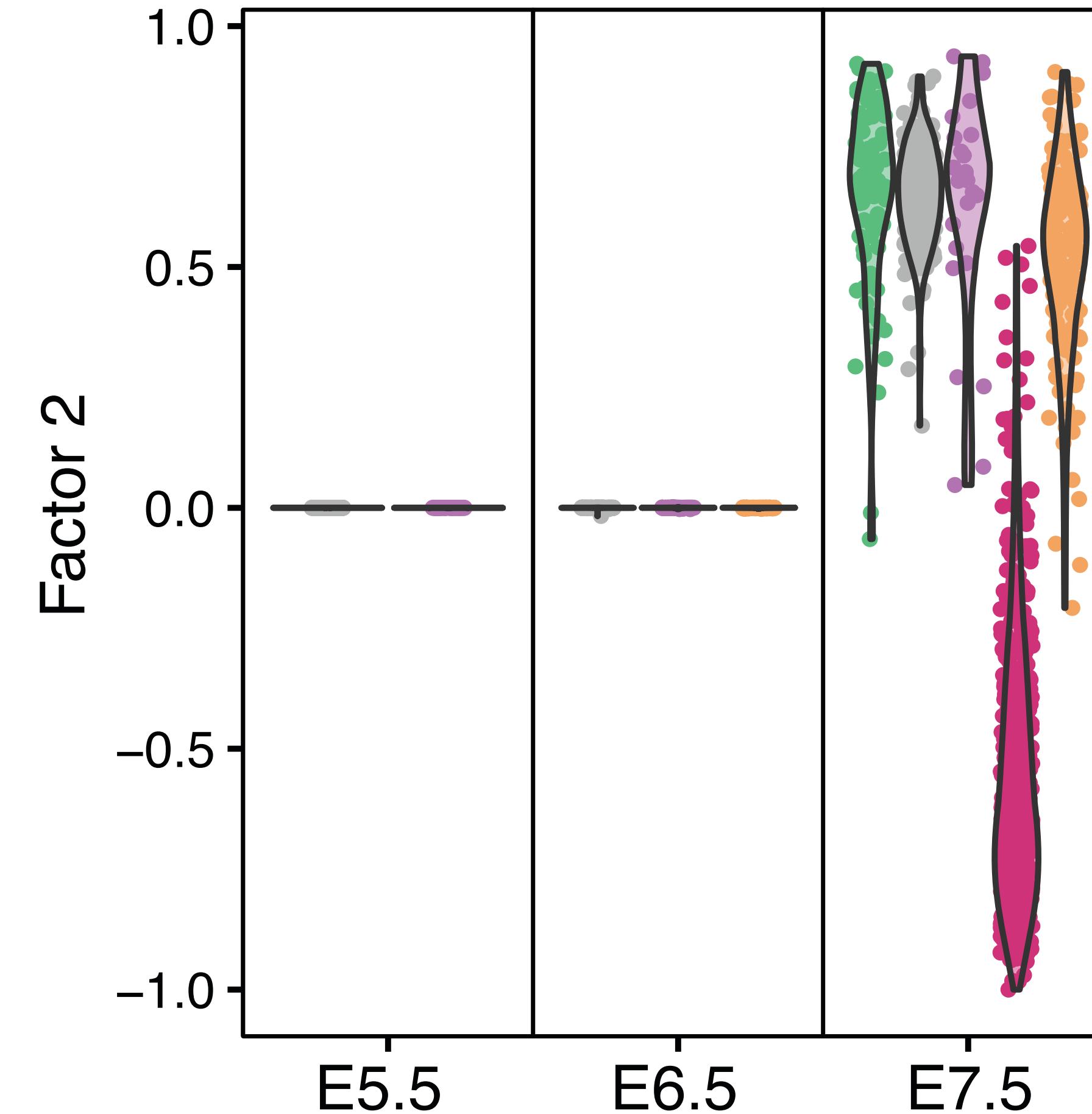
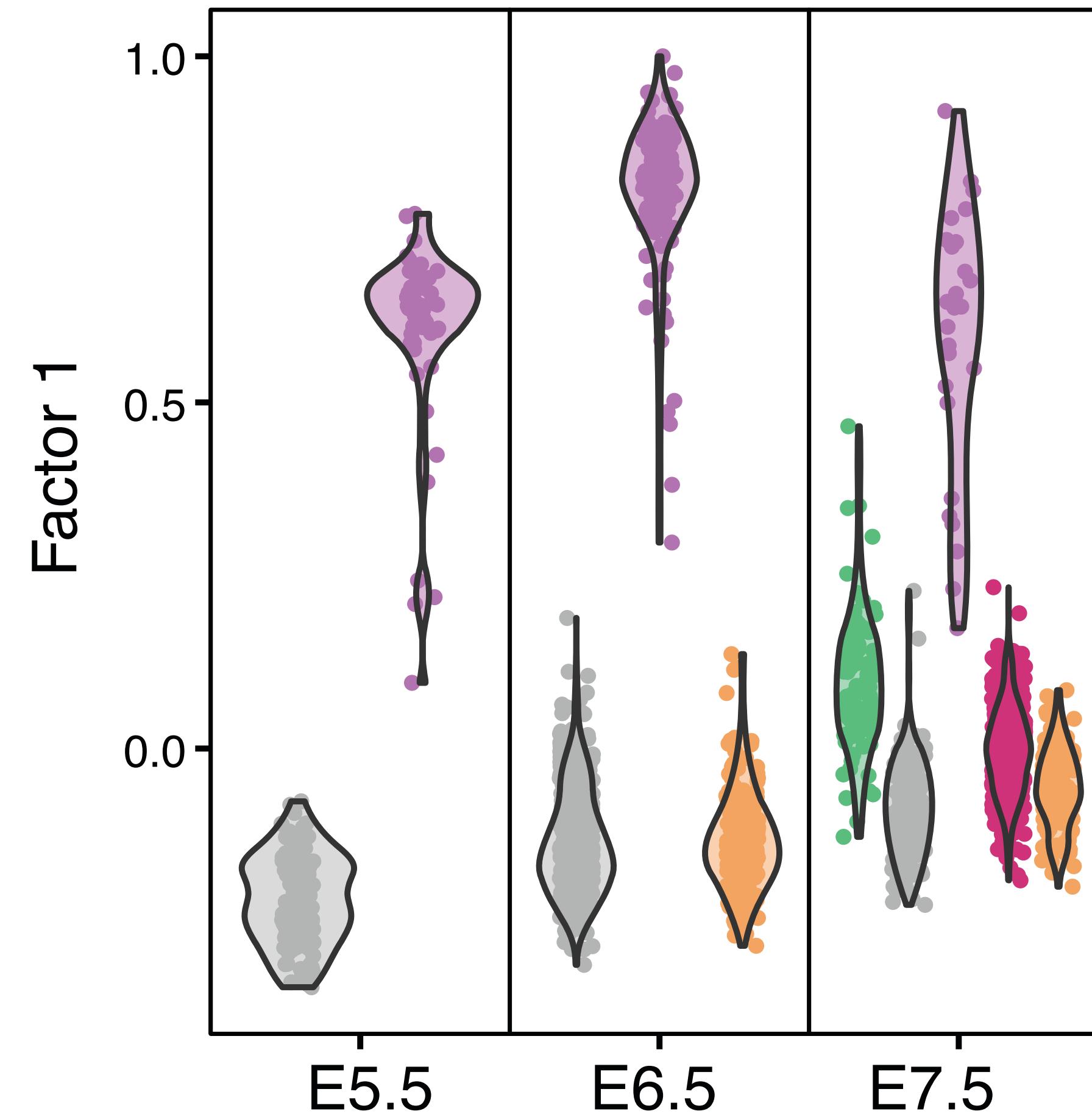
# Factor 2 ≈ Mesoderm



# Factor 2 ≈ mesoderm @ late stage

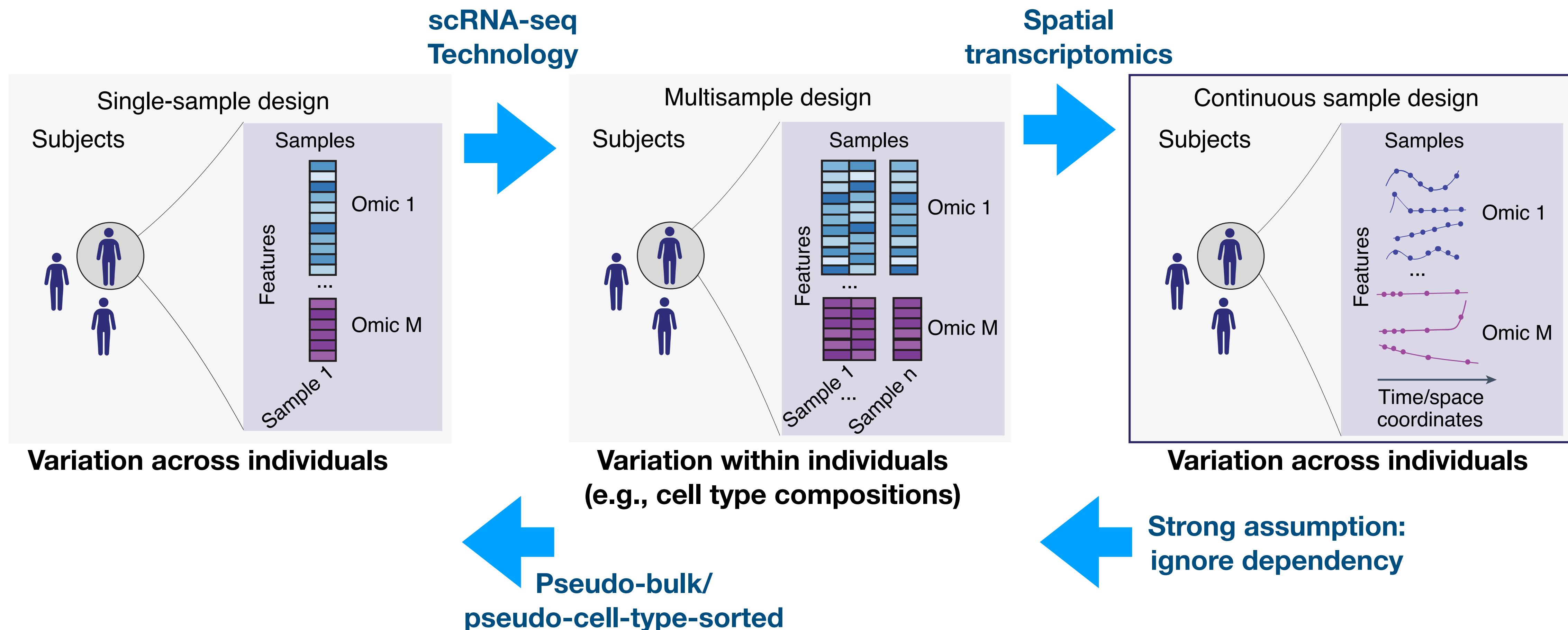


# A bit ambiguous interpretations

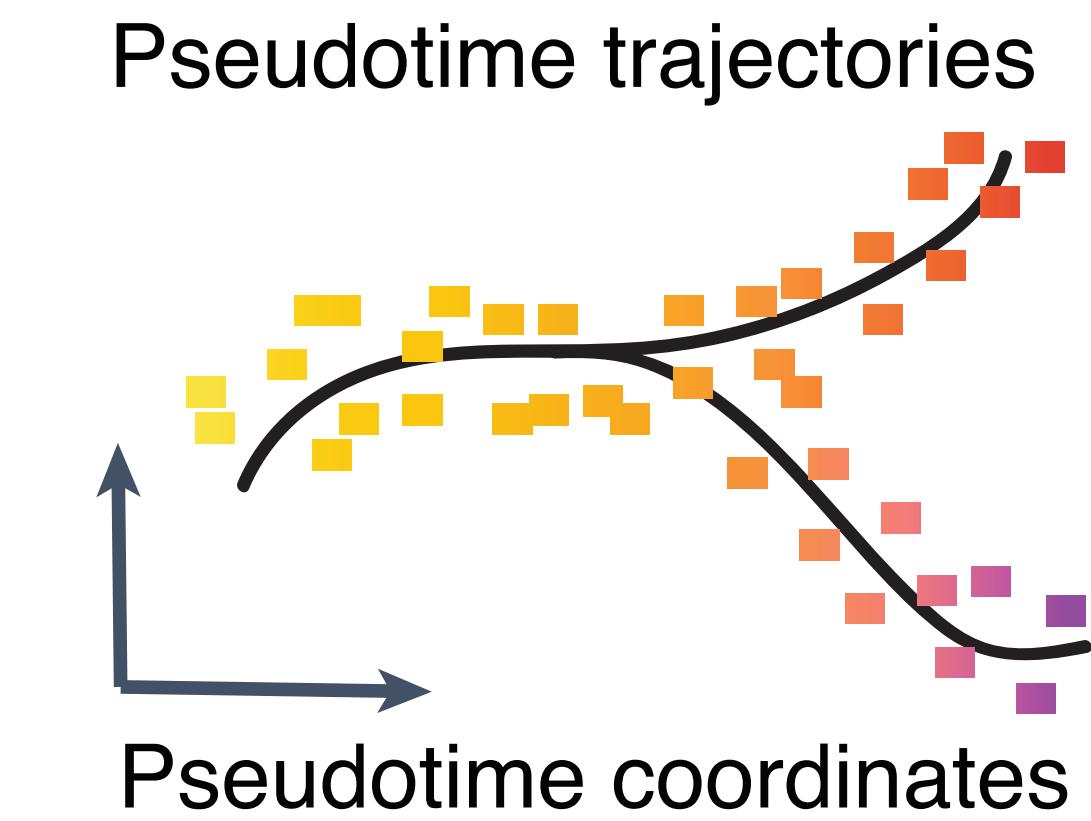
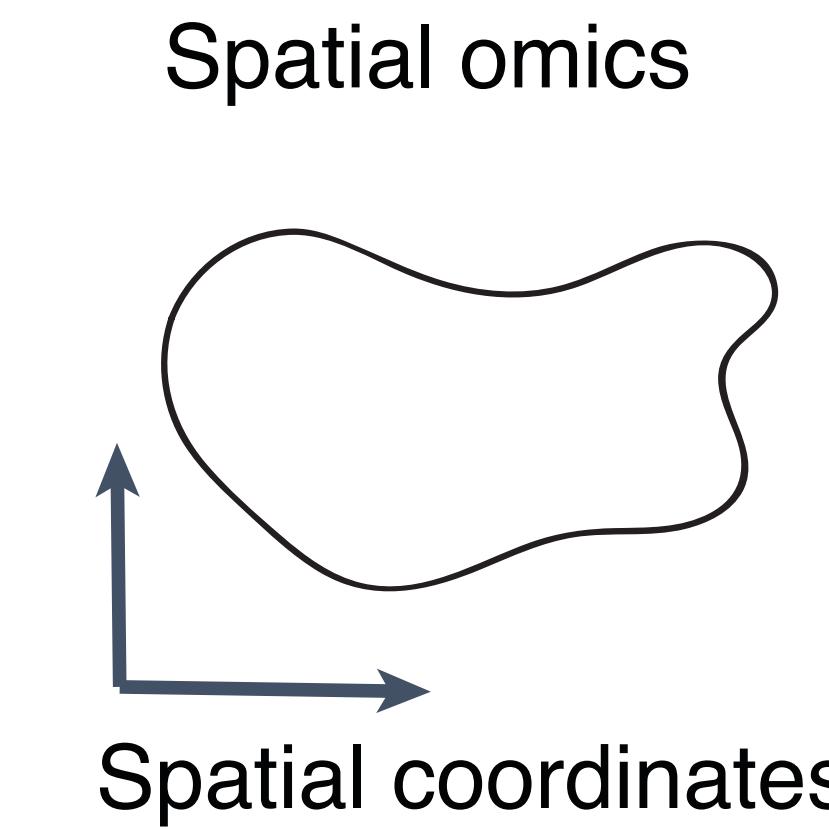
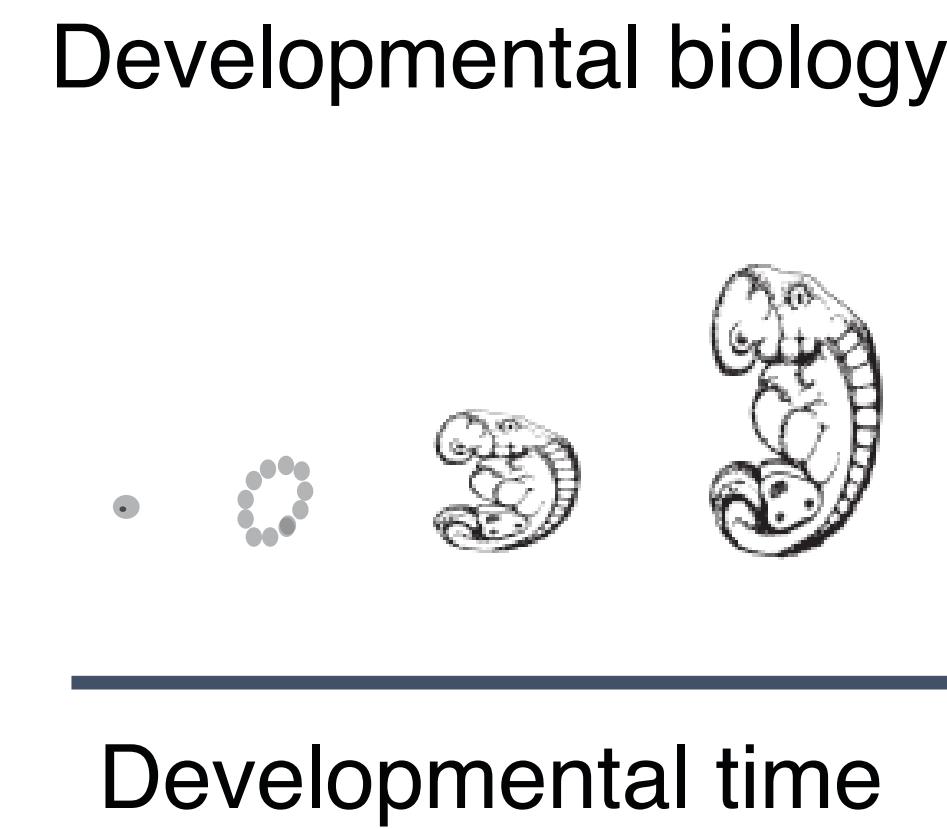
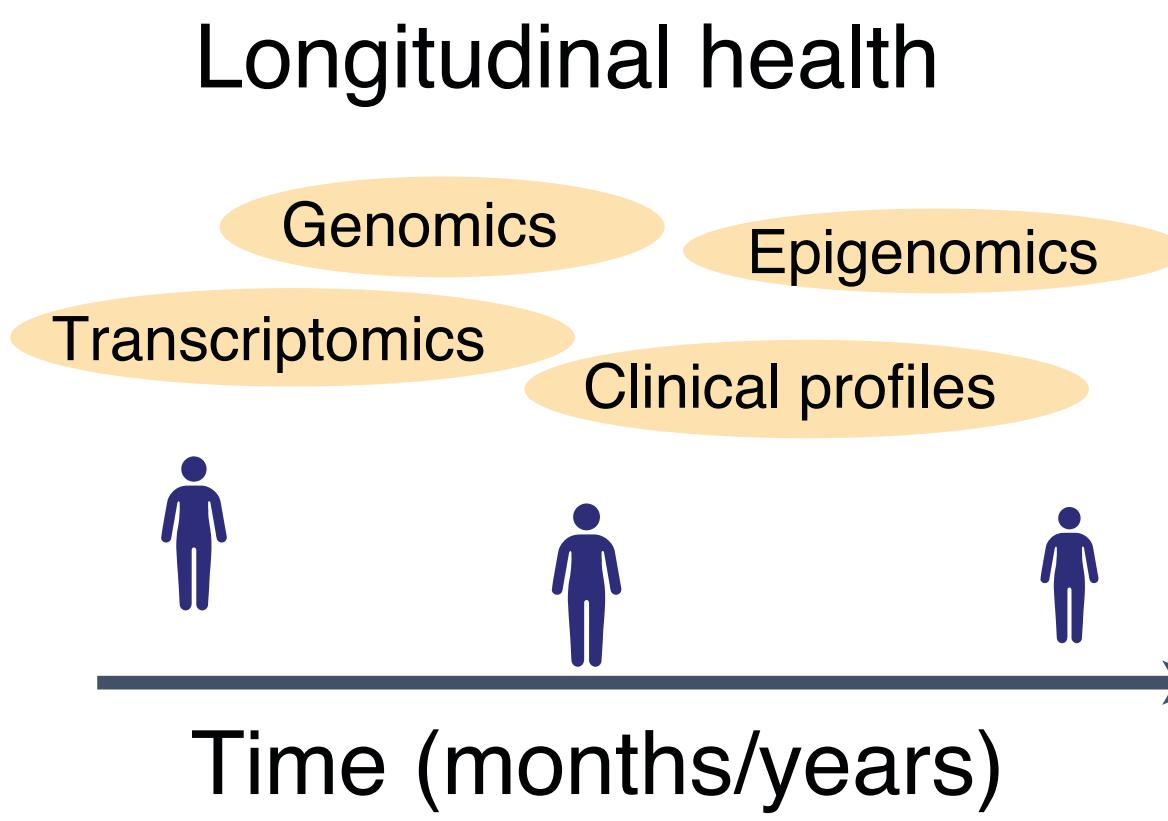


<https://biofam.github.io/MOFA2/>

# A new experimental design: many samples per individual



# Spatiotemporal dependency structures are embedded in biological systems



# How do we model dependencies/relationships?

Regression-based modeling

$$z = f(t) \text{ or } z = f(x)$$

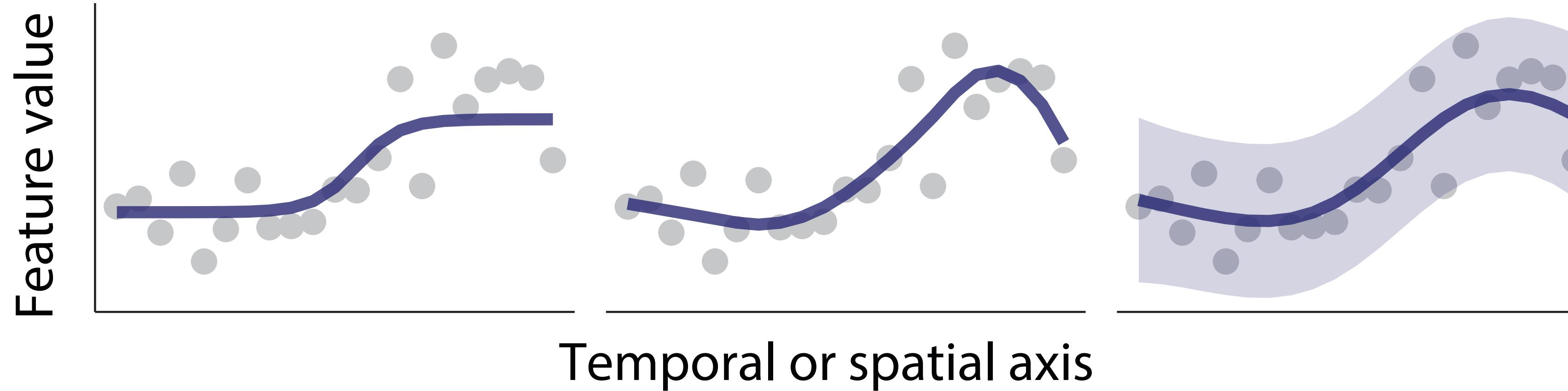
Parametric regression

e.g.,  $f(t) = (1 + e^{-t/a + b})^{-1}$     $f(t) = p_n(t), t \in [t_{n-1}, t_n]$

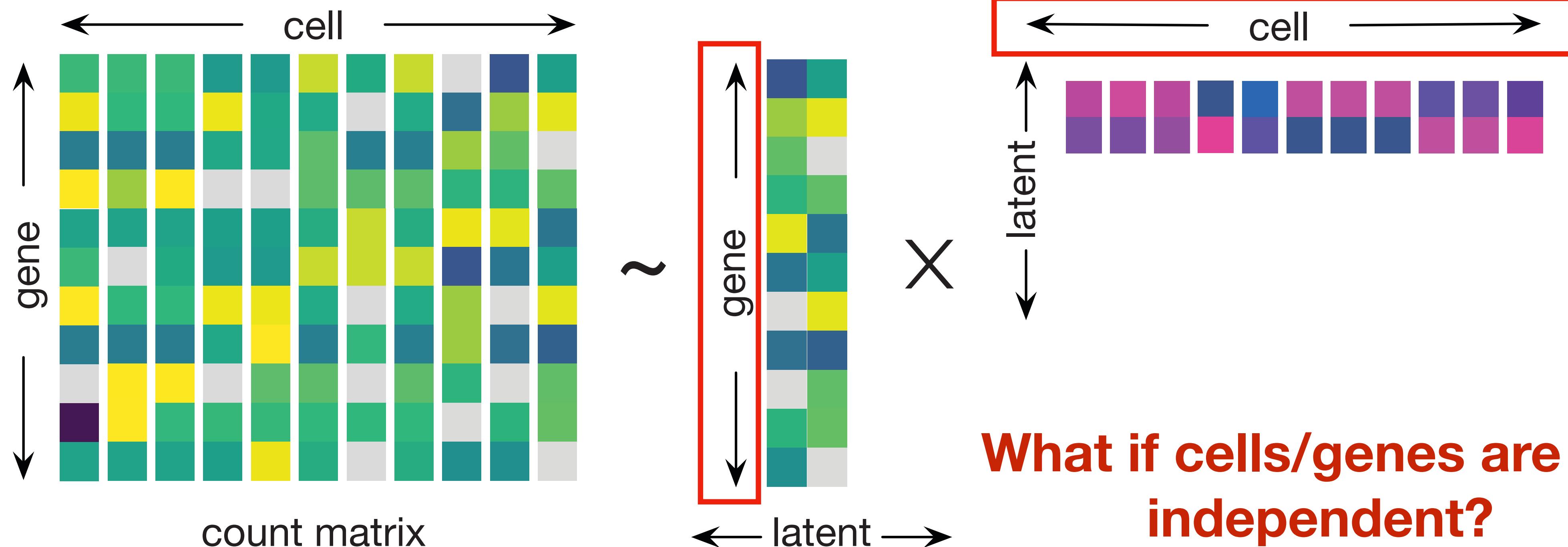
Splines

GP

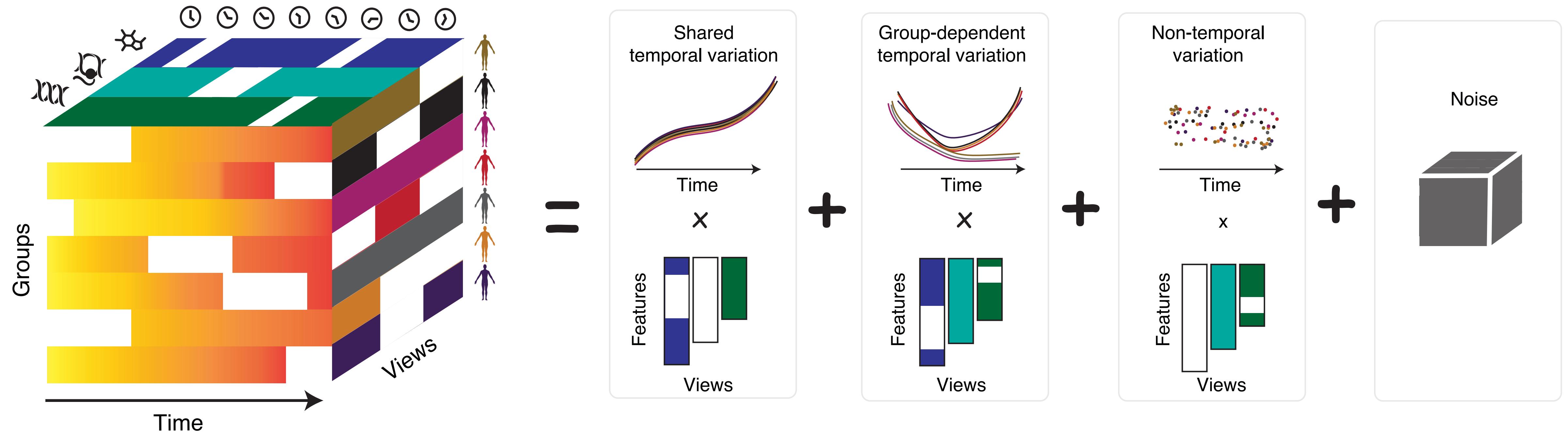
$$f \approx GP(\mu, \kappa)$$



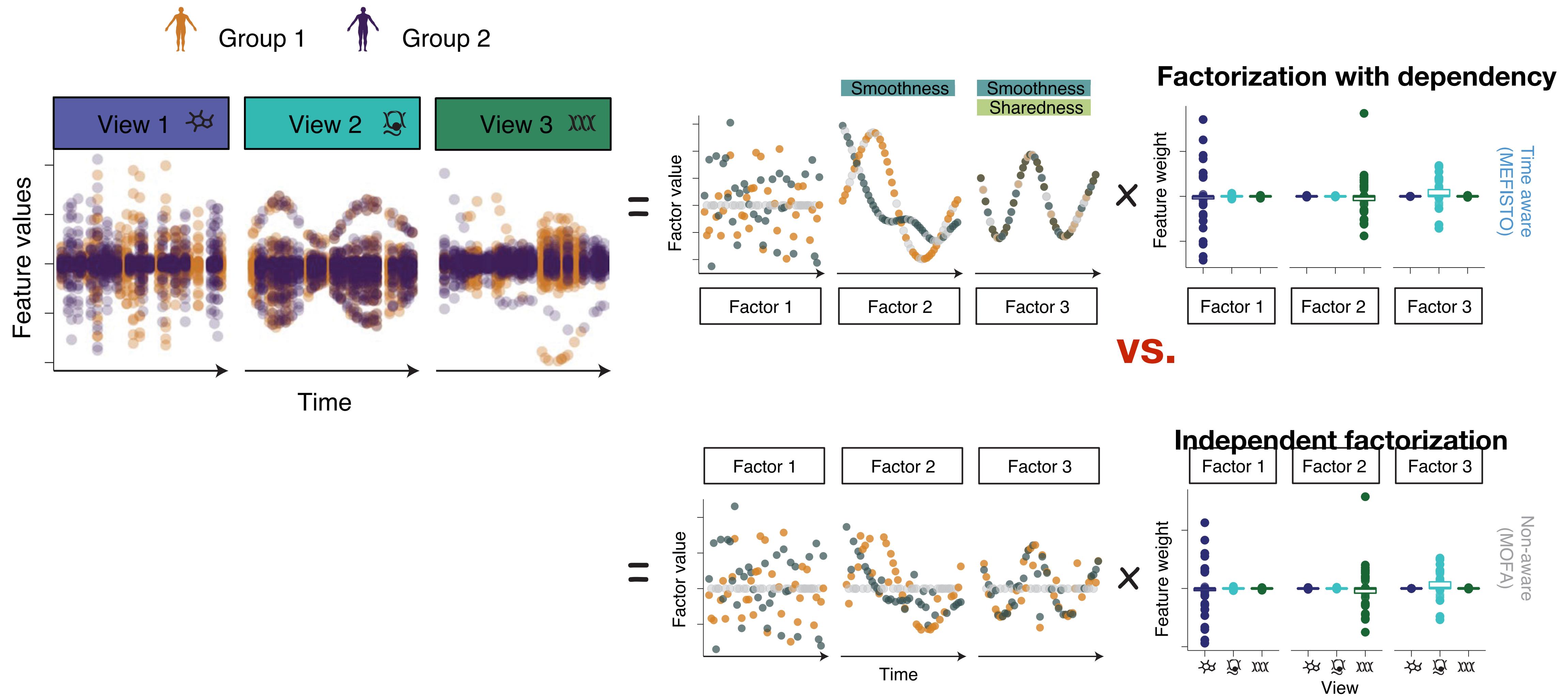
# Can we incorporate “smooth” temporal models in factorization settings?



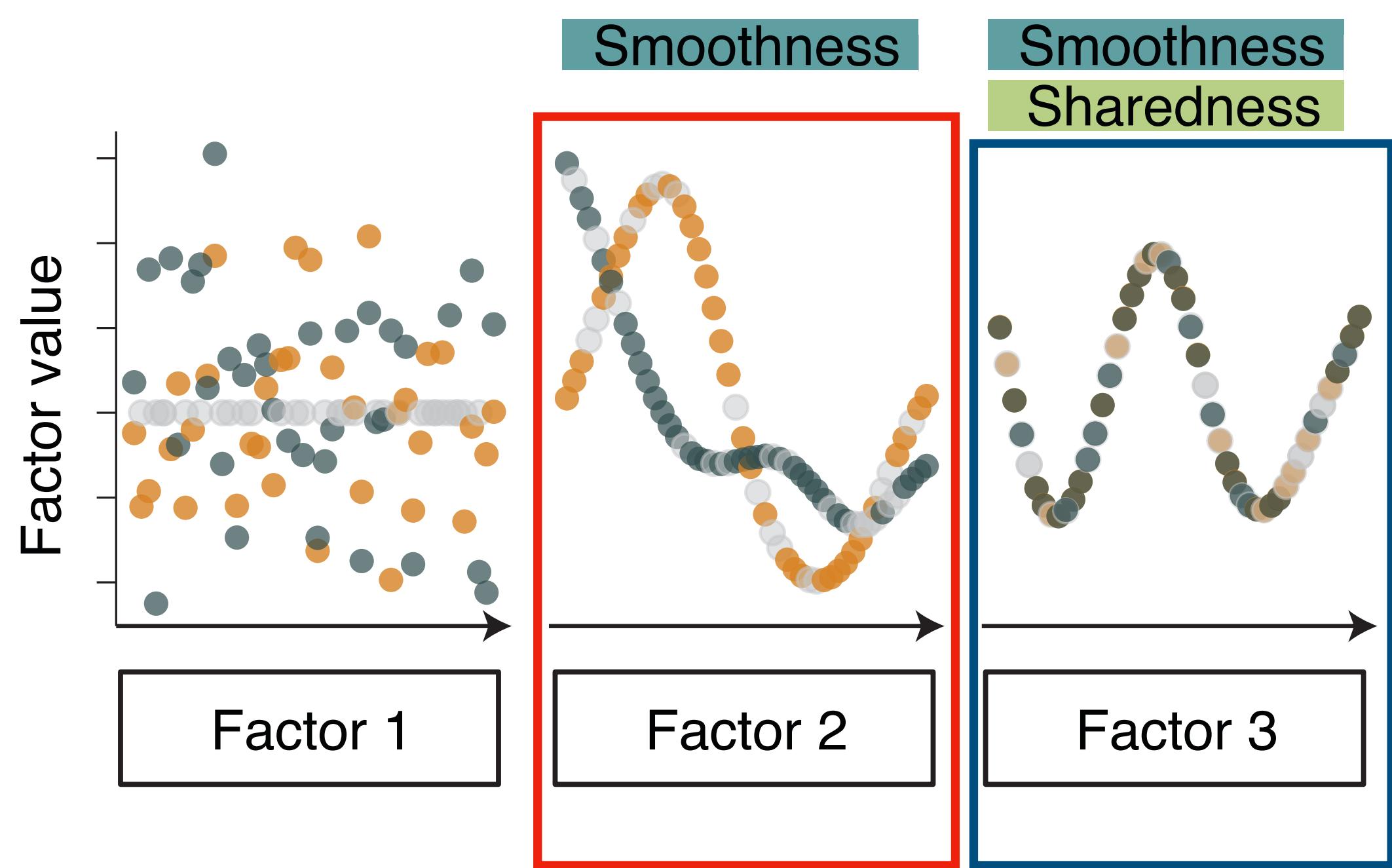
# MEFISTO: multi-view factor analysis + dependency modelling



# Smoothness assumption helps recover true latent factors and temporal dependency structures

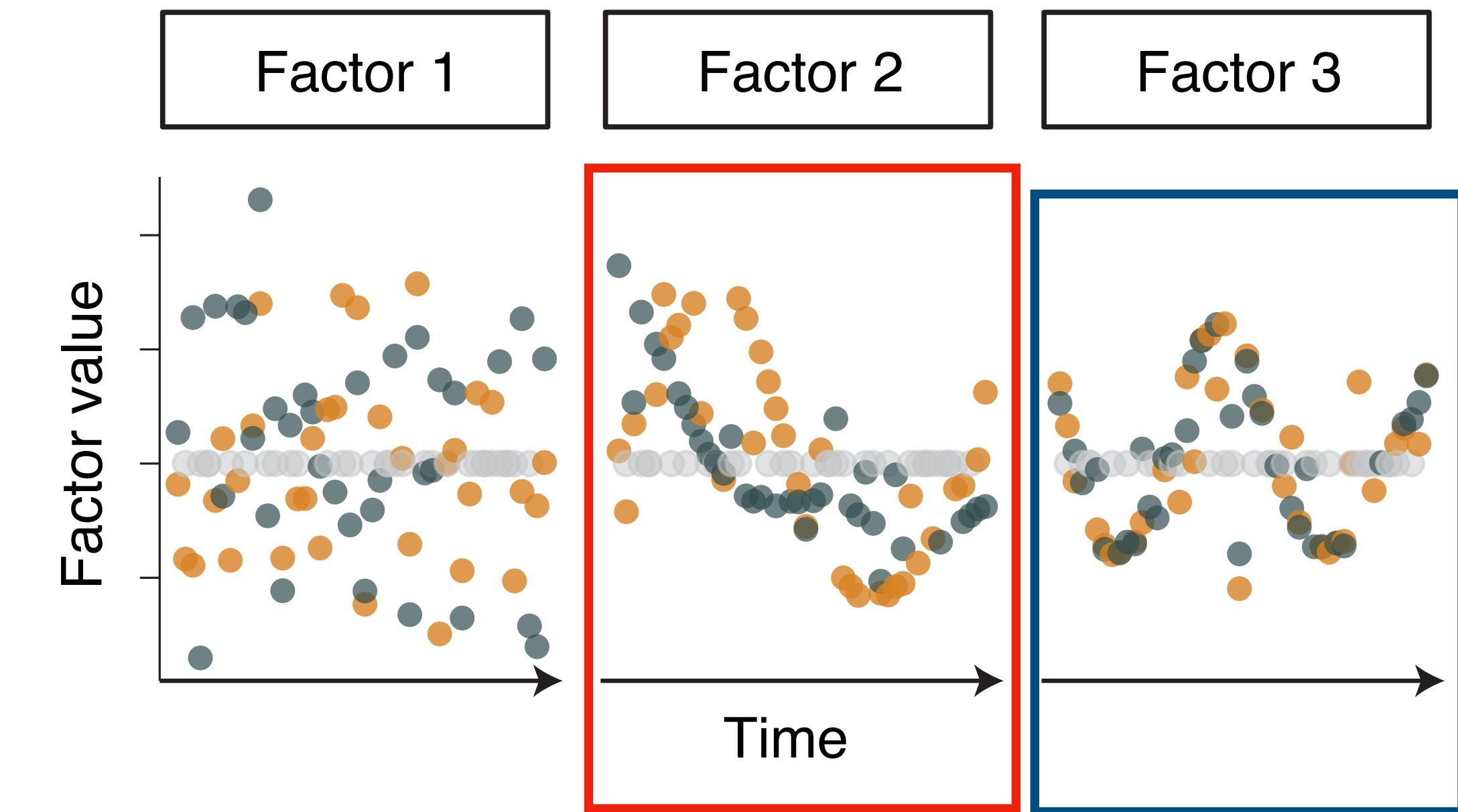


# Smoothness assumption helps recover true latent factors and temporal dependency structures



Factorization with dependency

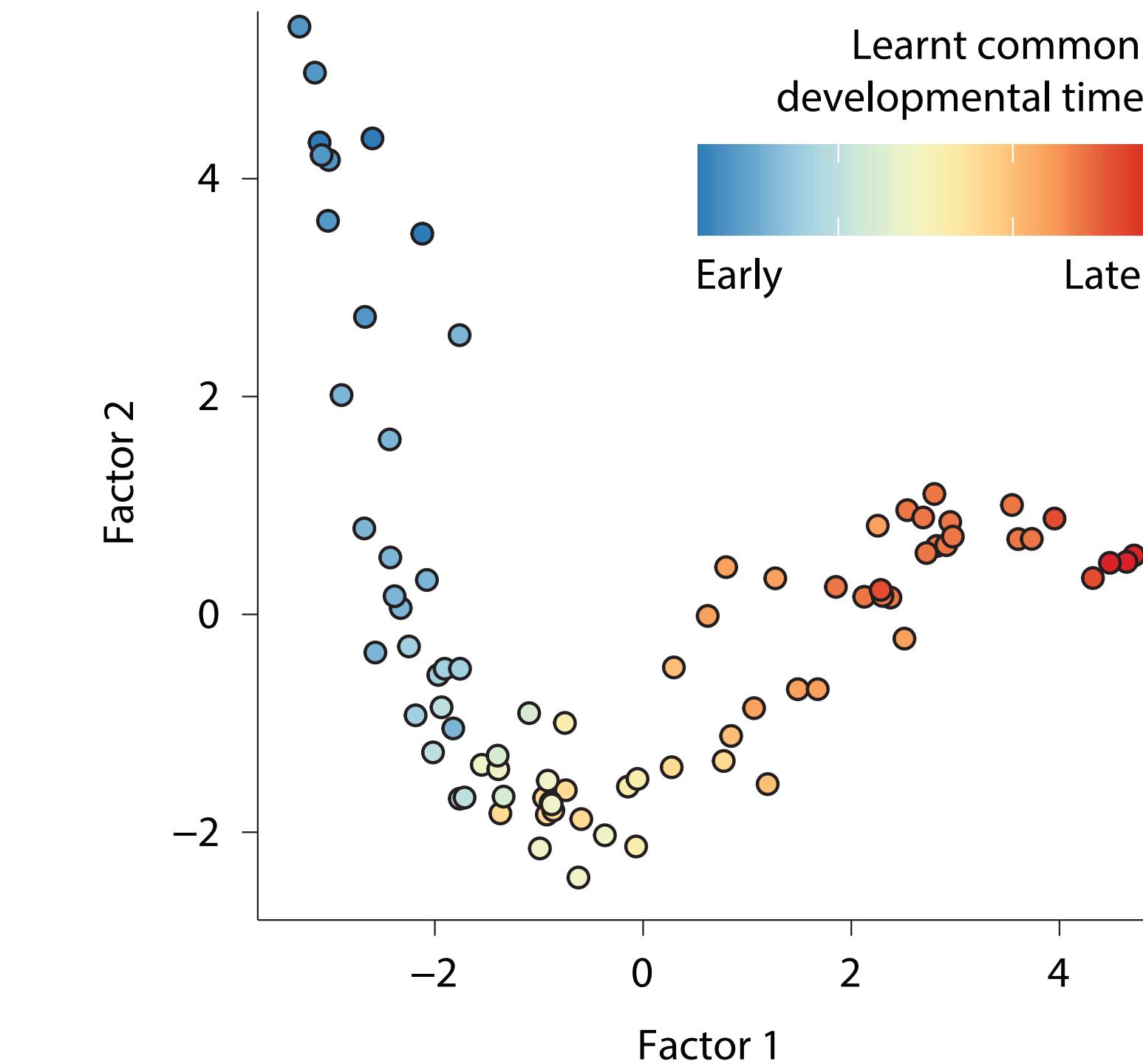
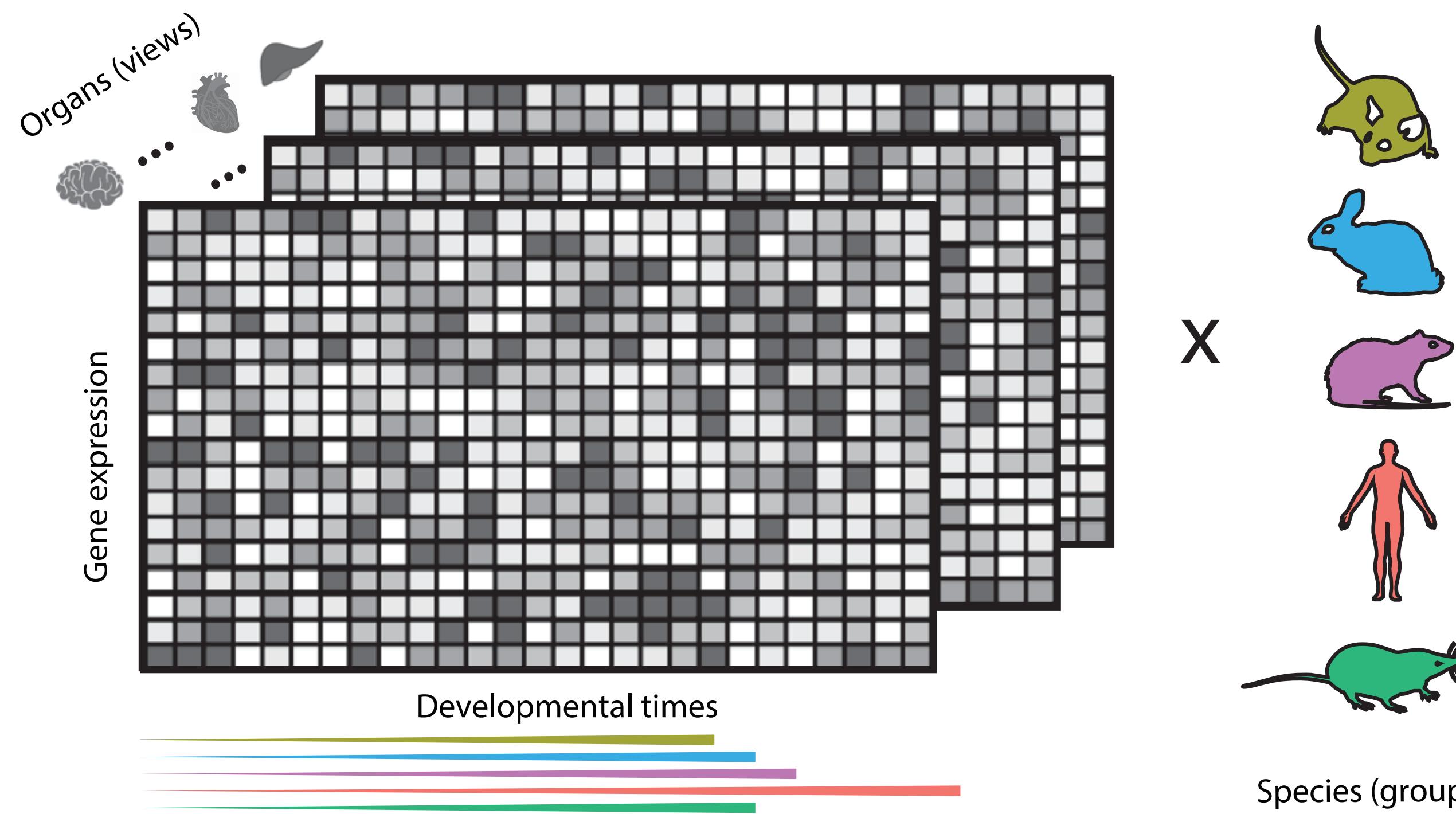
Note: the grey dots are unobserved



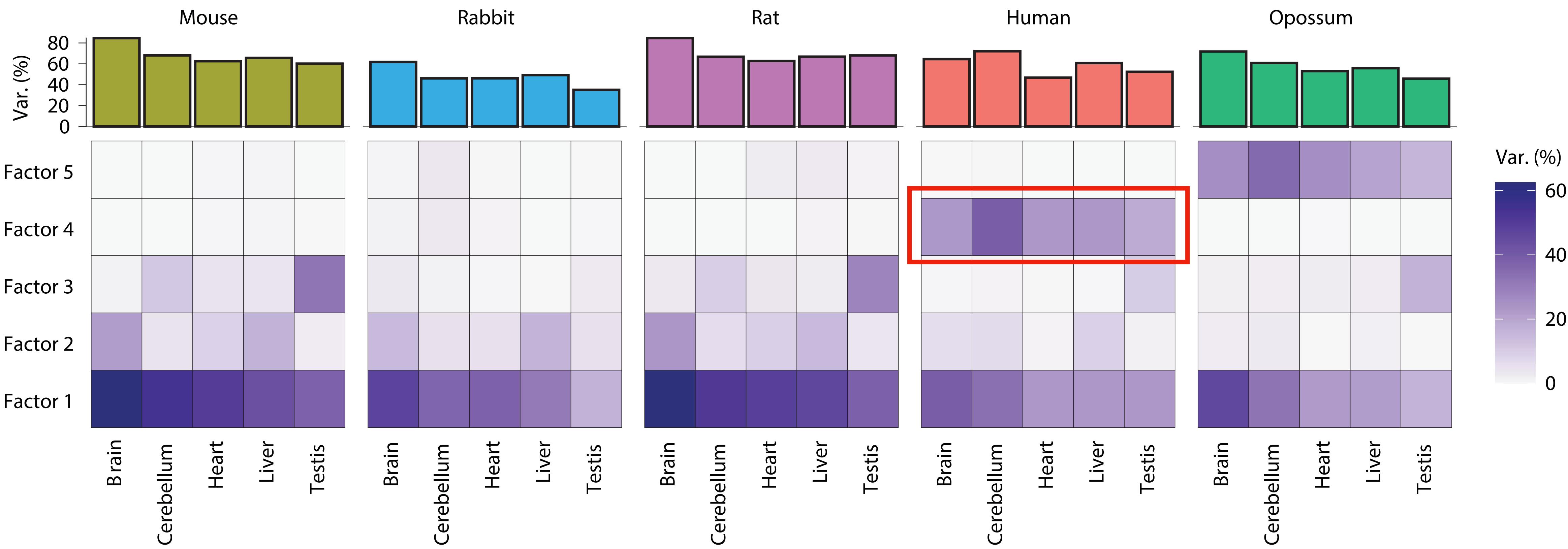
Independent factorization

Velten .. Stegle, Nature Methods (2022)

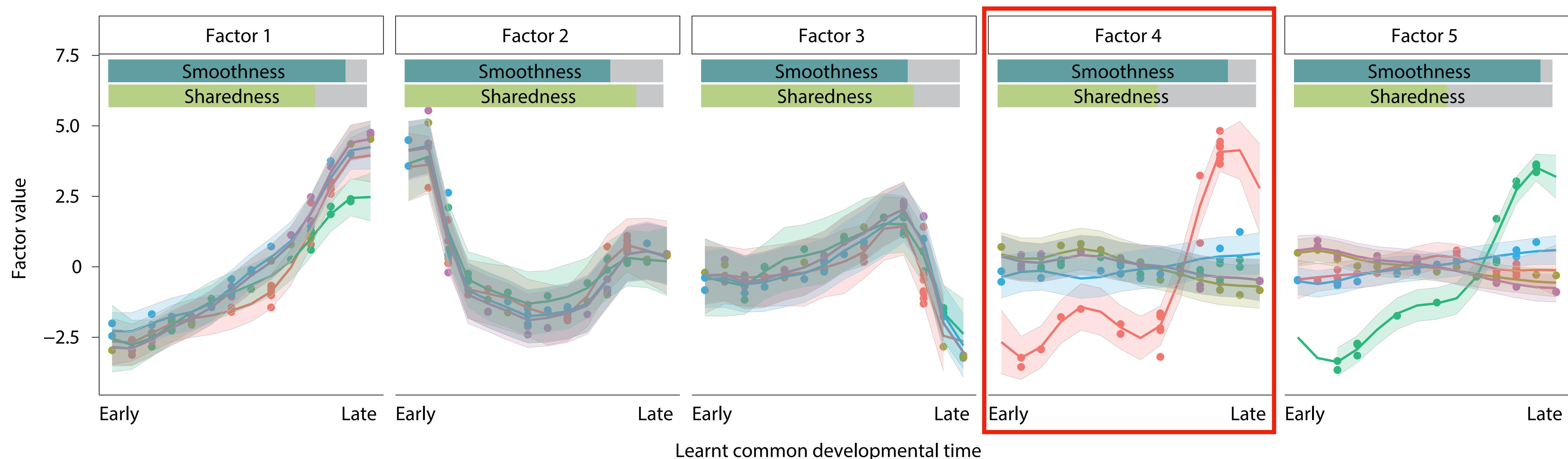
# MEFISTO to infer evolutionary gene programs



# What are human-specific gene factors/programs?



# Human-specific programs in the late developmental stages



<https://biofam.github.io/MOFA2/MEFISTO>

# MOFA vs. MEFISTO

$$\mathbf{Y}^{(g,m)} = \mathbf{Z}^{(g)} \mathbf{W}^{(m)^\top} + \boldsymbol{\epsilon}^{(g,m)}$$

MEFISTO

$$Z_{nk}^{(g)} = f_k(\mathbf{c}_n^{(g)}) + \boldsymbol{\epsilon}_{nk}^{(g)}$$

Gaussian  
Process  
(function  
approx)

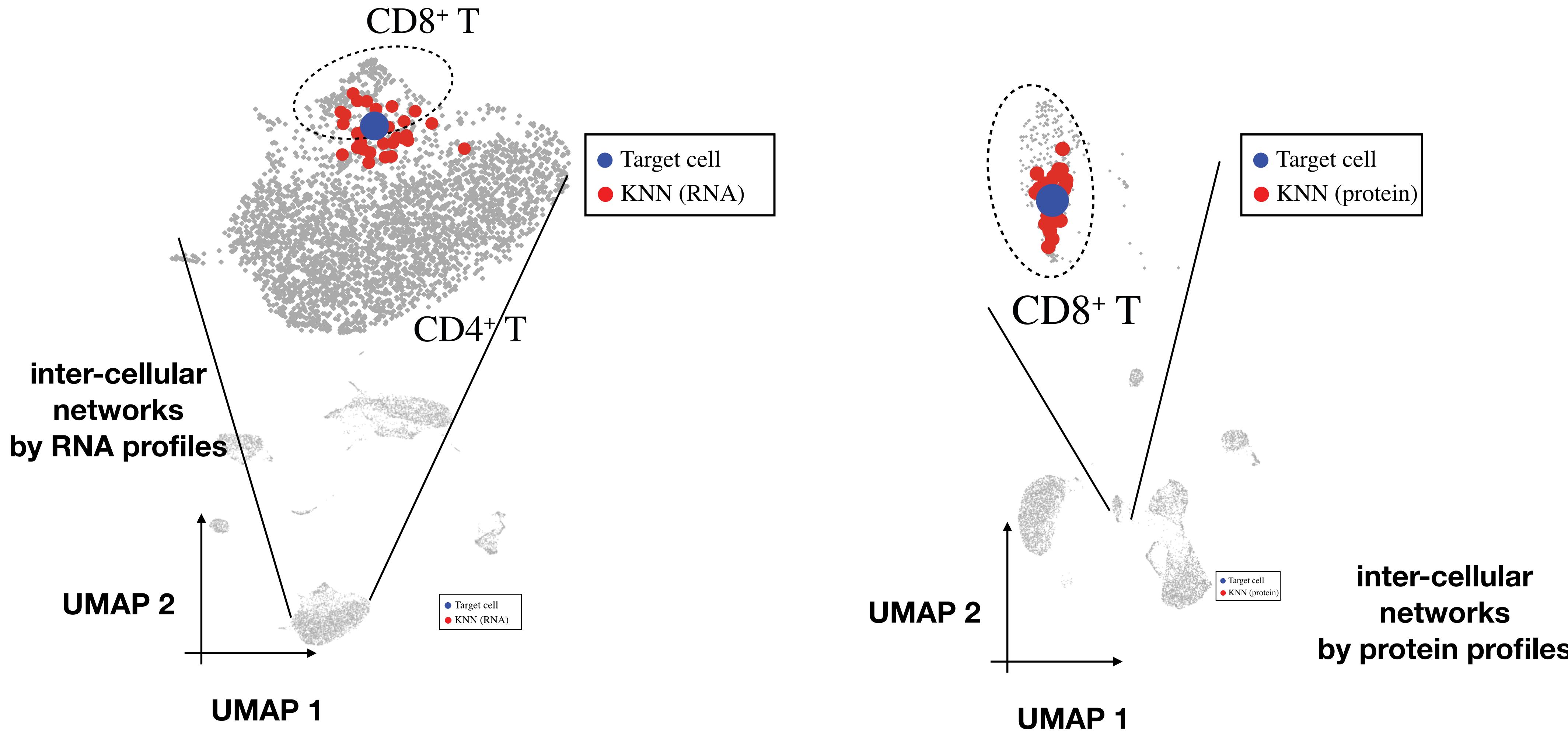
covariate  
manifesting  
the underlying  
dependency

# **Discussion: Pros and cons for MOFA**

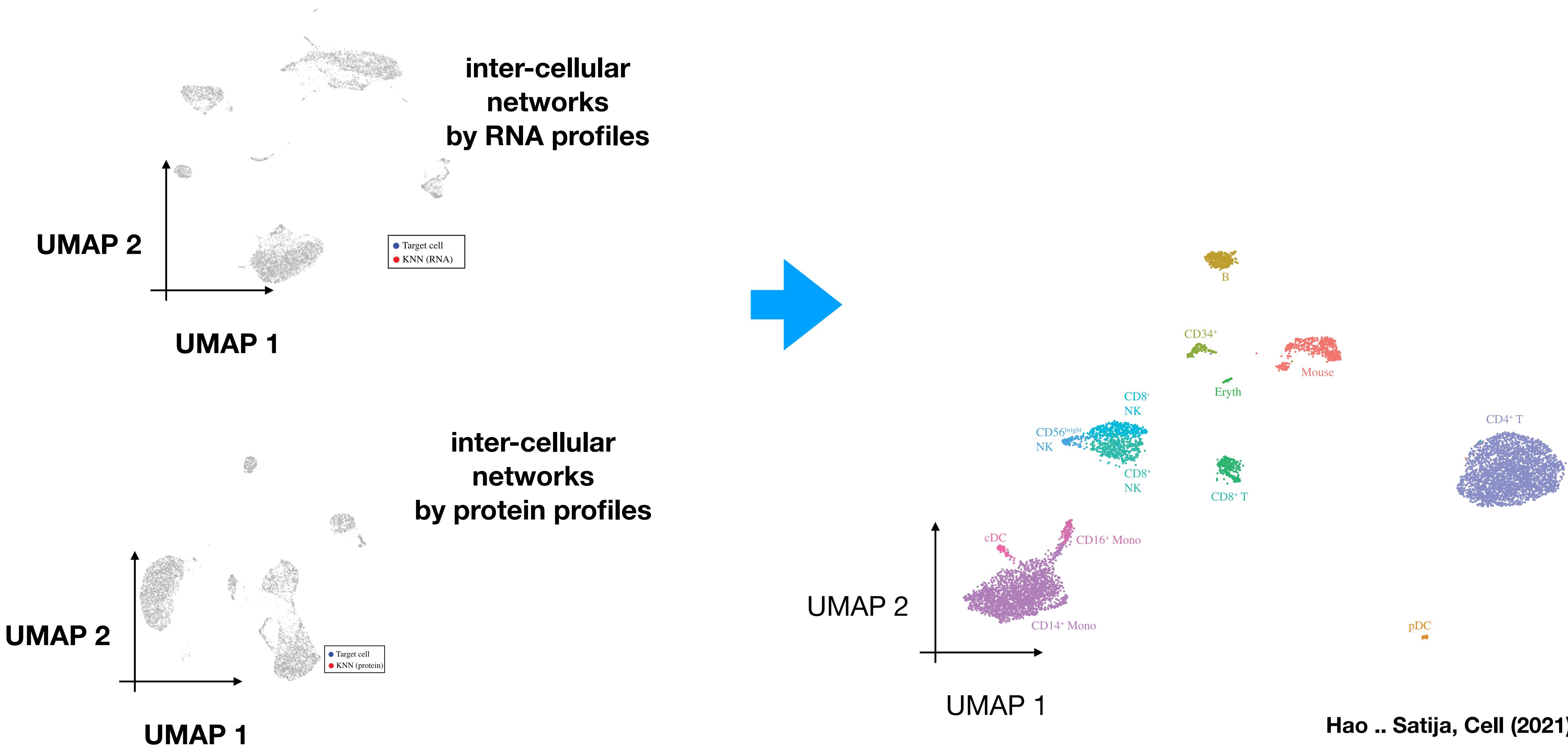
# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms

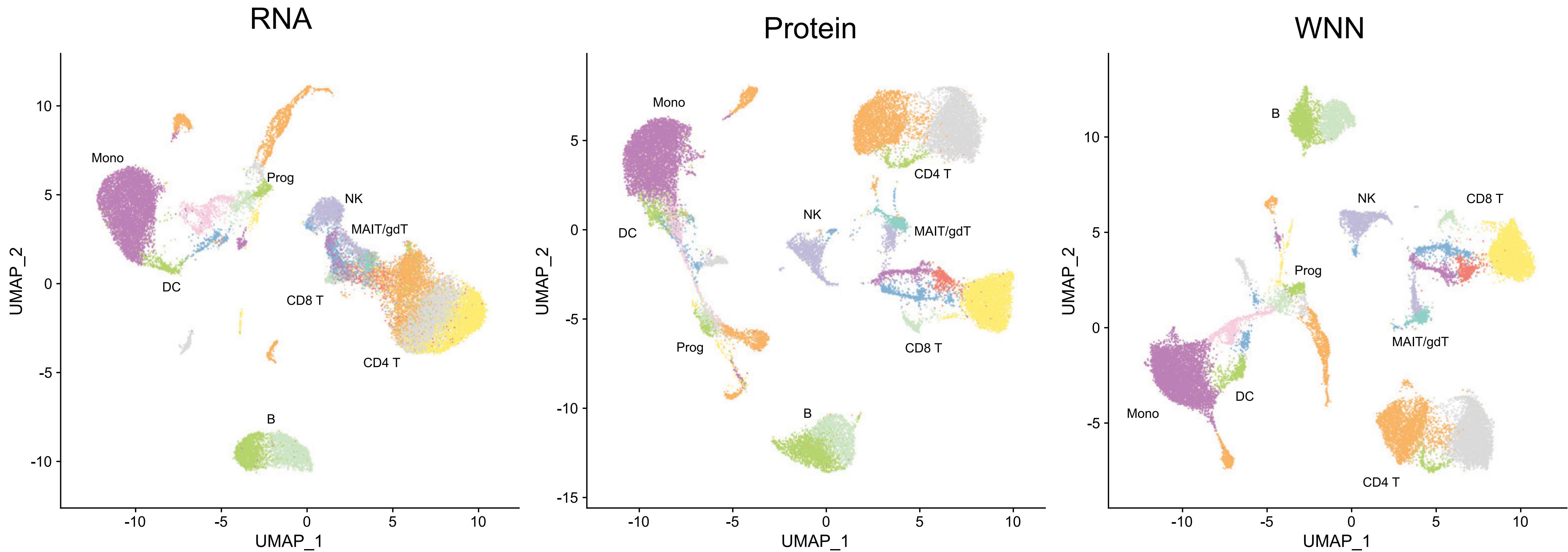
# Can we integrate multiple types of data using multiple types of cell-cell networks?



# Reconcile two weighted networks to build one consensus network (borrow info from each other)

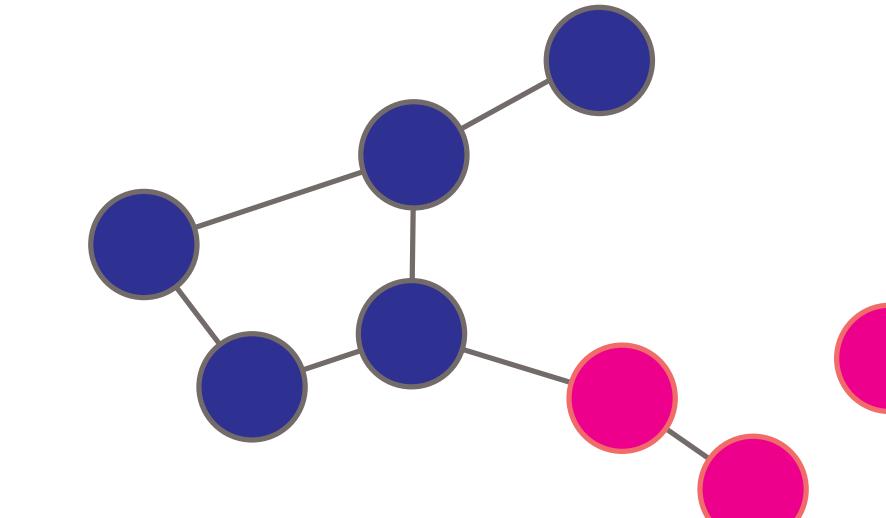


# Network-based integration works well if the goal of the analysis is to identify cell groups

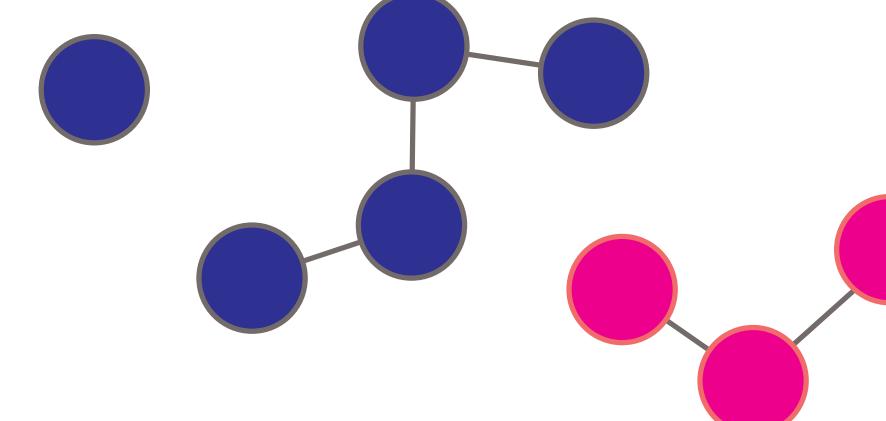


# What is the goal of this type of integration?

Cell-cell interaction  
network in data #1

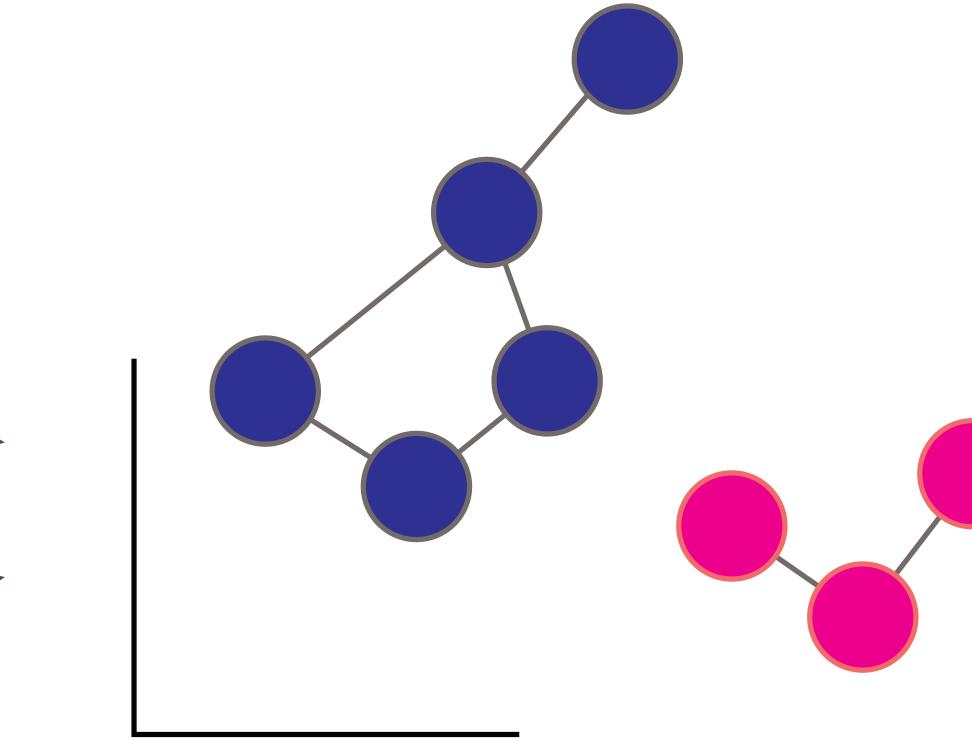


Cell-cell interaction  
network in data #2



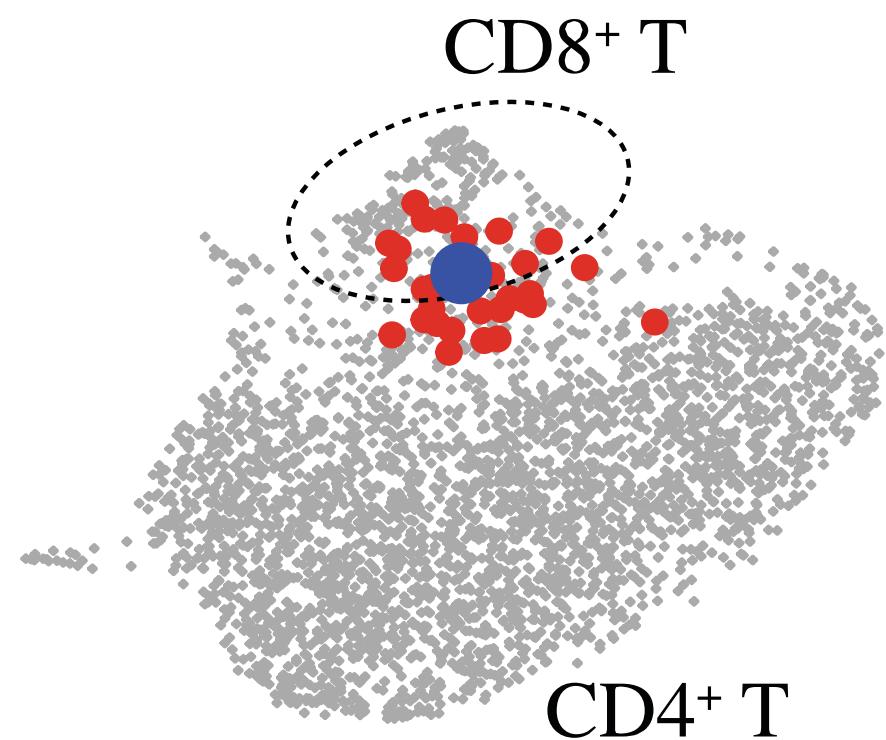
Goal:

1. Fill in missing info
2. Smooth out outliers



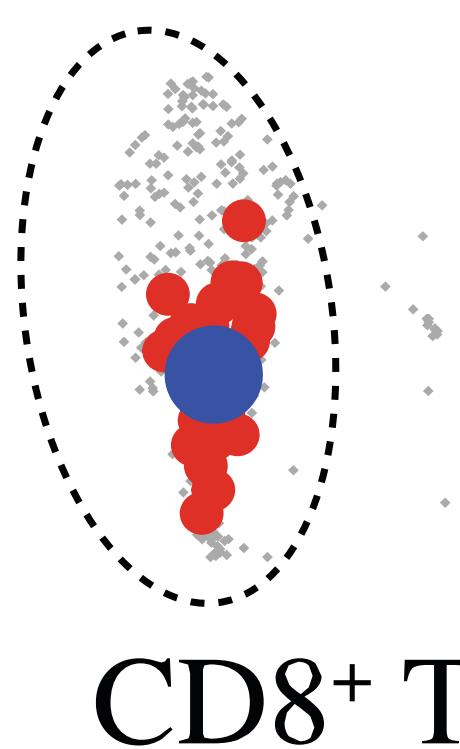
Consensus  
of two results

# K-nearest neighbour can predict gene expressions and protein activities (in the same domain)



● Target cell  
● KNN (RNA)

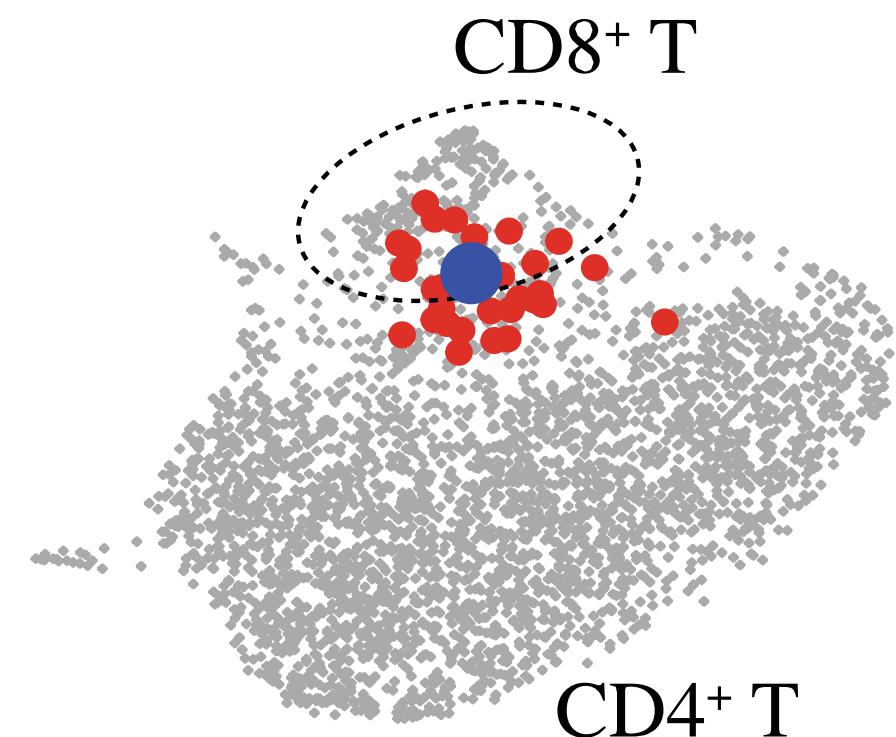
$$\hat{r}_{i,knn_r} = \frac{\sum_{j=1}^k r_{knn_{r,j}}}{k} : \text{prediction of RNA profile for cell } i, \text{ based on RNA neighbors}$$



● Target cell  
● KNN (protein)

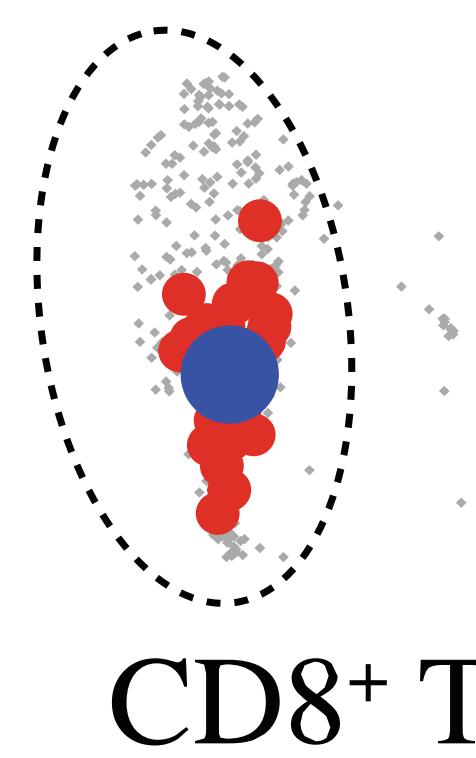
$$\hat{p}_{i,knn_p} = \frac{\sum_{j=1}^k p_{knn_{p,j}}}{k} : \text{prediction of protein profile for cell } i, \text{ based on protein neighbors}$$

# If the interaction was consistent, we would share the same neighbourhoods



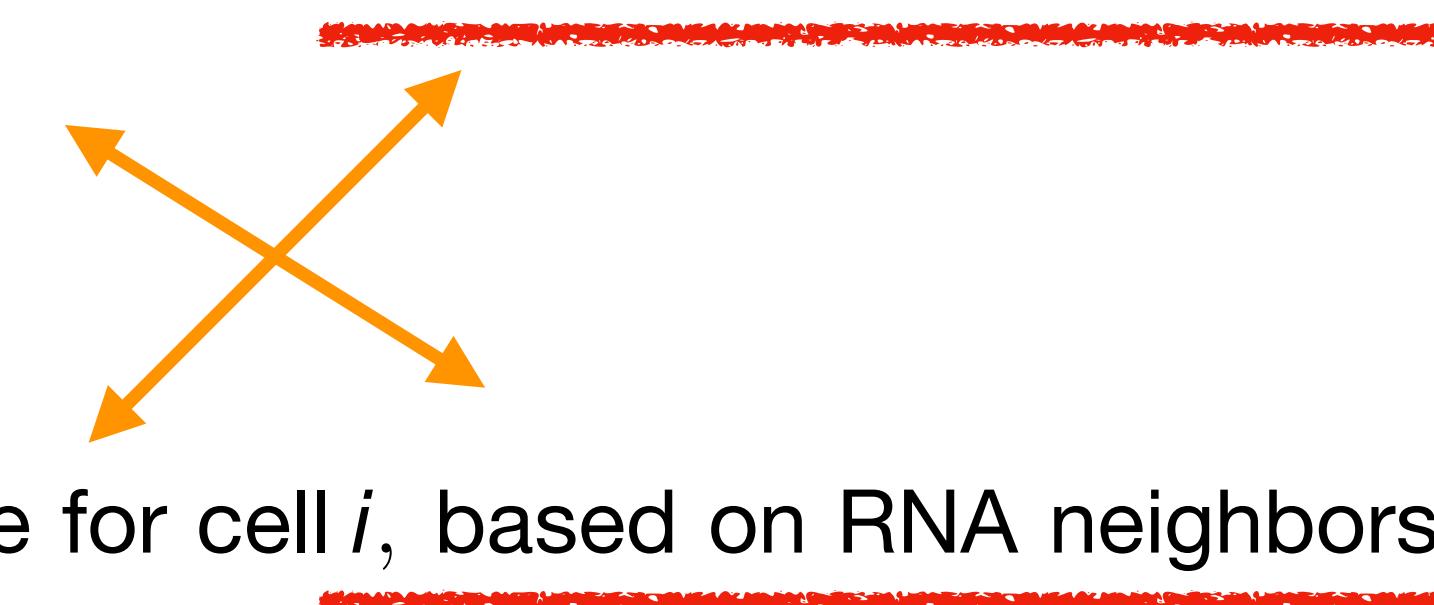
$$\hat{r}_{i,knn_p} = \frac{\sum_{j=1}^k r_{knn_{p,i,j}}}{k}$$

: prediction of RNA profile for cell  $i$ , based on protein neighbors

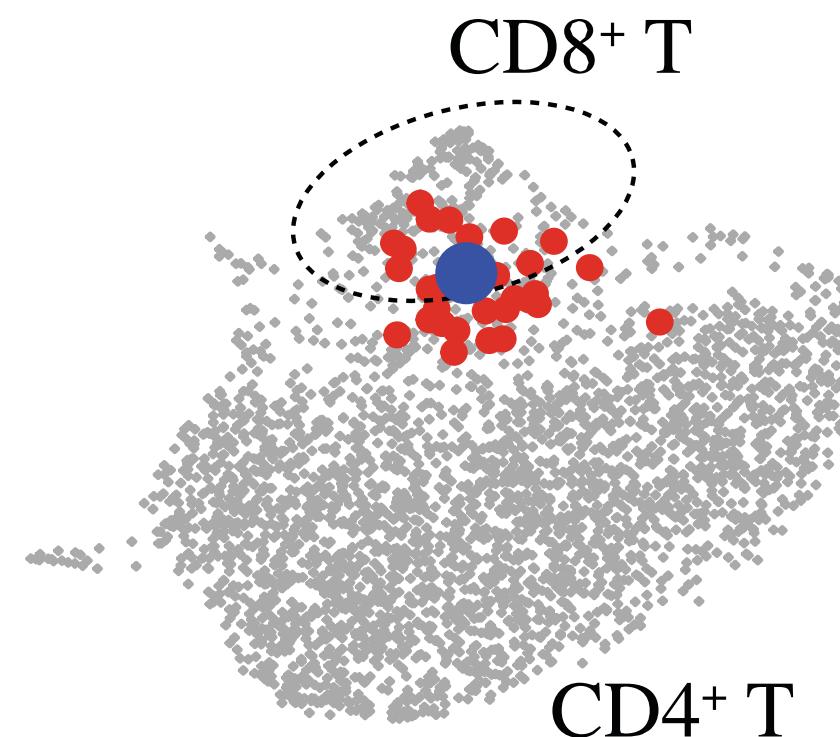


$$\hat{p}_{i,knn_r} = \frac{\sum_{j=1}^k p_{knn_{r,i,j}}}{k}$$

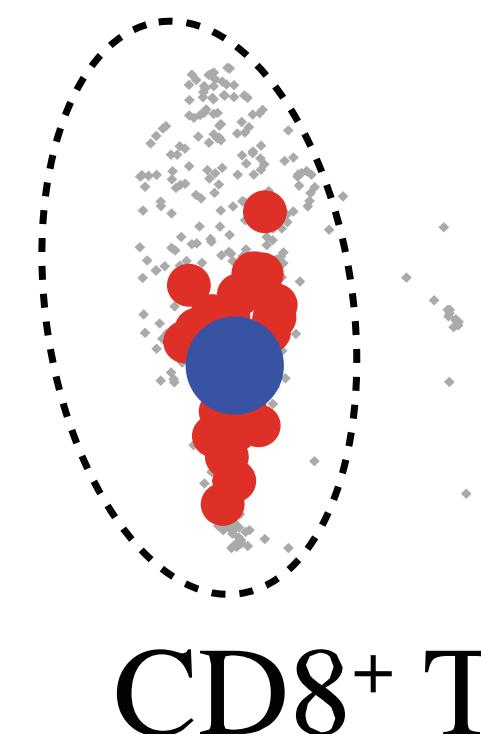
: prediction of protein profile for cell  $i$ , based on RNA neighbors



# How do you put them together? Take the weighted average over multiple types of edges



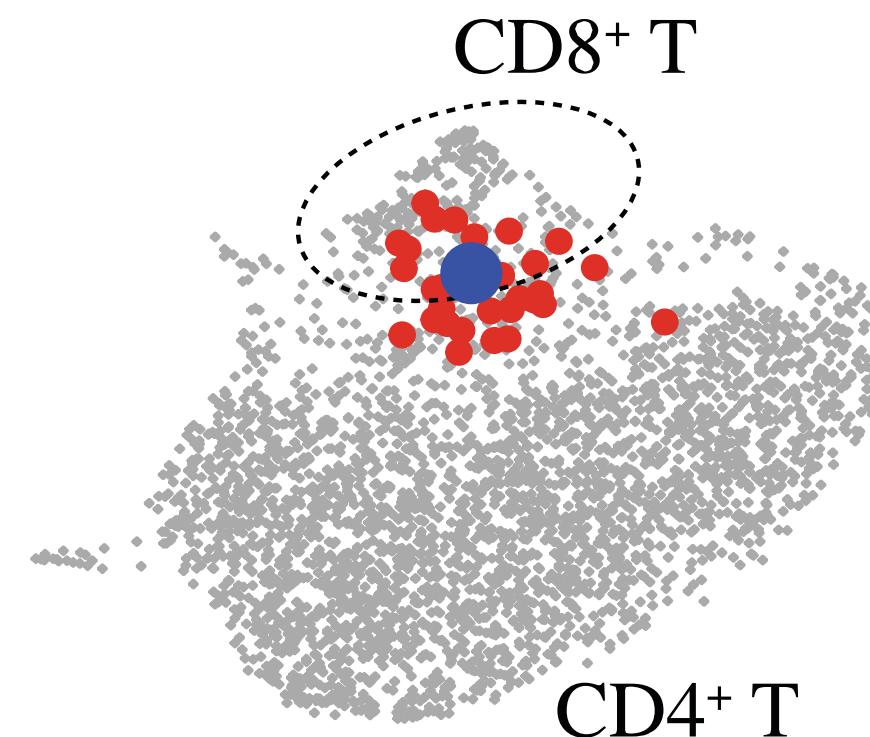
$$s_{rna}(i) = \frac{\theta_{rna} \left( r_i, \hat{r}_{i,knn_r} \right)}{\theta_{rna} \left( r_i, \hat{r}_{i,knn_p} \right) + \epsilon}, \quad s_{protein}(i) = \frac{\theta_{protein} \left( p_i, \hat{p}_{i,knn_p} \right)}{\theta_{protein} \left( p_i, \hat{p}_{i,knn_r} \right) + \epsilon}$$



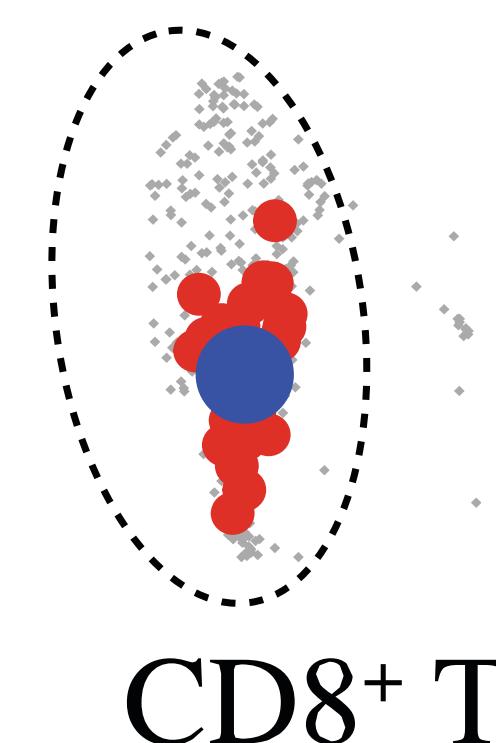
How much RNA weights  
can be explained by  
RNA or imputed by protein?

How much protein weights  
can be explained by  
protein or imputed by RNA?

# How do you put them together? Take the weighted average over multiple types of edges

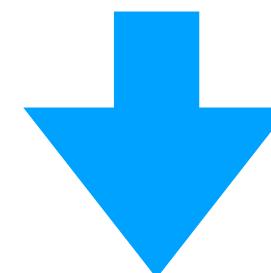


- Target cell
- KNN (RNA)

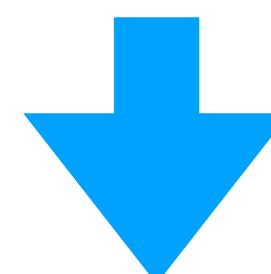


- Target cell
- KNN (protein)

$$s_{rna}(i) = \frac{\theta_{rna}(r_i, \hat{r}_{i,knn_r})}{\theta_{rna}(r_i, \hat{r}_{i,knn_p}) + \varepsilon}, \quad s_{protein}(i) = \frac{\theta_{protein}(p_i, \hat{p}_{i,knn_p})}{\theta_{protein}(p_i, \hat{p}_{i,knn_r}) + \varepsilon}$$



$$w_{rna}(i) = \frac{e^{s_{rna}(i)}}{e^{s_{rna}(i)} + e^{s_{protein}(i)}}, \quad w_{protein}(i) = \frac{e^{s_{protein}(i)}}{e^{s_{rna}(i)} + e^{s_{protein}(i)}}$$



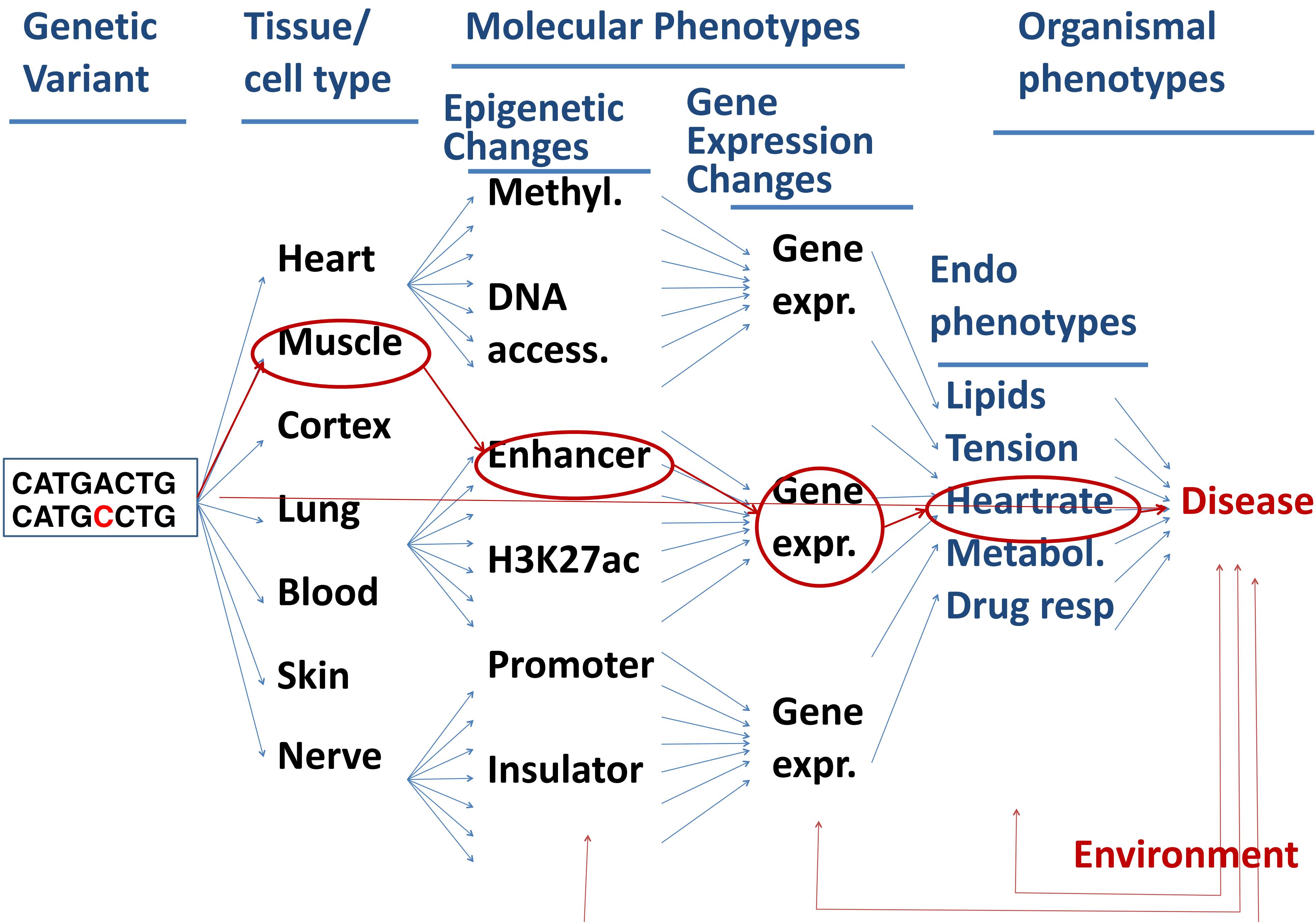
$$\theta_{weighted}(i, j) = w_{rna}(i)\theta_{rna}(r_i, r_j) + w_{protein}(i)\theta_{protein}(p_i, p_j)$$

# **Pros and cons for Weighted Network Approach**

# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms

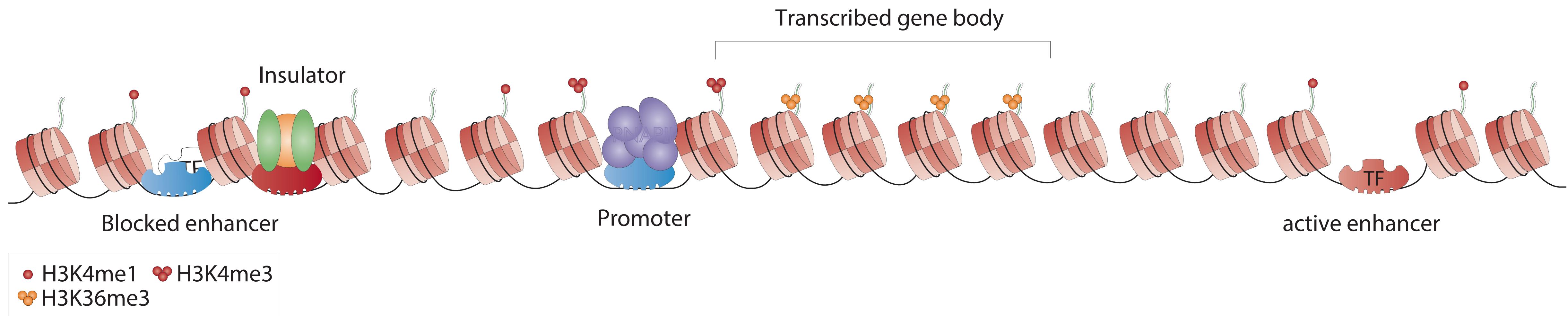
Where is  
**CCTCTGTG**  
**TCGGA**



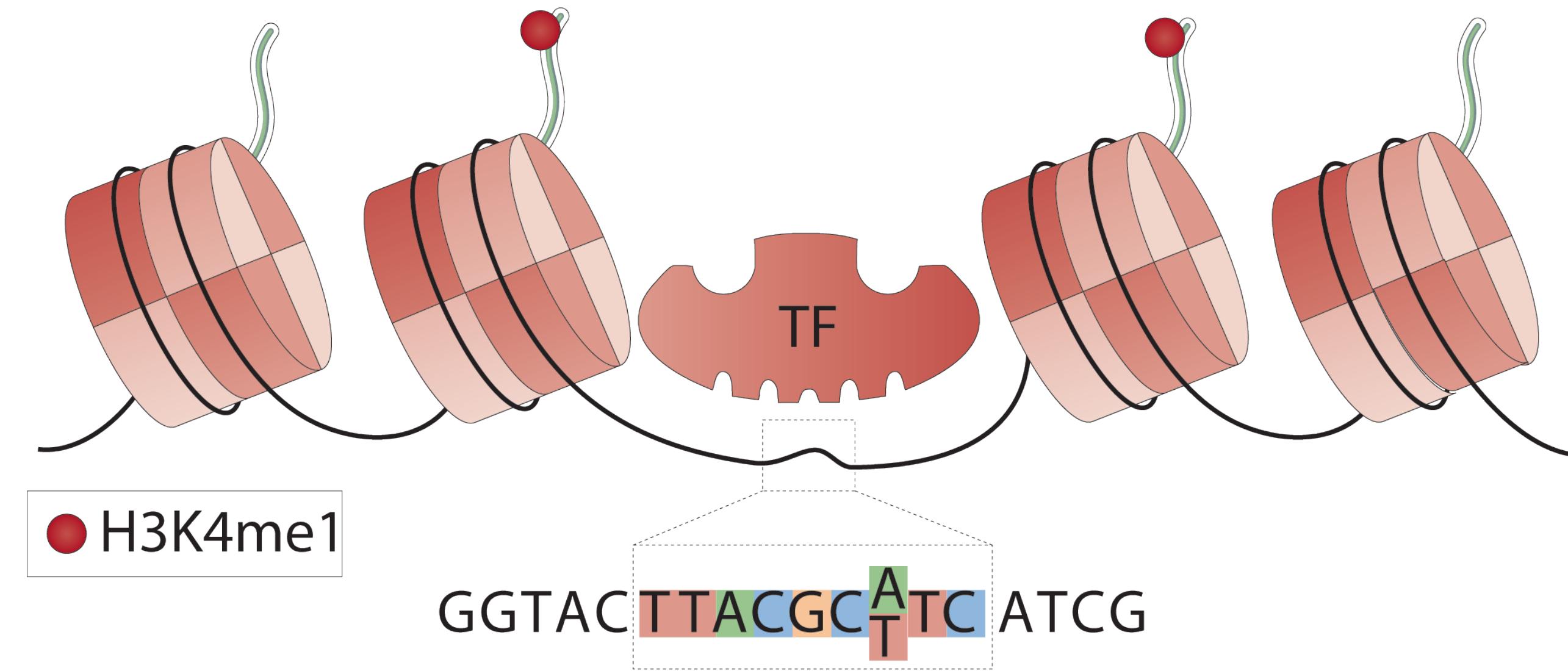
Slide credit: Manolis Kellis

**Feedback from environment / disease state**

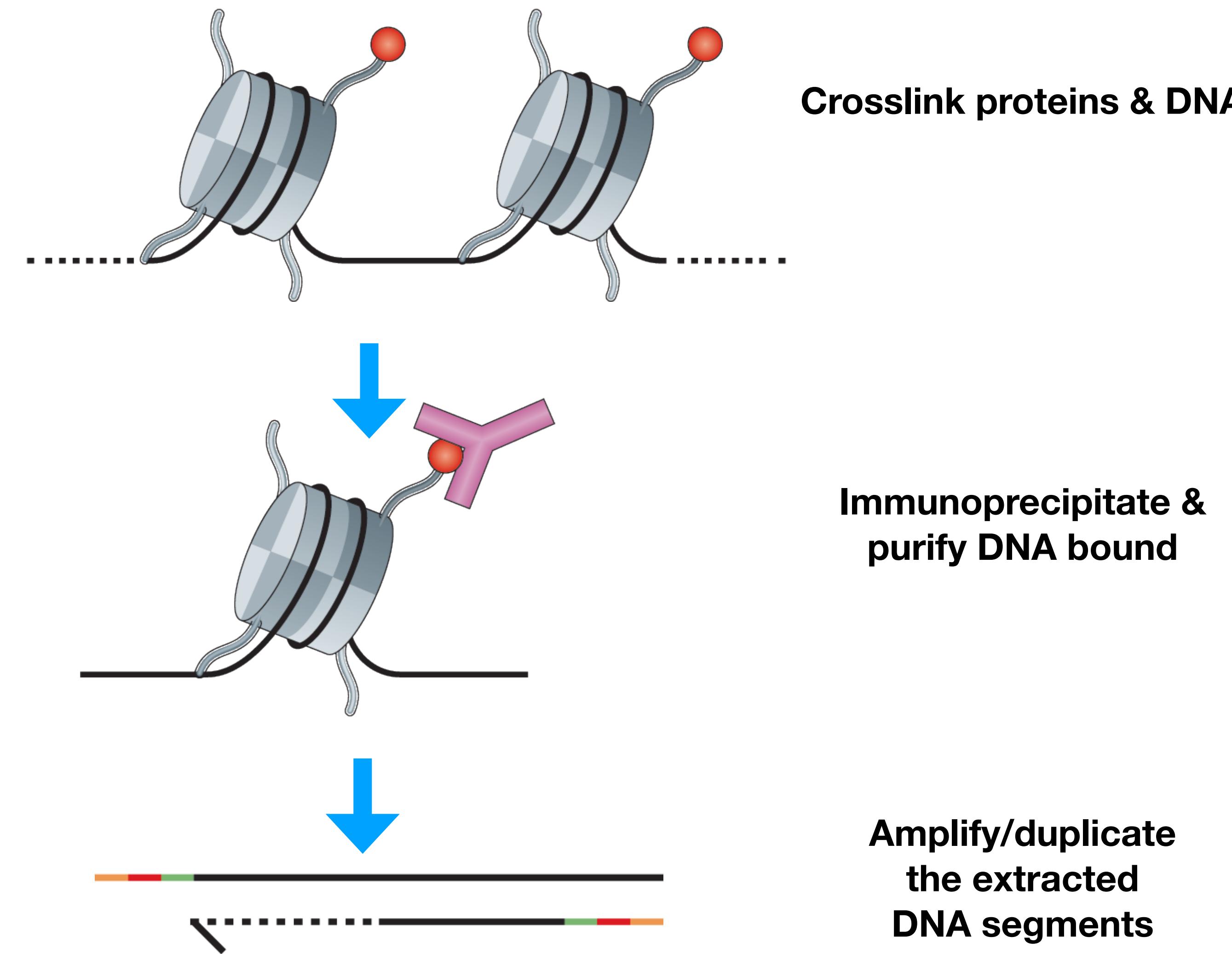
# Different histone code marks different types of regulatory elements



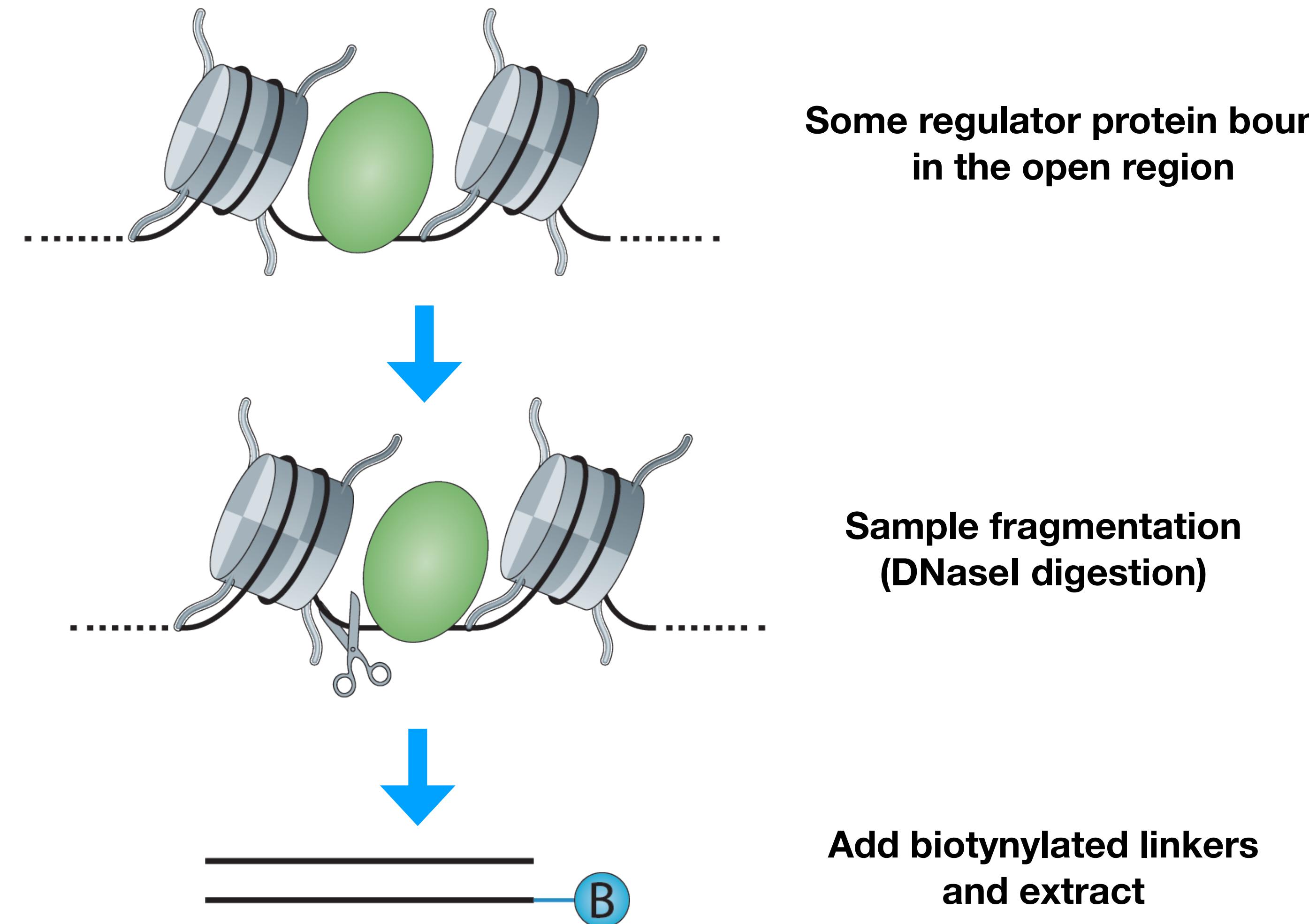
# One of the goals is to understand the logic of non-coding DNA sequences



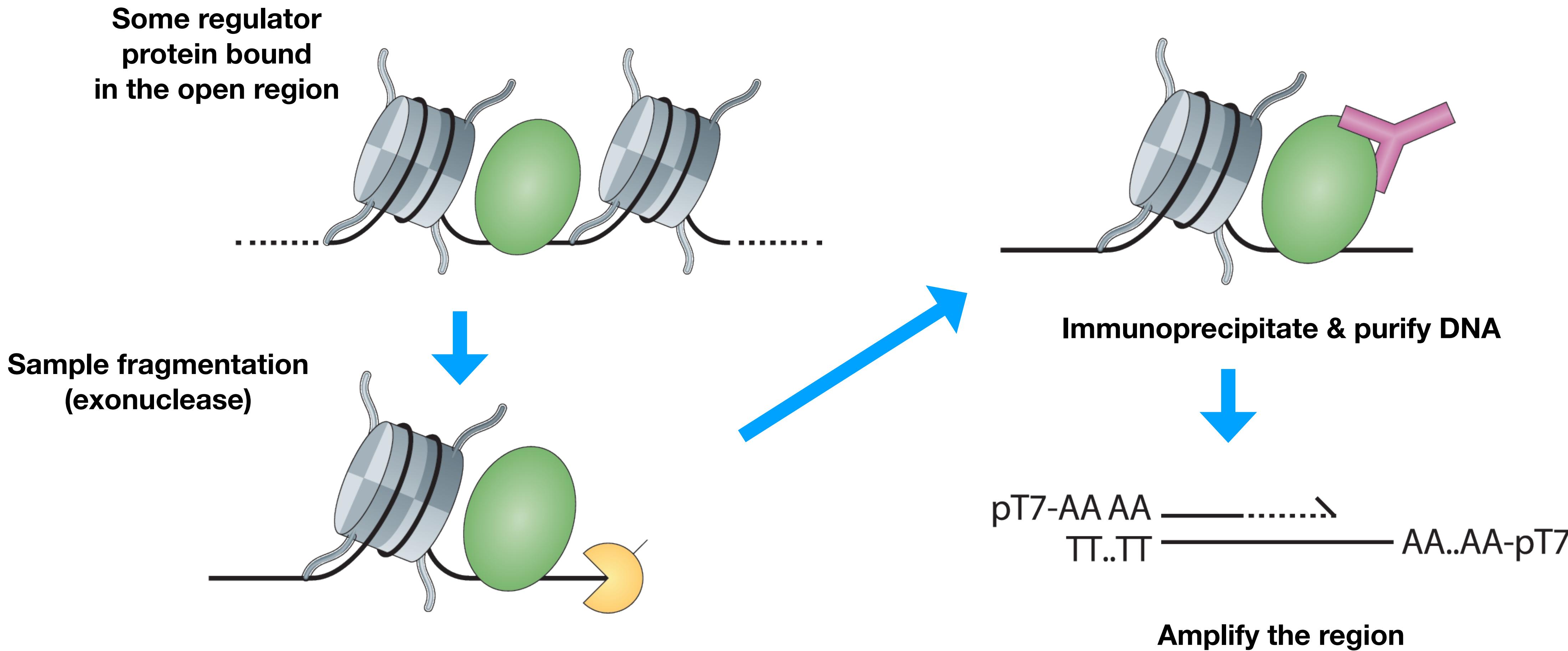
# ChIP-seq: Quantifying histone modifications



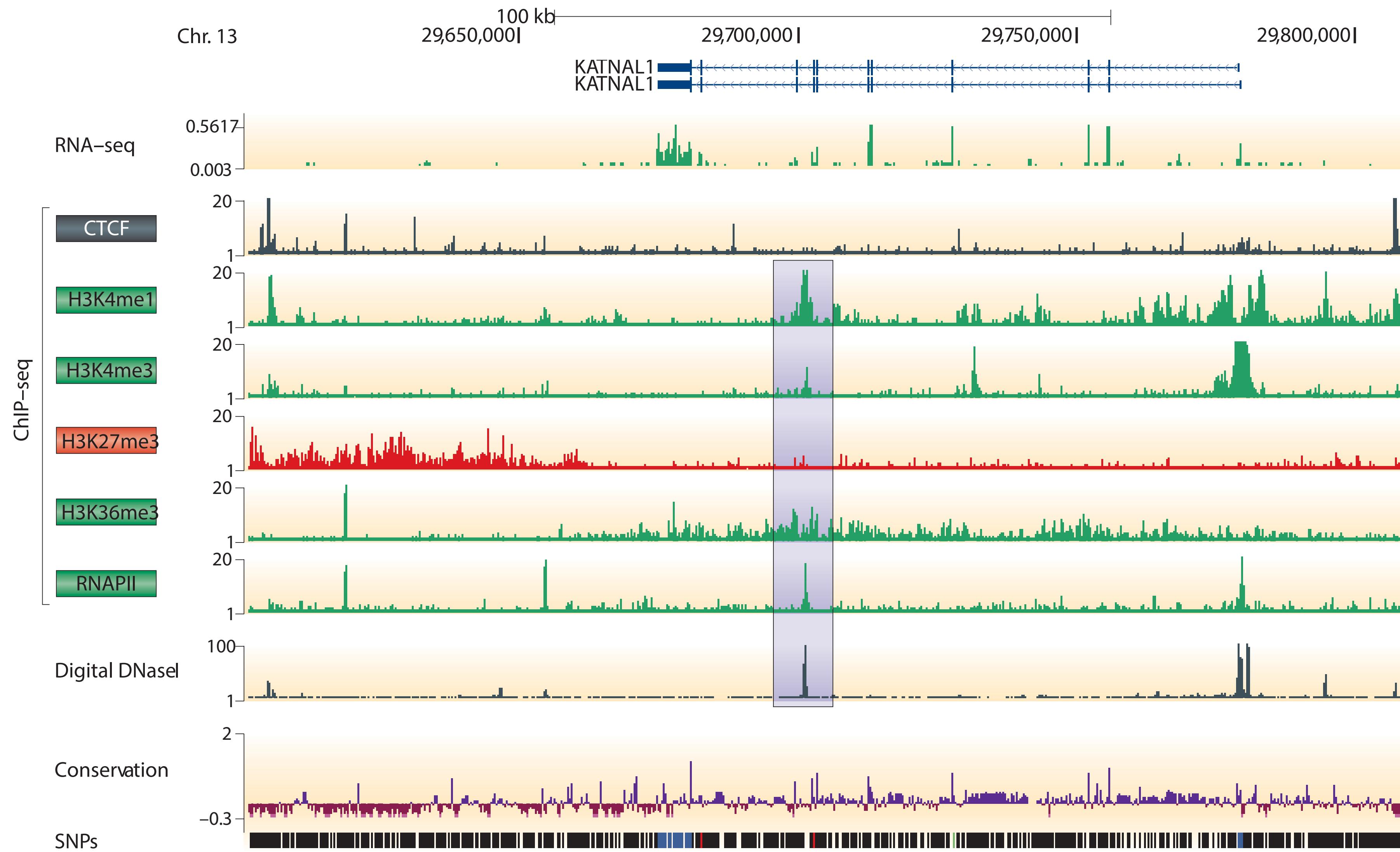
# DNase-seq: Quantifying DNA accessibility



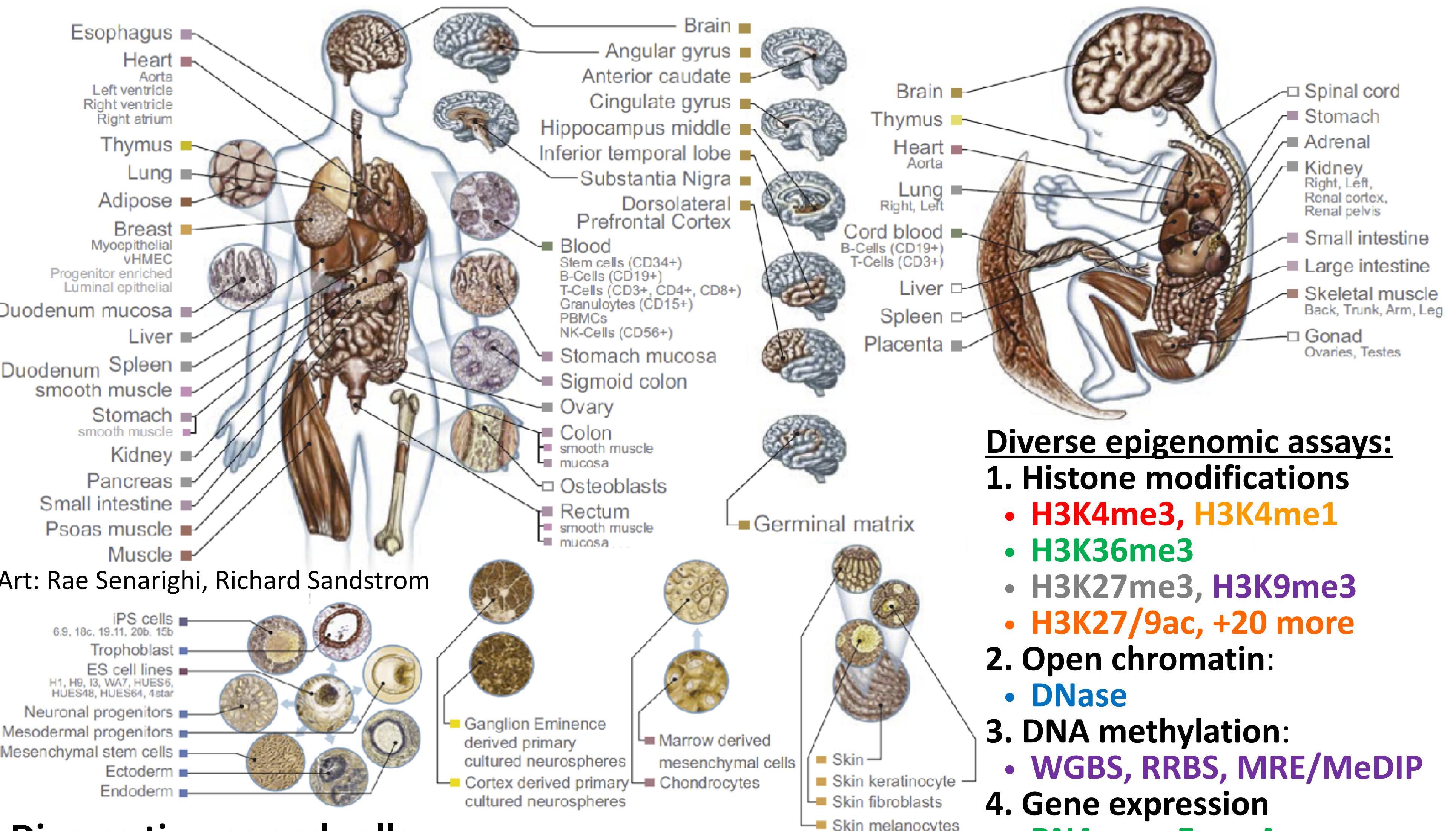
# ChIP-seq: Quantifying regions bound by a particular protein



# An example of multiple ChIP-seq tracks



# Epigenomics Roadmap across 100+ tissues/cell types



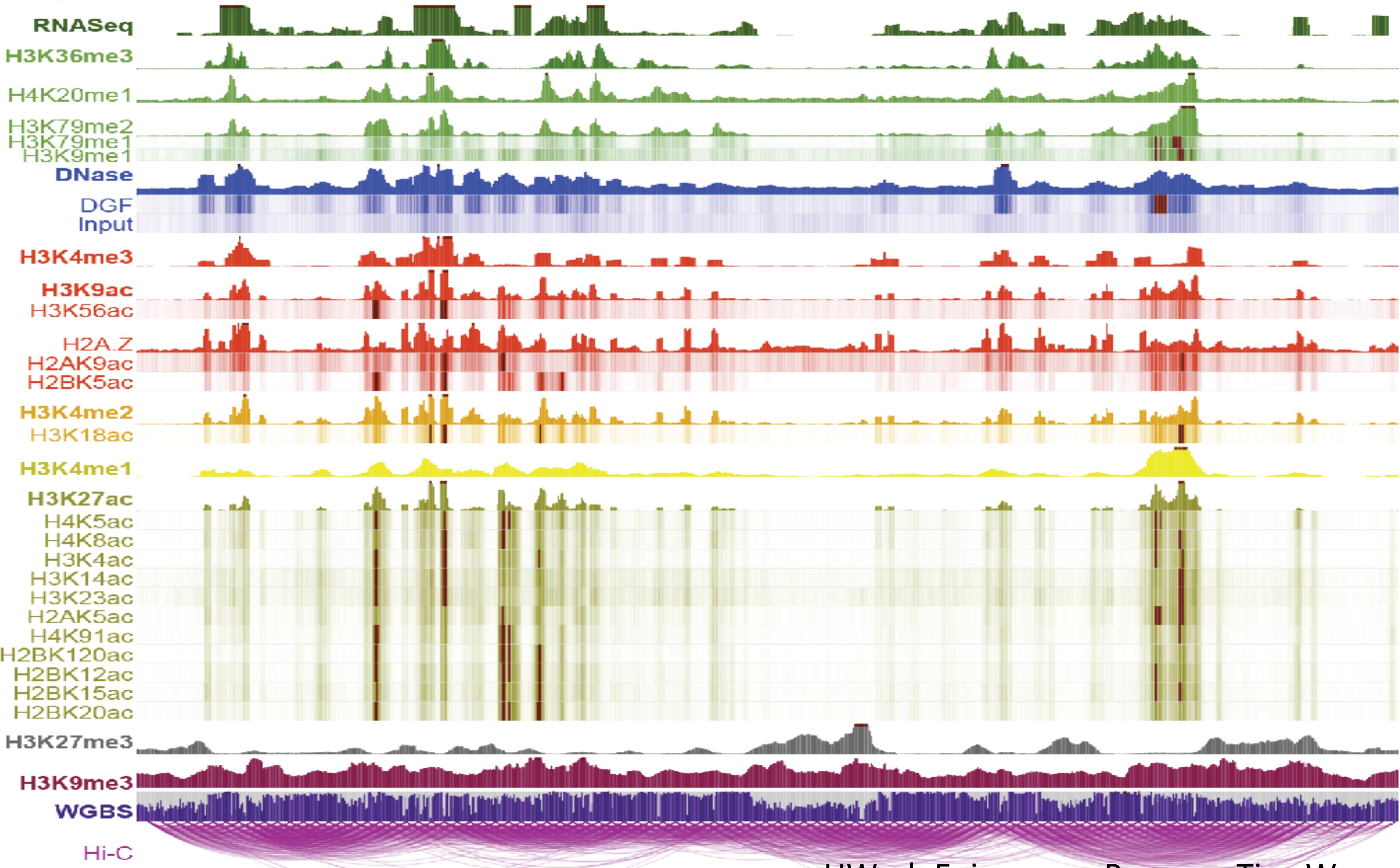
## Diverse tissues and cells:

1. Adult tissues and cells (brain, muscle, heart, digestive, skin, adipose, lung, blood...)
2. Fetal tissues (brain, skeletal muscle, heart, digestive, lung, cord blood...)
3. ES cells, iPS, differentiated cells (meso/endo/ectoderm, neural, mesench, trophobl)

## Diverse epigenomic assays:

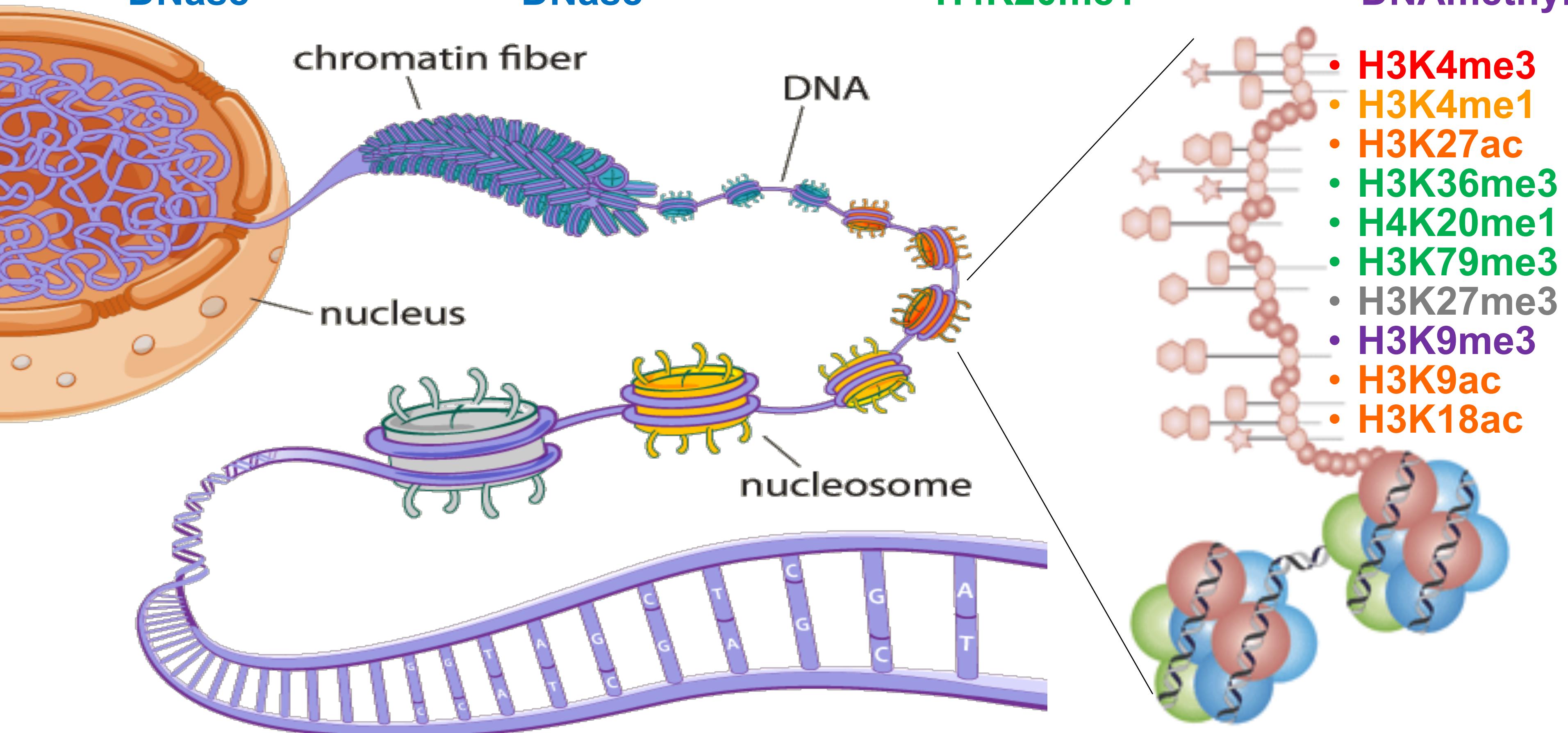
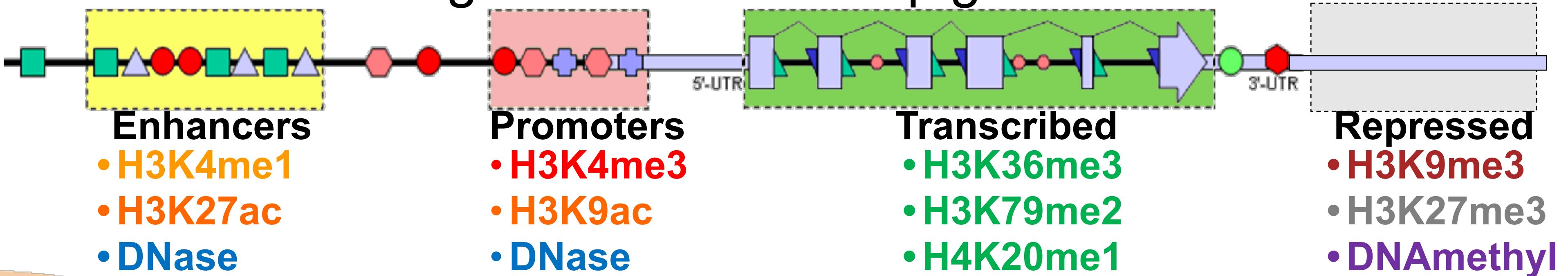
1. Histone modifications
  - H3K4me3, H3K4me1
  - H3K36me3
  - H3K27me3, H3K9me3
  - H3K27/9ac, +20 more
2. Open chromatin:
  - DNase
3. DNA methylation:
  - WGBS, RRBS, MRE/MeDIP
4. Gene expression
  - RNA-seq, Exon Arrays

# Deep sampling of 9 reference epigenomes (e.g. IMR90)



UWash Epigenome Browser, Ting Wang

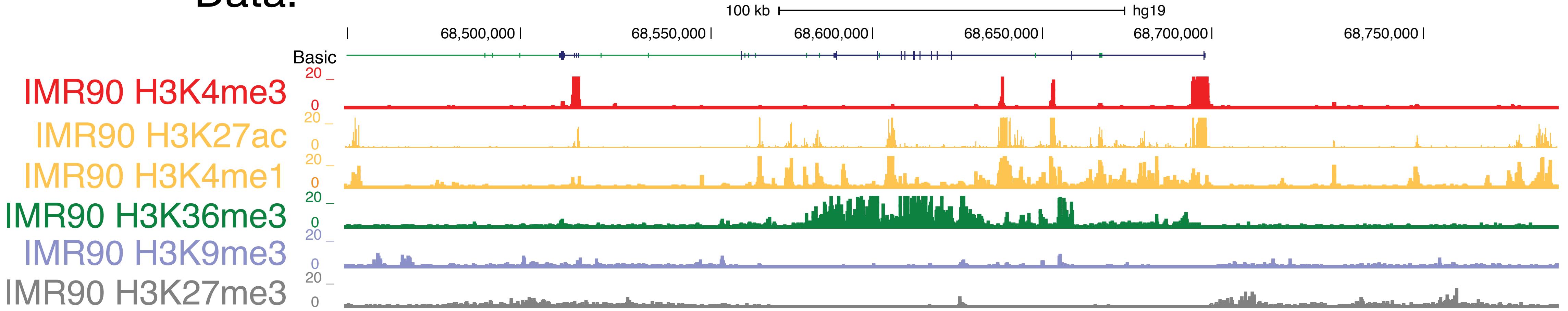
# Diverse chromatin signatures encode epigenomic state



- 100s of known modifications, many new still emerging
- Systematic mapping using ChIP-, Bisulfite-, DNase-Seq

# Goal: aggregate multiple ChIP-seq tracks

Data:



Unknown hidden states (epigenetic status; colour)

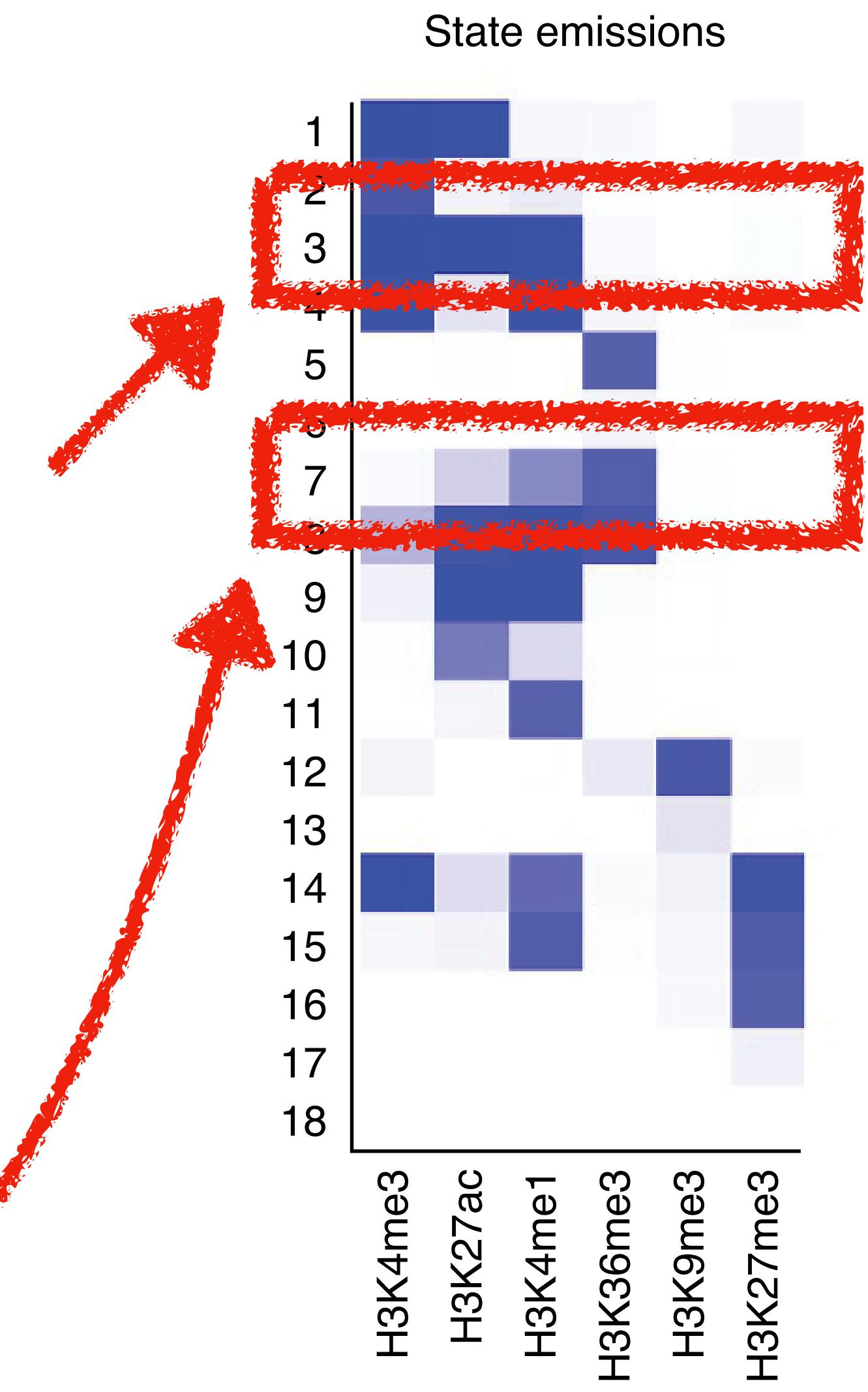
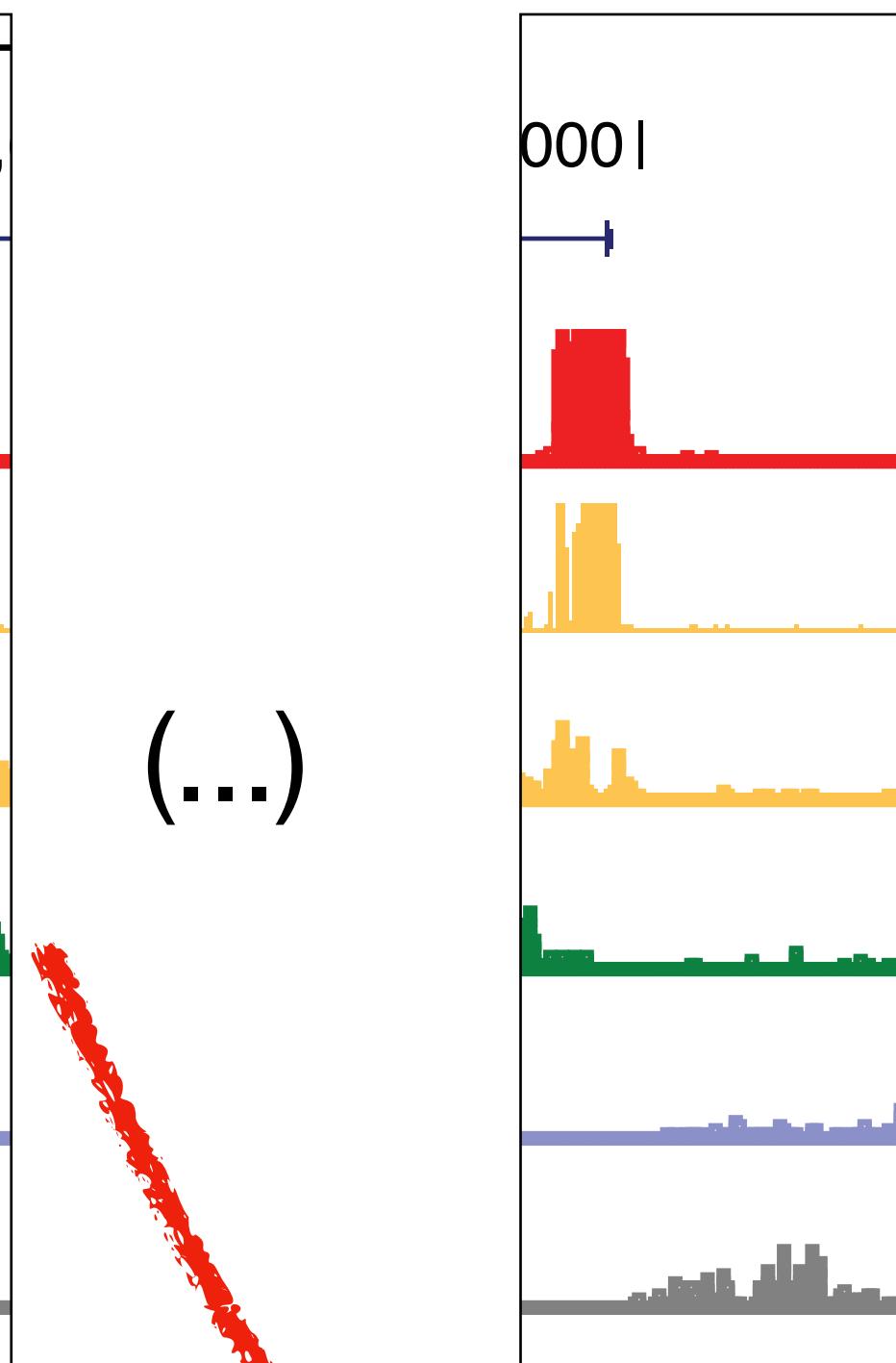
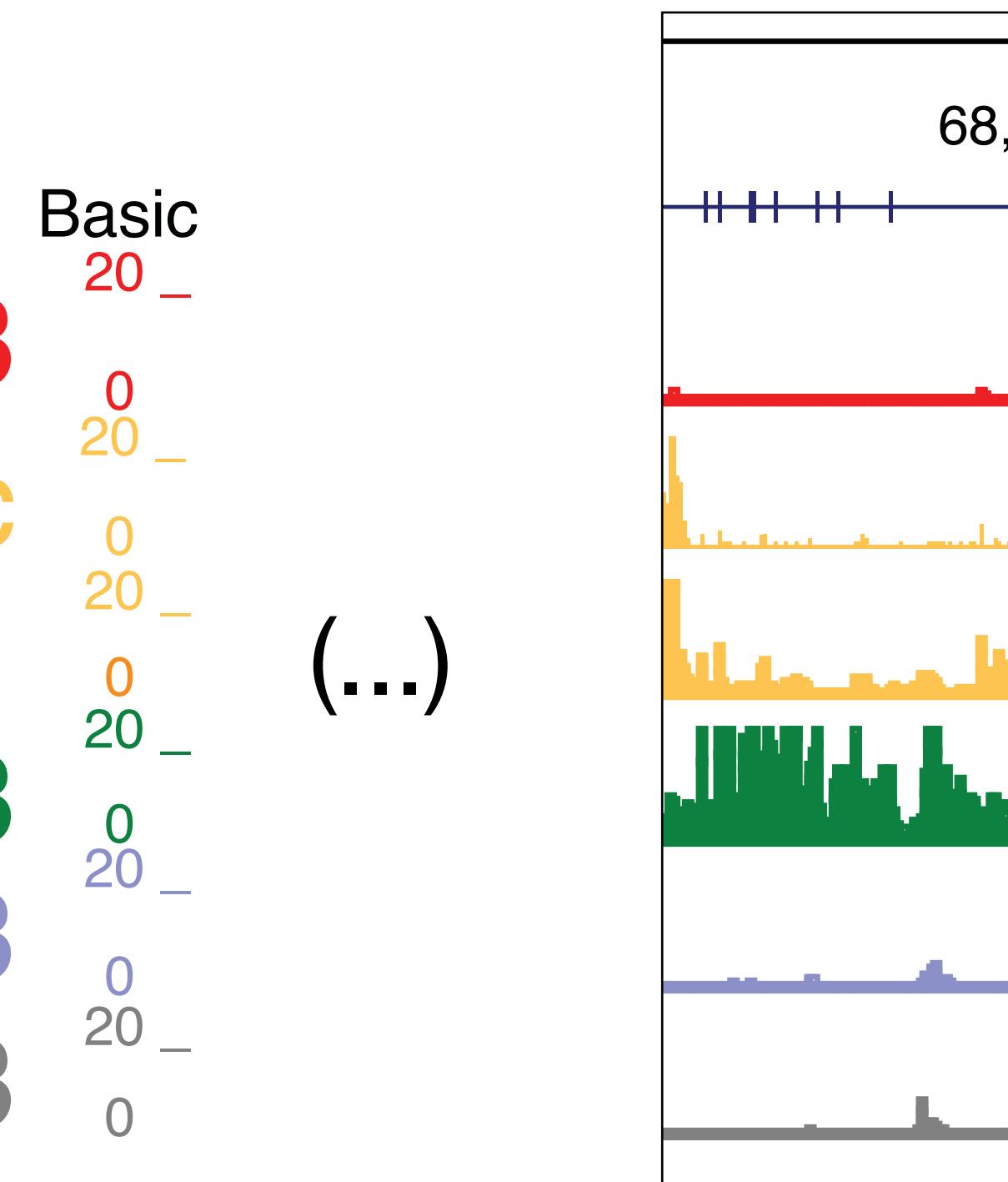


# ChromHMM (hidden Markov model)

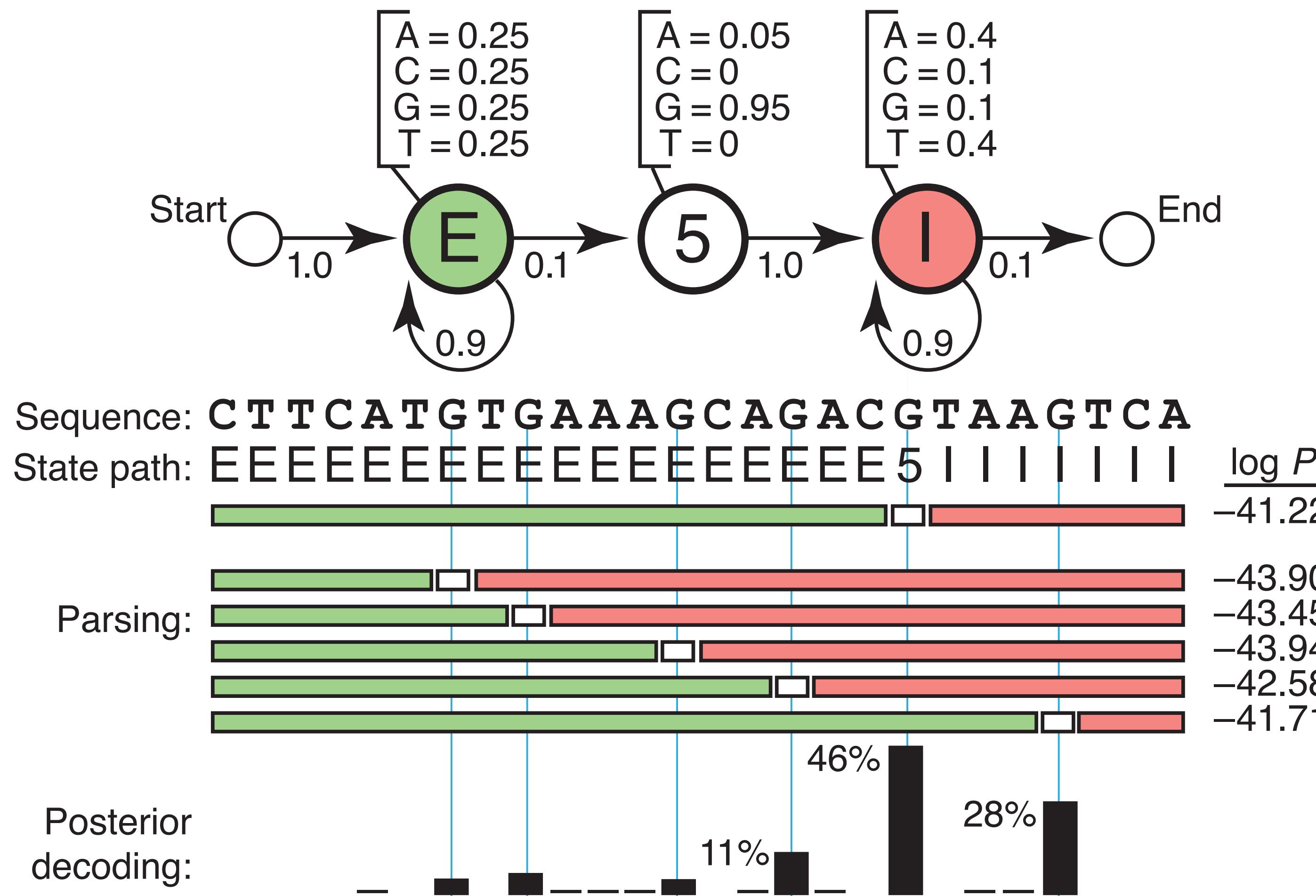
# Data:

IMR90 H3K4me3  
IMR90 H3K27ac  
IMR90 H3K4me1  
IMR90 H3K36me3  
IMR90 H3K9me3  
IMR90 H3K27me3

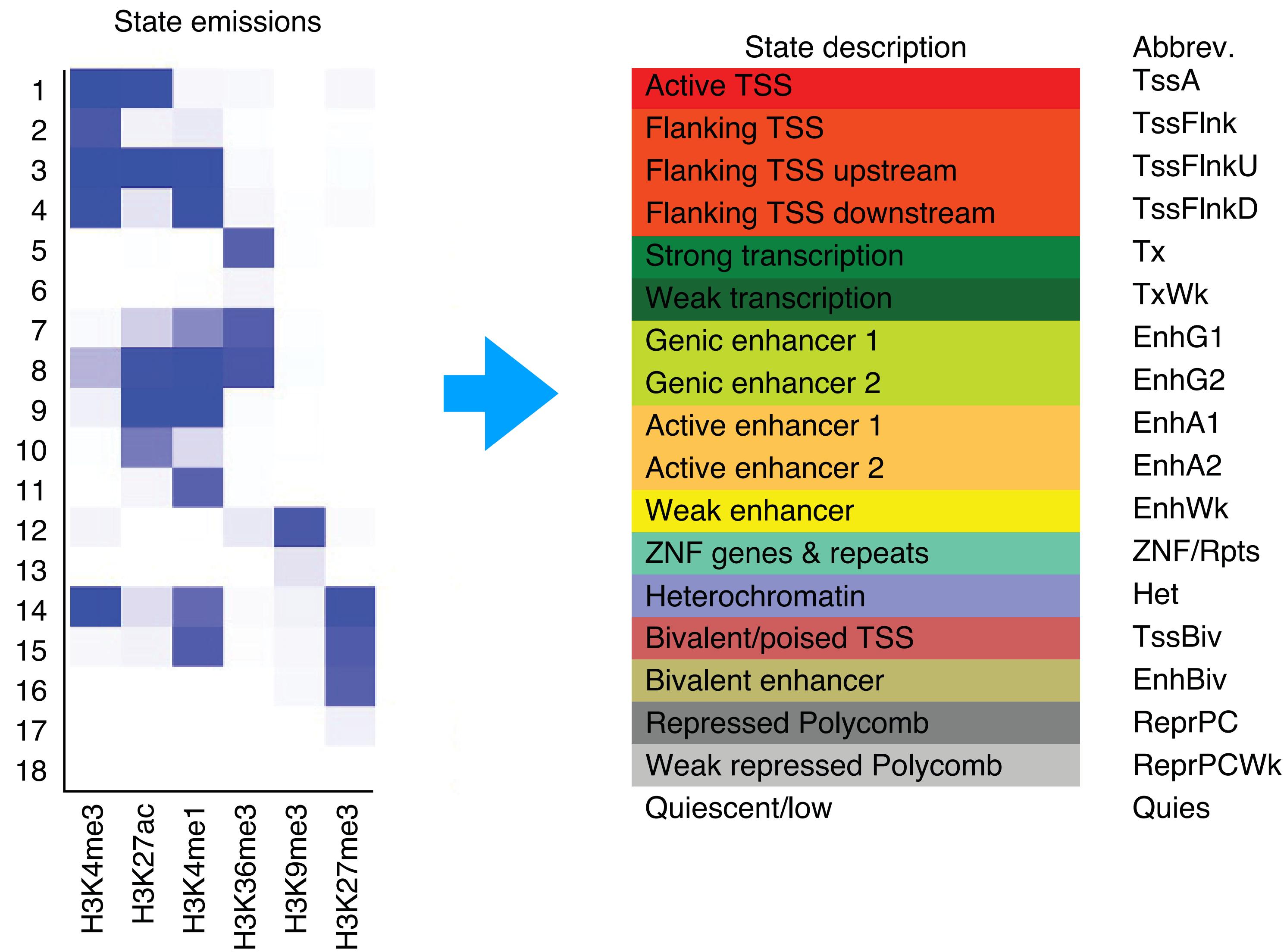
# Systematic annotation:



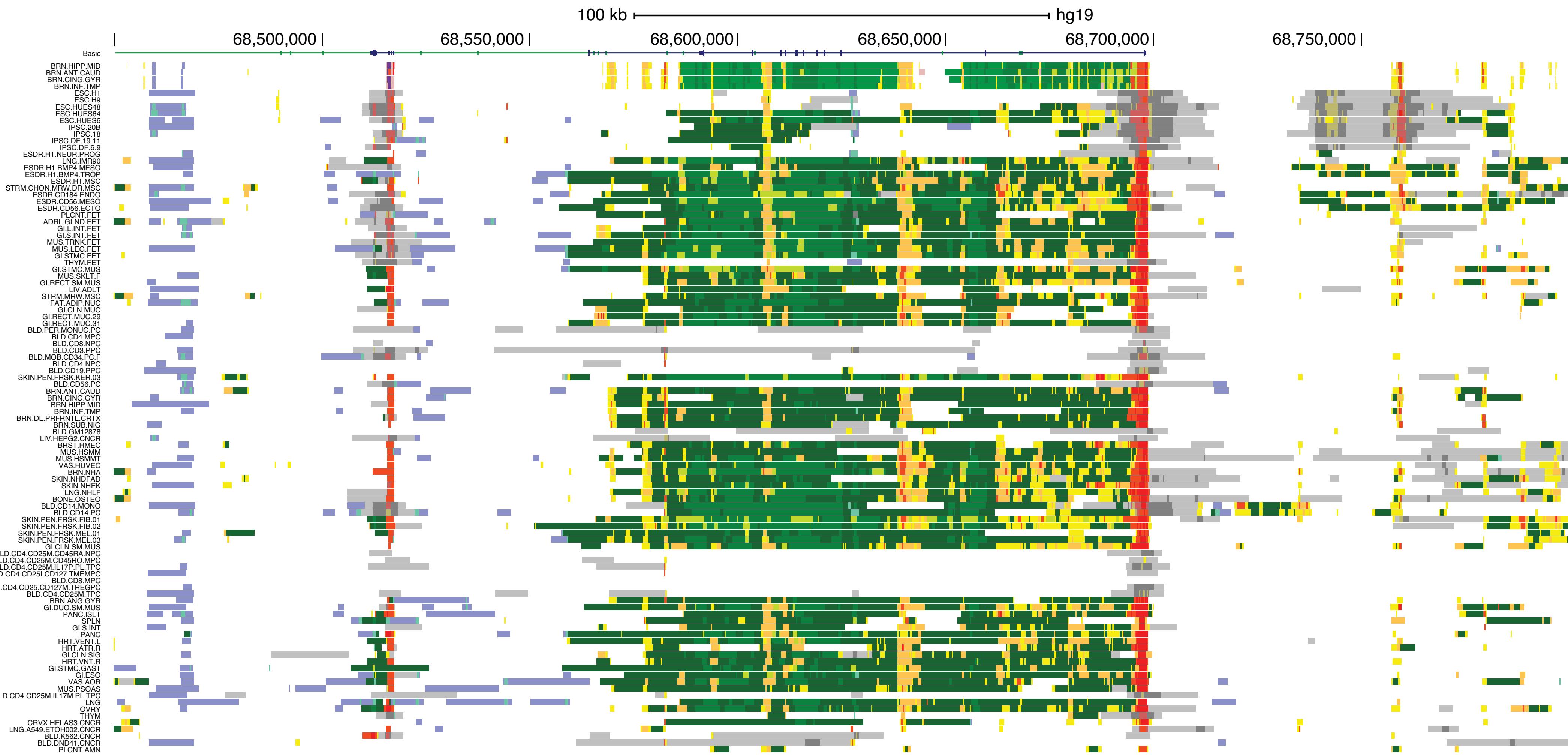
# HMM has been used in multiple bioinformatics (especially DNA sequence analysis)



# ChromHMM (hidden Markov model)



# ChromHMM: multiple ChIP-seq signal tracks → systematic epigenome annotation

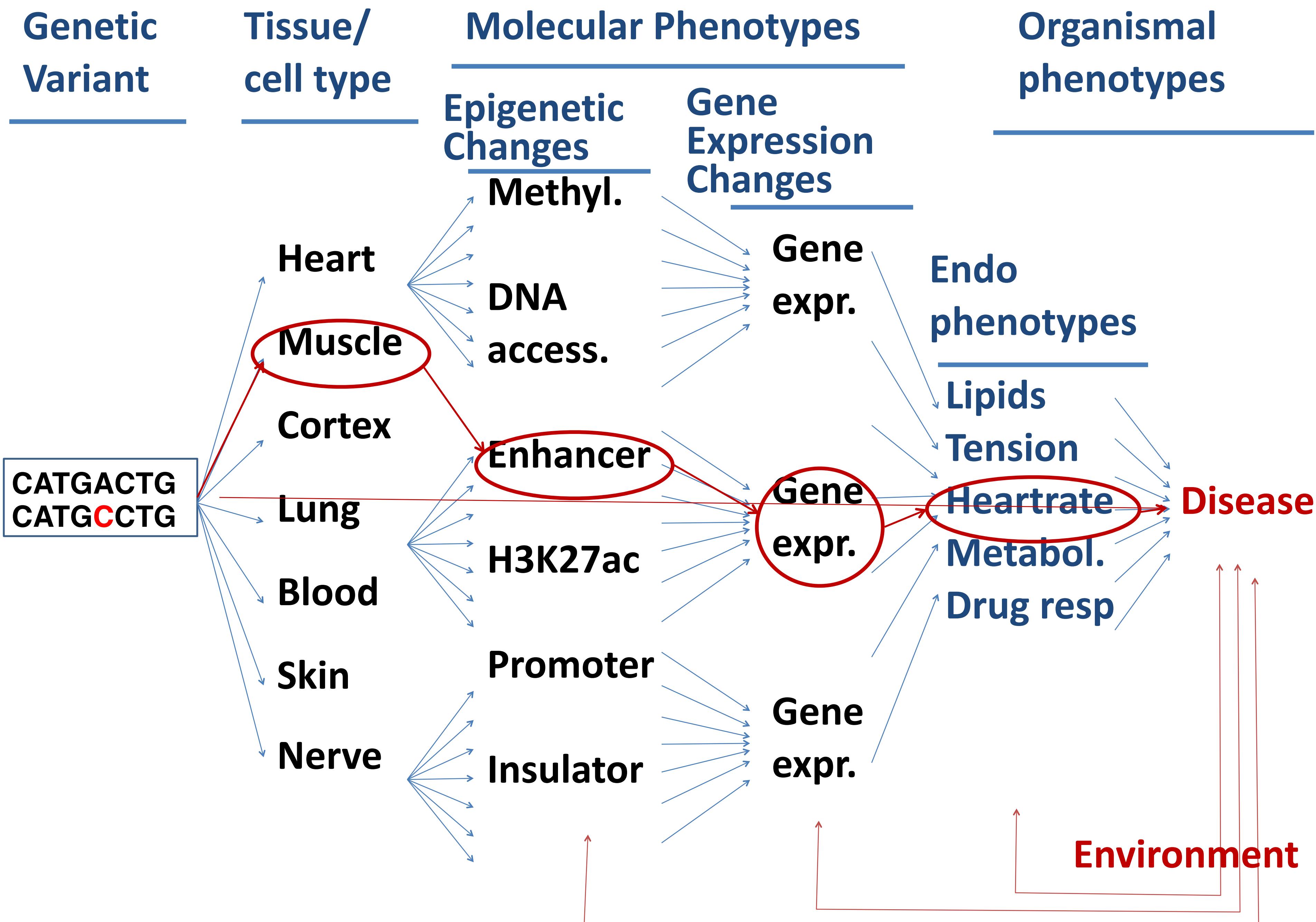


# Discussion: ChromHMM

Do we have other alternatives?

# Today's lecture: Multiomics data integration

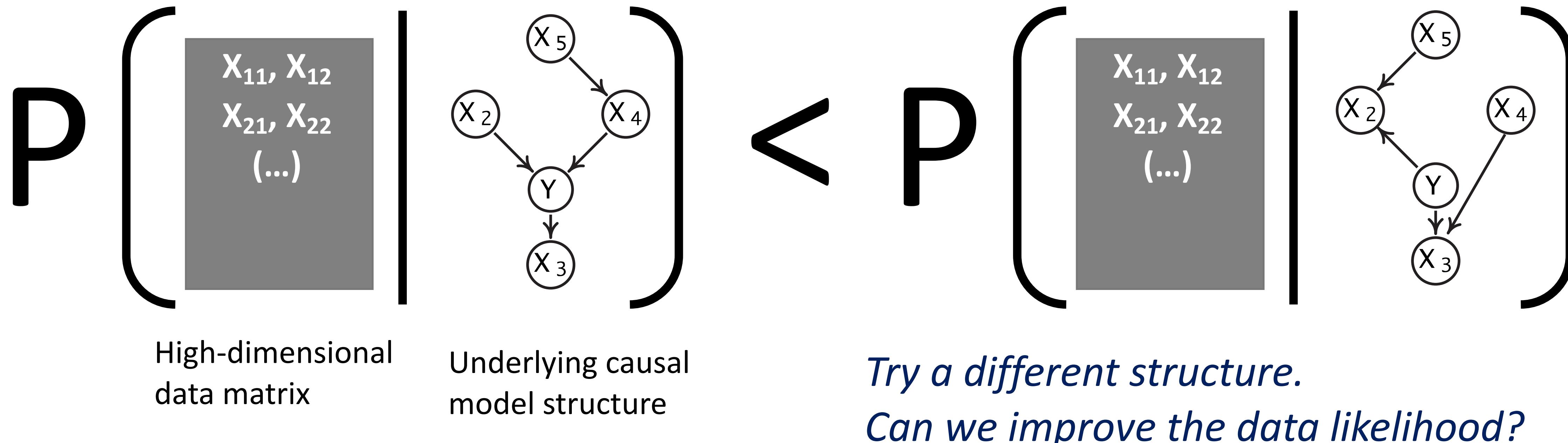
- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms



Slide credit: Manolis Kellis

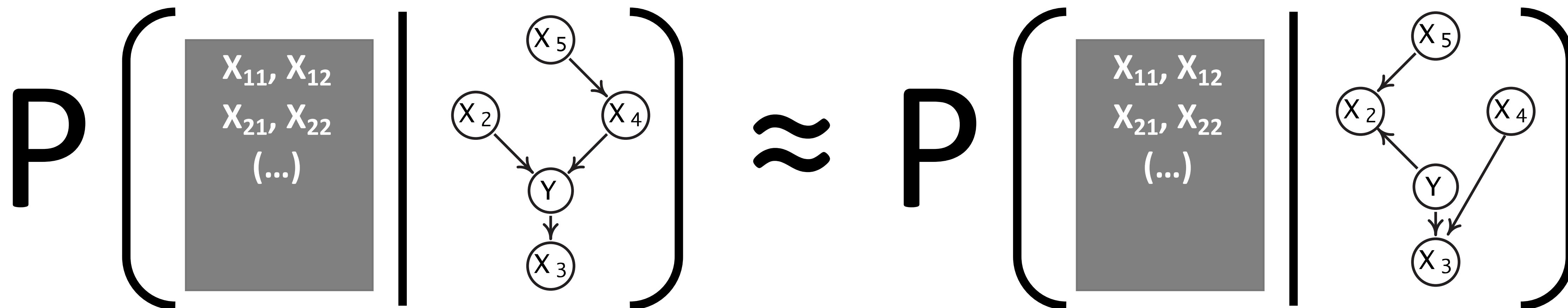
**Feedback from environment / disease state**

# Isn't it about finding the most probable structure?



“Large-sample learning of Bayesian Network is NP-hard”,  
Chickering .. Heckerman, UAI (2002)

# A traditional structural learning may not achieve identifiability ... (Even if we increase N)

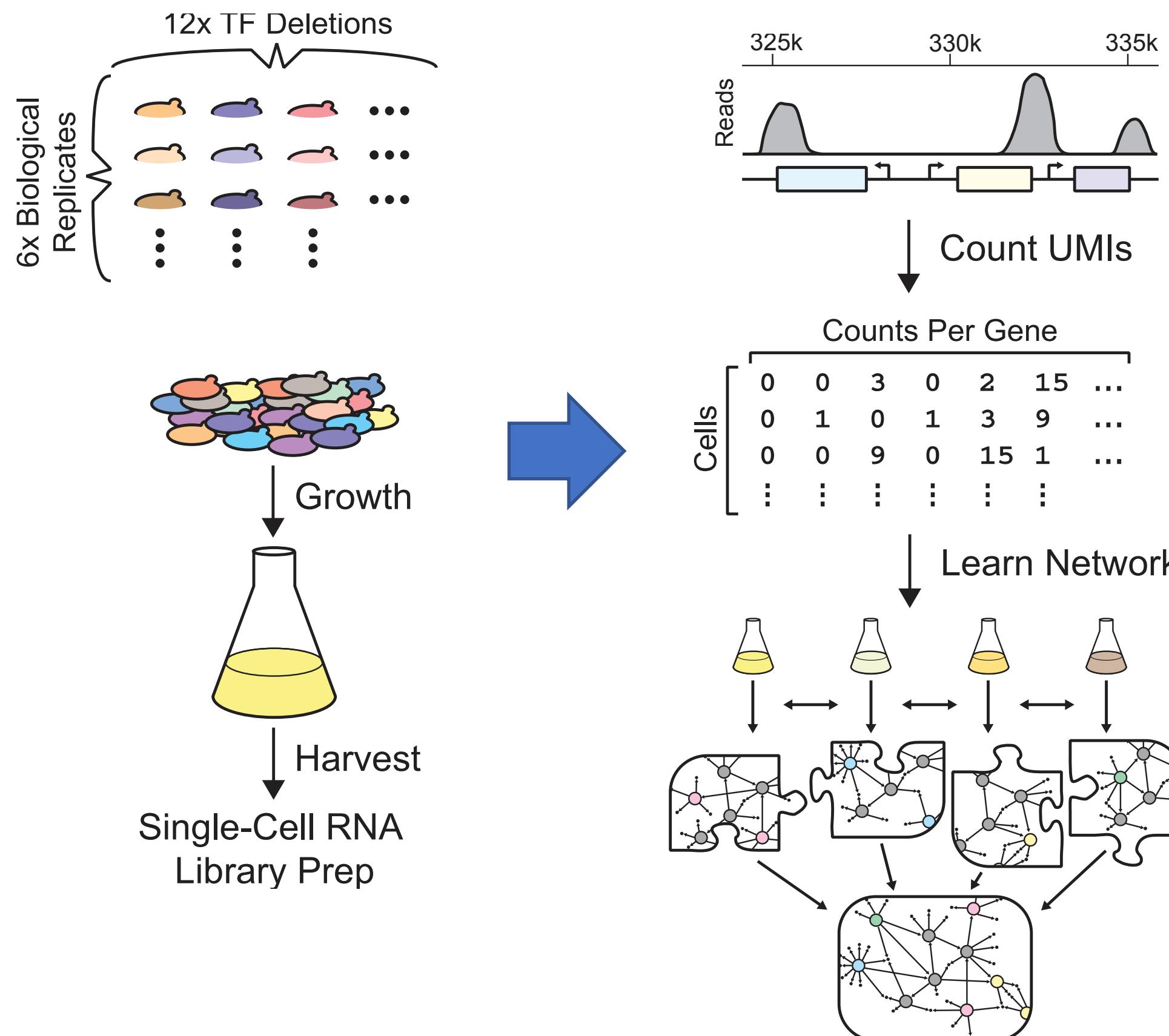


High-dimensional  
data matrices

Underlying causal  
model structure

*What if our probability score is  
more or less the same for  
different structures?*

# Can we borrow the awesome power of genetics?



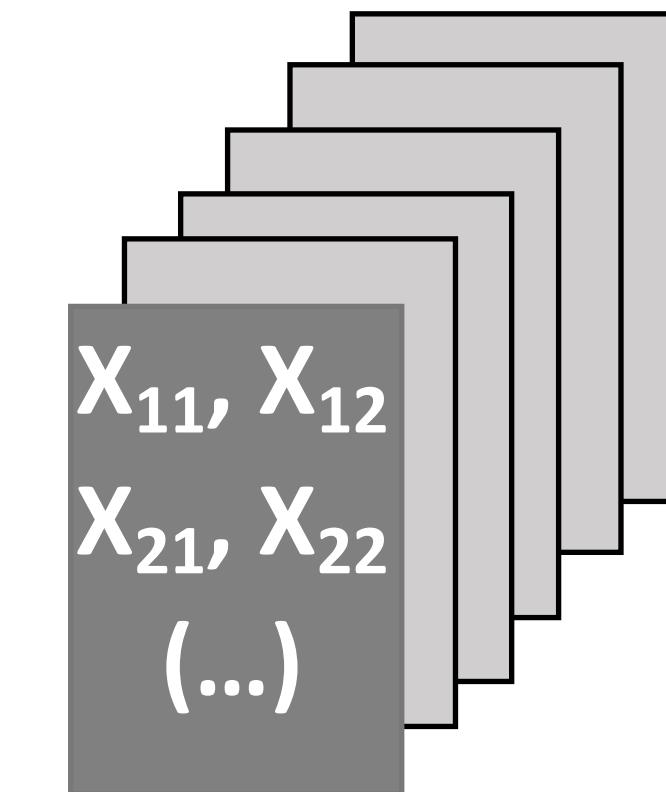
**Early 2000's**

$$\begin{matrix} X_{11}, X_{12}, X_{13} \\ X_{21}, X_{22}, X_{23} \\ \dots \end{matrix}$$

A single or only a handful of conditions

*Many computational biologists gave up on this NP-hard problem... learning BN from data*

**After perturb-seq**

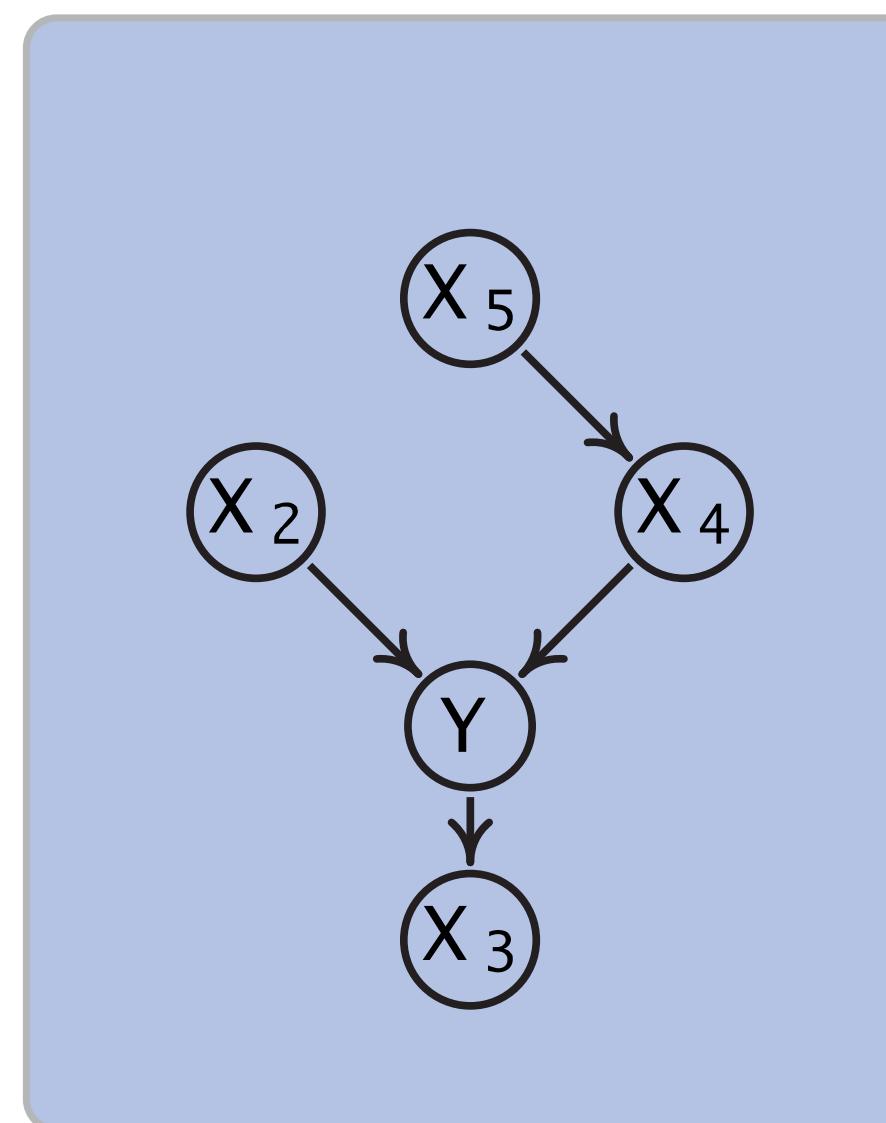


Massive (random) perturbation assays → hundreds of data matrices

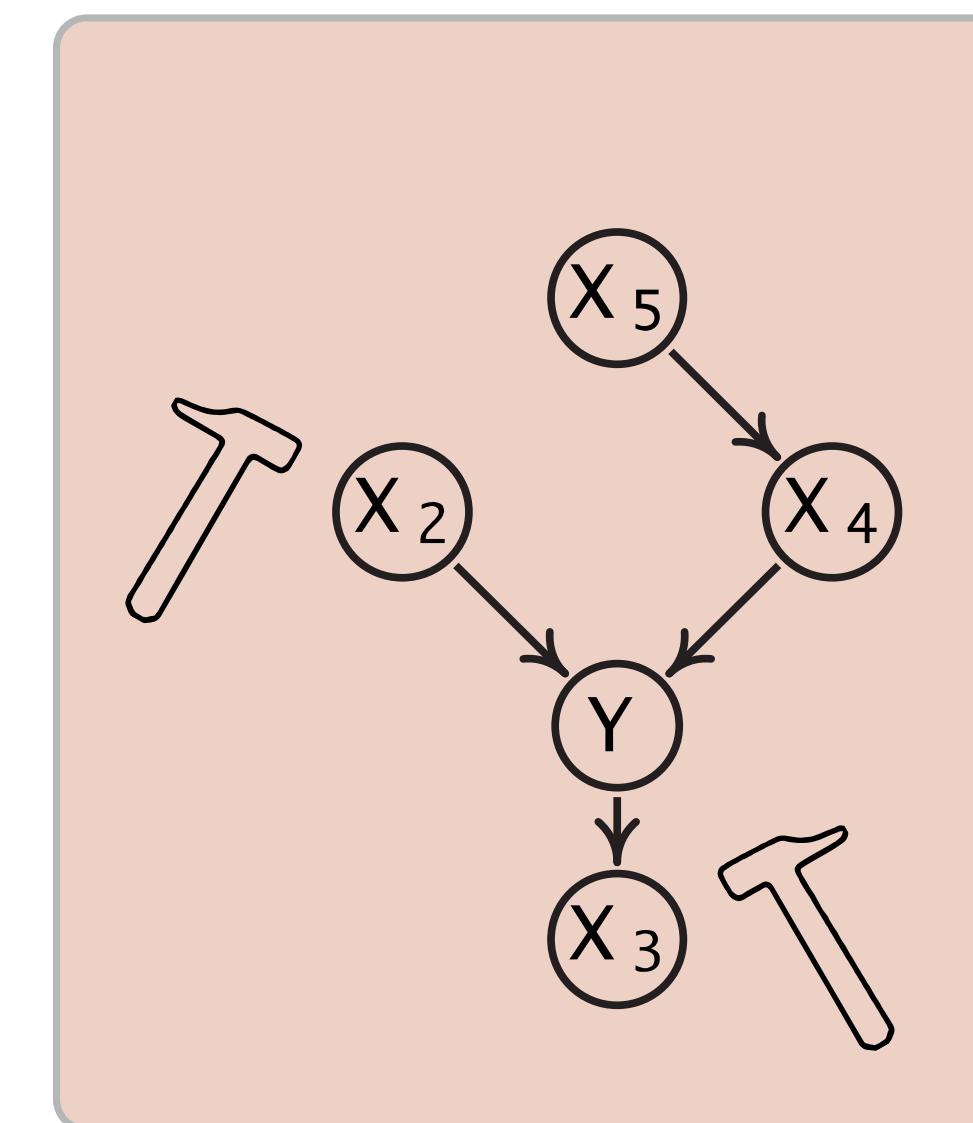
# Causal Discovery by Invariant Predictions

## The premise of Causation by Invariance:

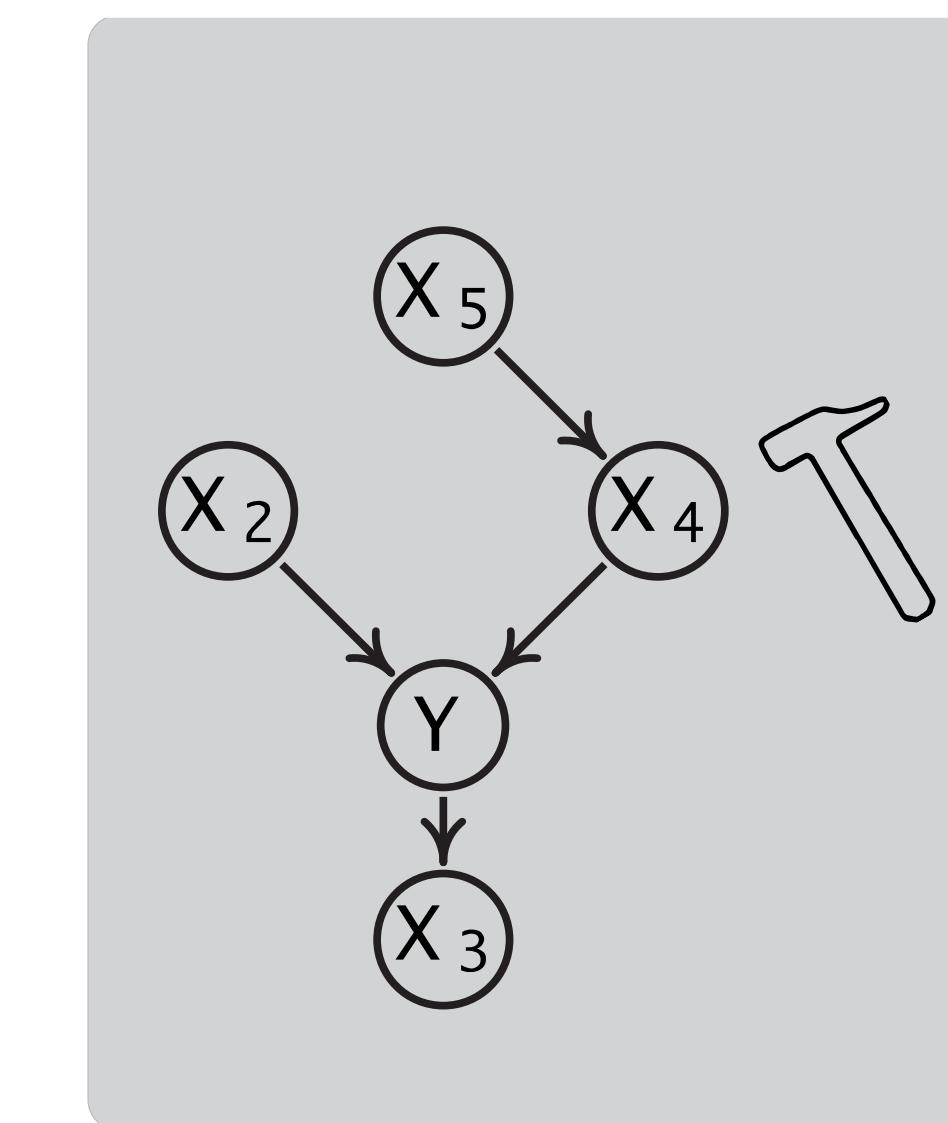
There is a causal structure that remains *invariant* across multiple experiments (perturbation-invariant)



No perturbation

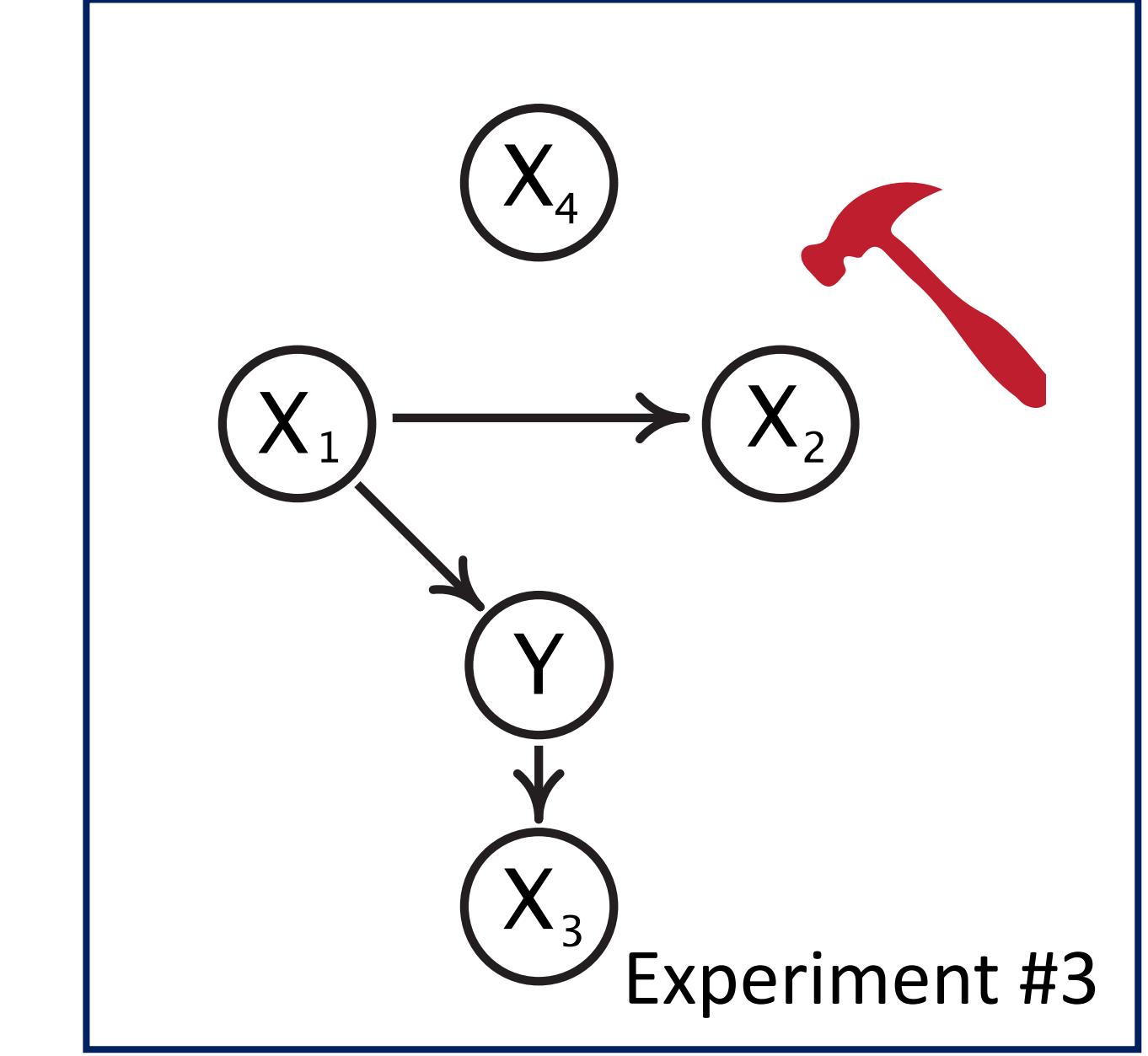
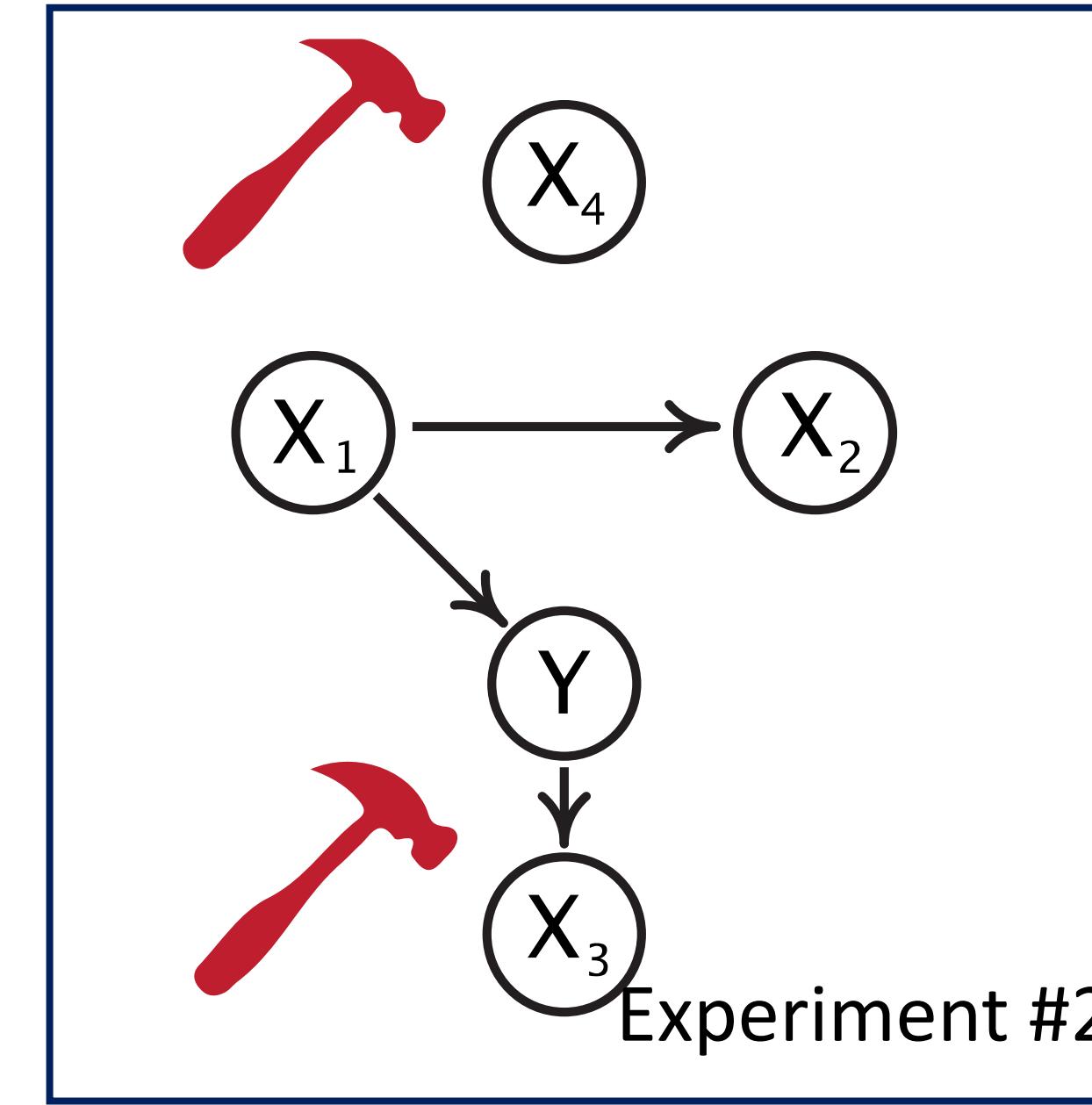
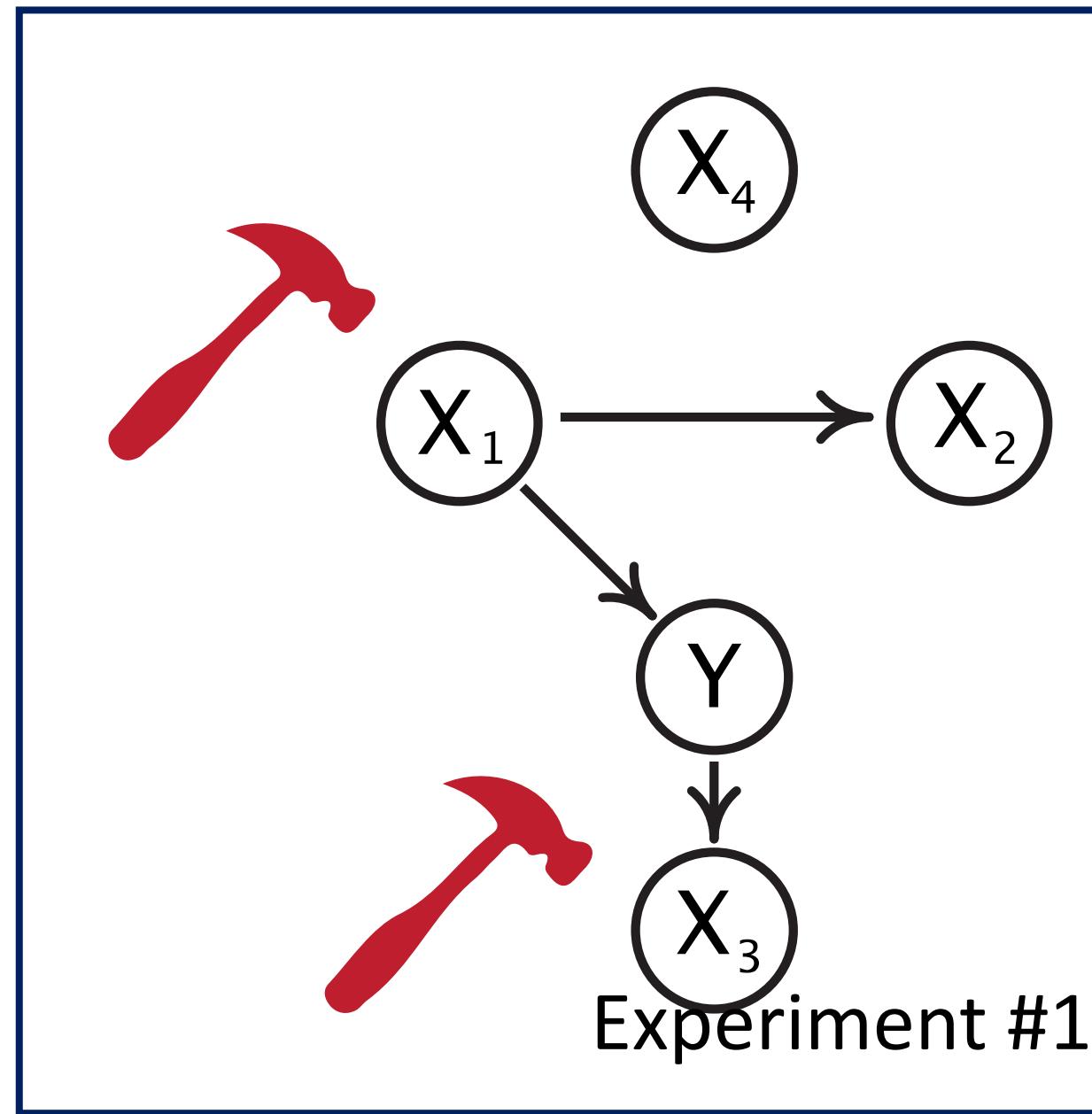


Gene knock-out #2, #3



Gene knock-out #4

# Causal Discovery by Invariance (toy example)



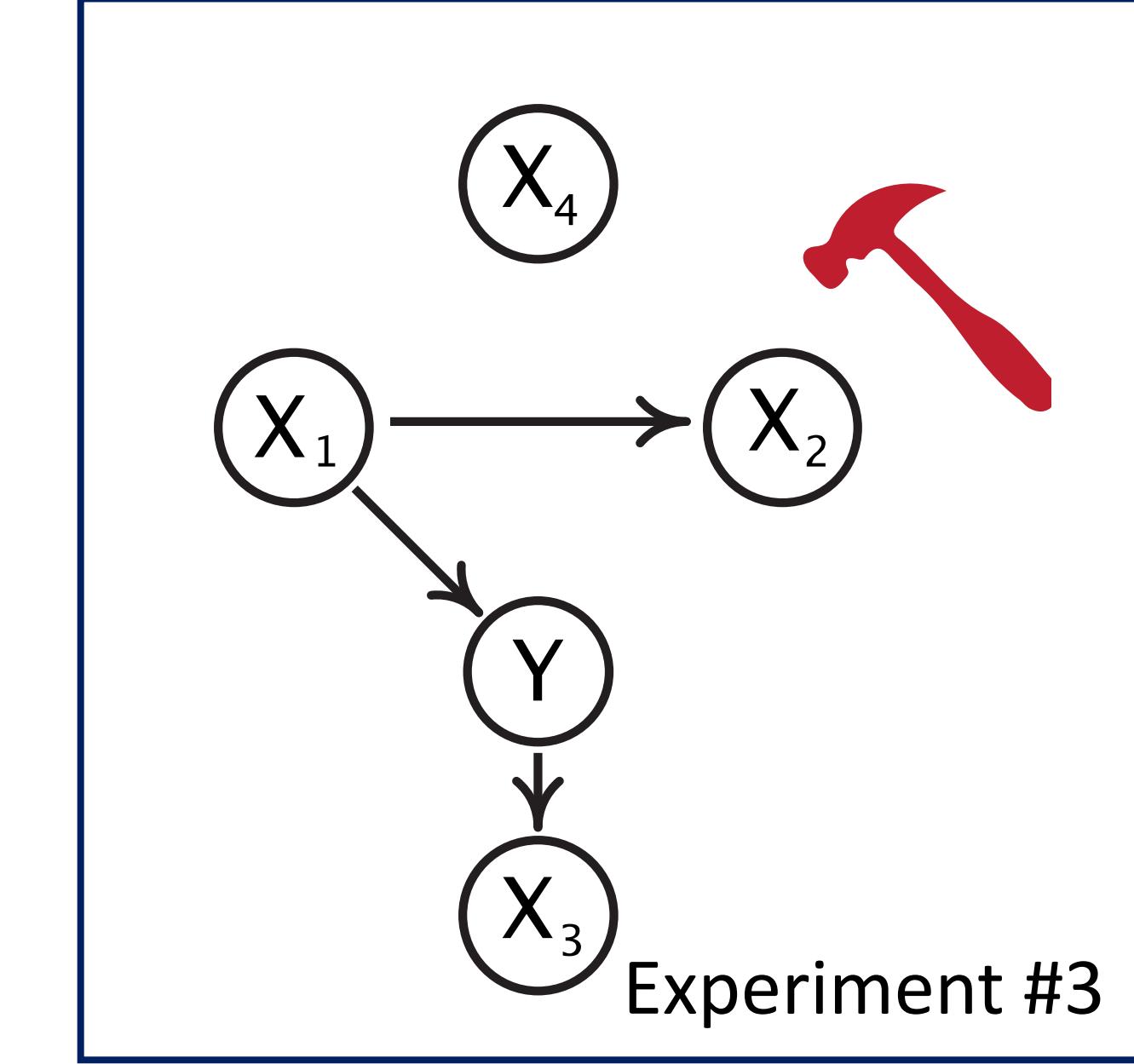
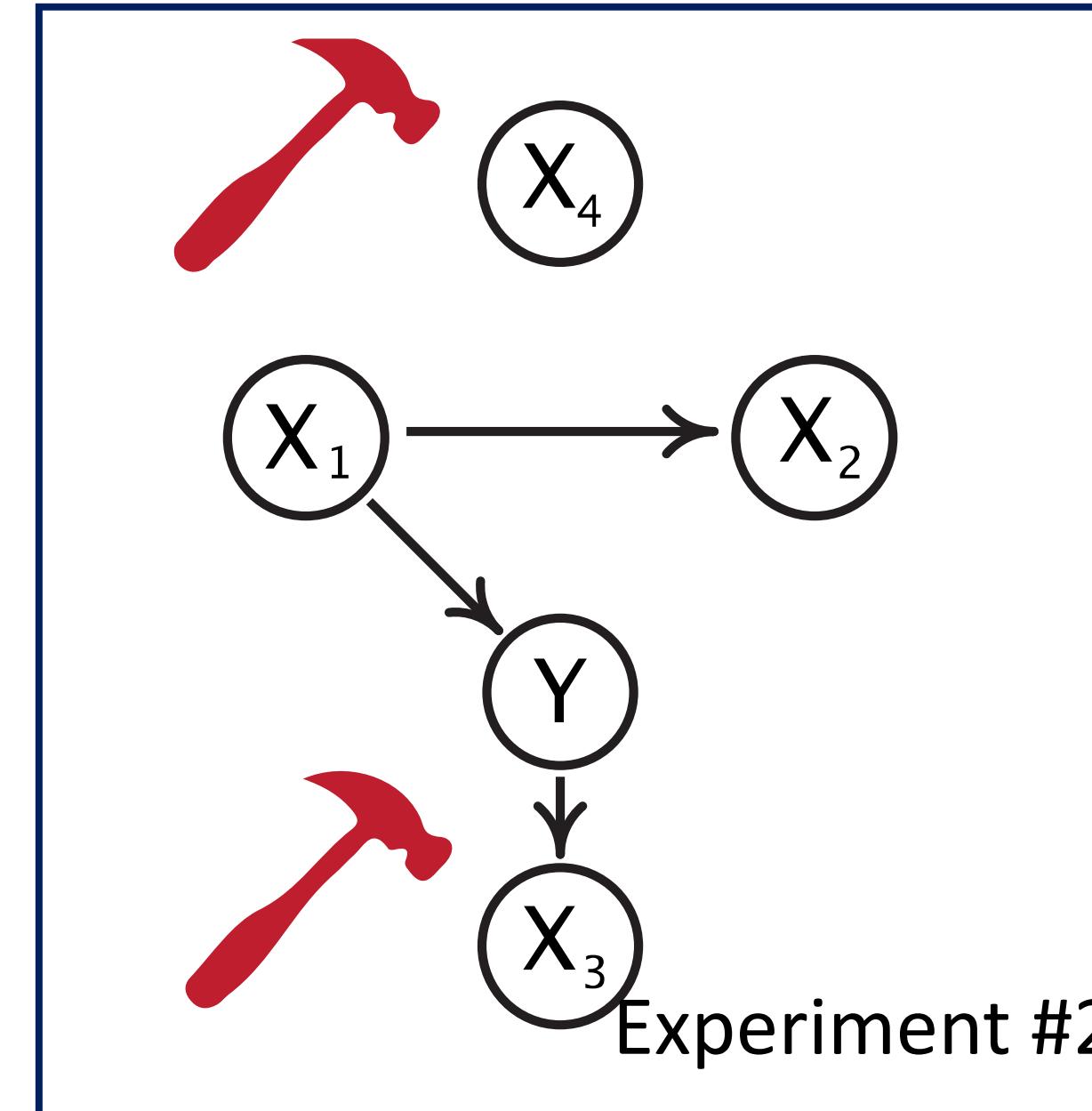
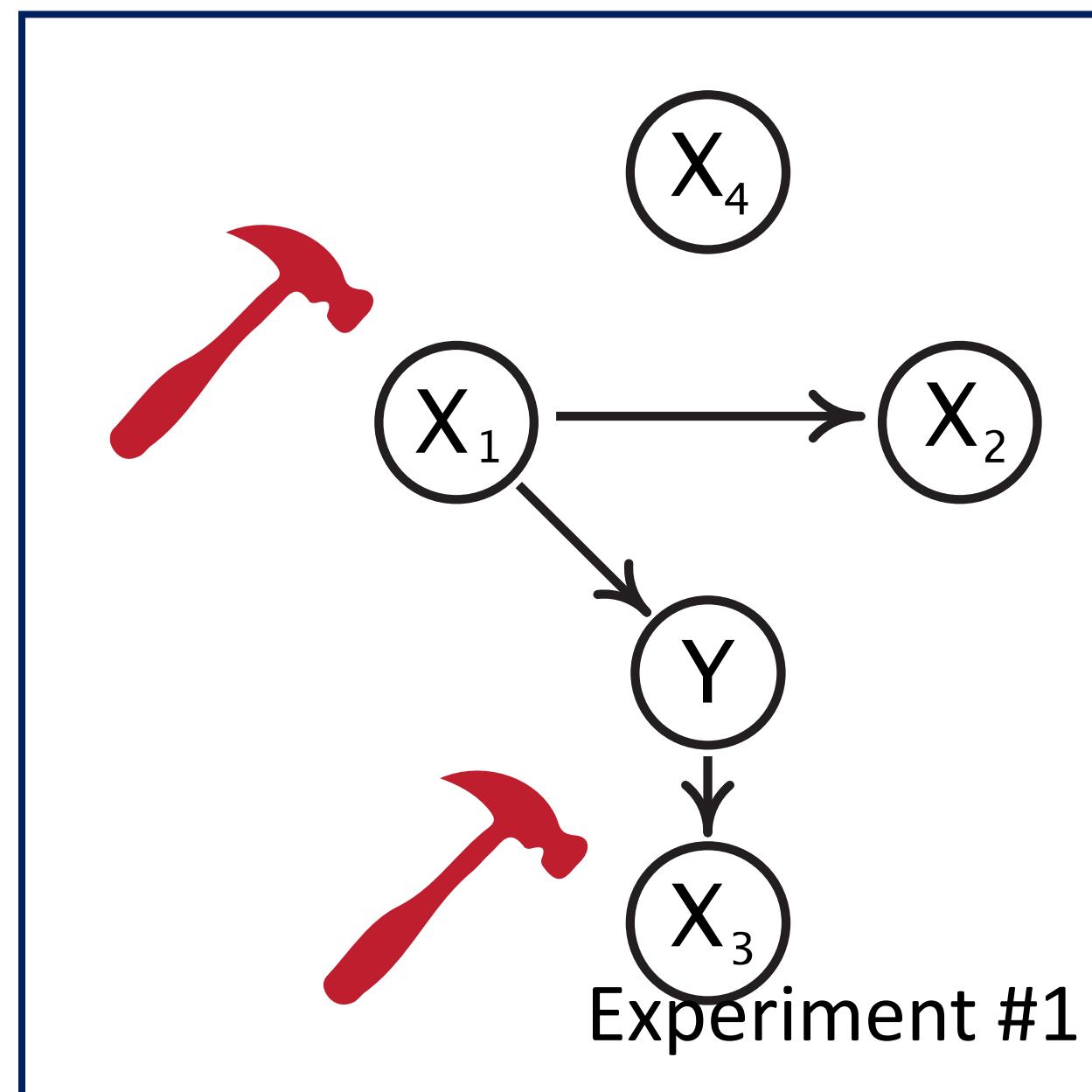
**Knock-out**  $\text{do}(X=0)$

# Causal Discovery by Invariance $\approx$ excluding non-parental variables

Q. Does this model hold the invariance condition?

$$Y = X_1 + \varepsilon$$

do( $X=0$ )



Observation:

$$Y = \varepsilon_y$$

$$X_1 = 0$$

$$Y = X_1 + \varepsilon_y$$

$$X_1 = \varepsilon_1$$

$$Y = X_1 + \varepsilon_y$$

$$X_1 = \varepsilon_1$$

# Does this model make sense in all the experiments?

Observation:

$$Y = \varepsilon_y$$

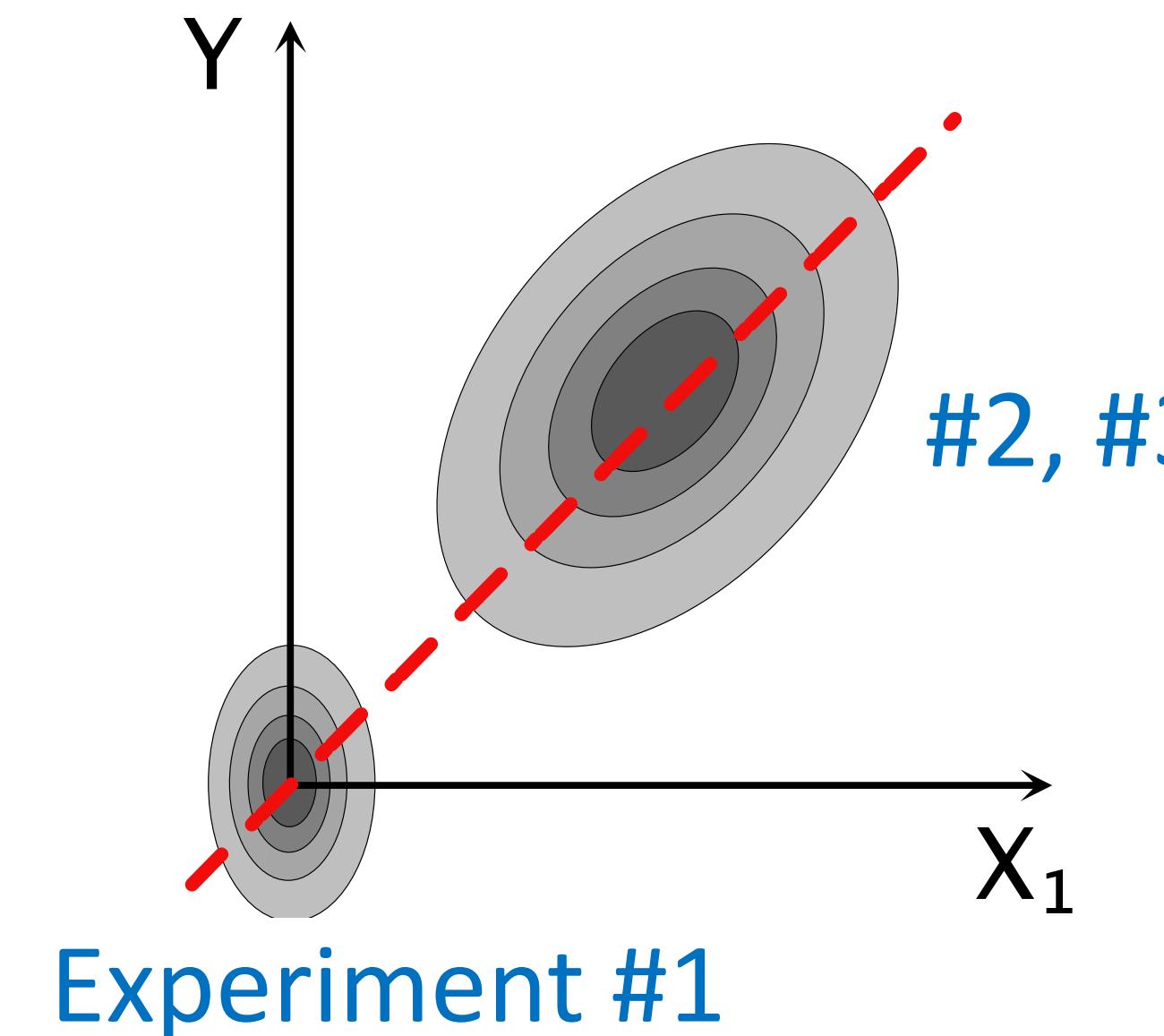
$$Y = X_1 + \varepsilon_y$$

$$Y = X_1 + \varepsilon_y$$

Experiment #1

Experiment #2

Experiment #3



Regression Residuals:

$$\begin{array}{ll} X_1 + \varepsilon_y \sim 0 & \#1 \\ X_1 + \varepsilon_y \sim X_1 & \#2 \\ X_1 + \varepsilon_y \sim X_1 & \#3 \end{array}$$

The same distribution  
of the residuals

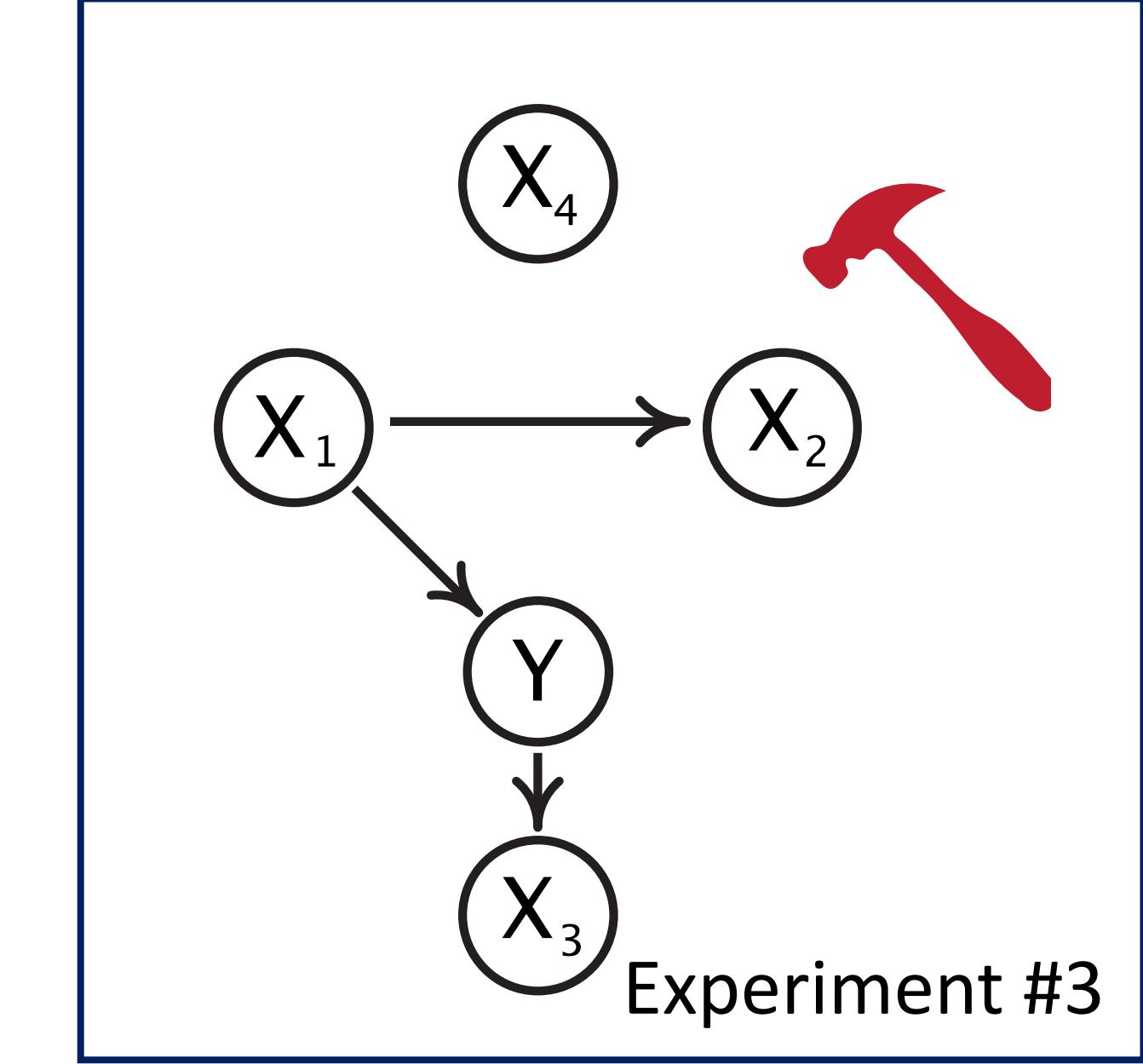
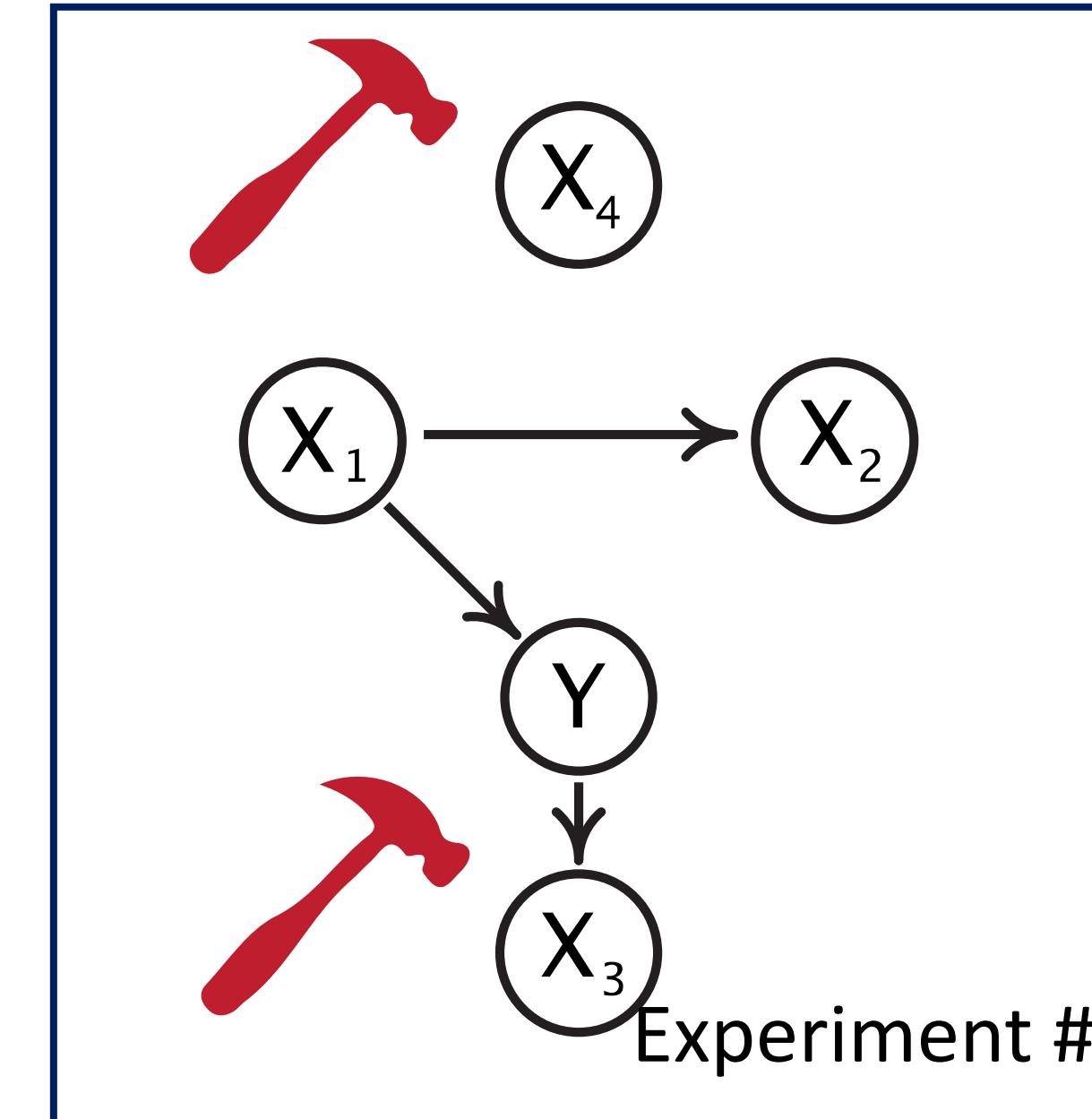
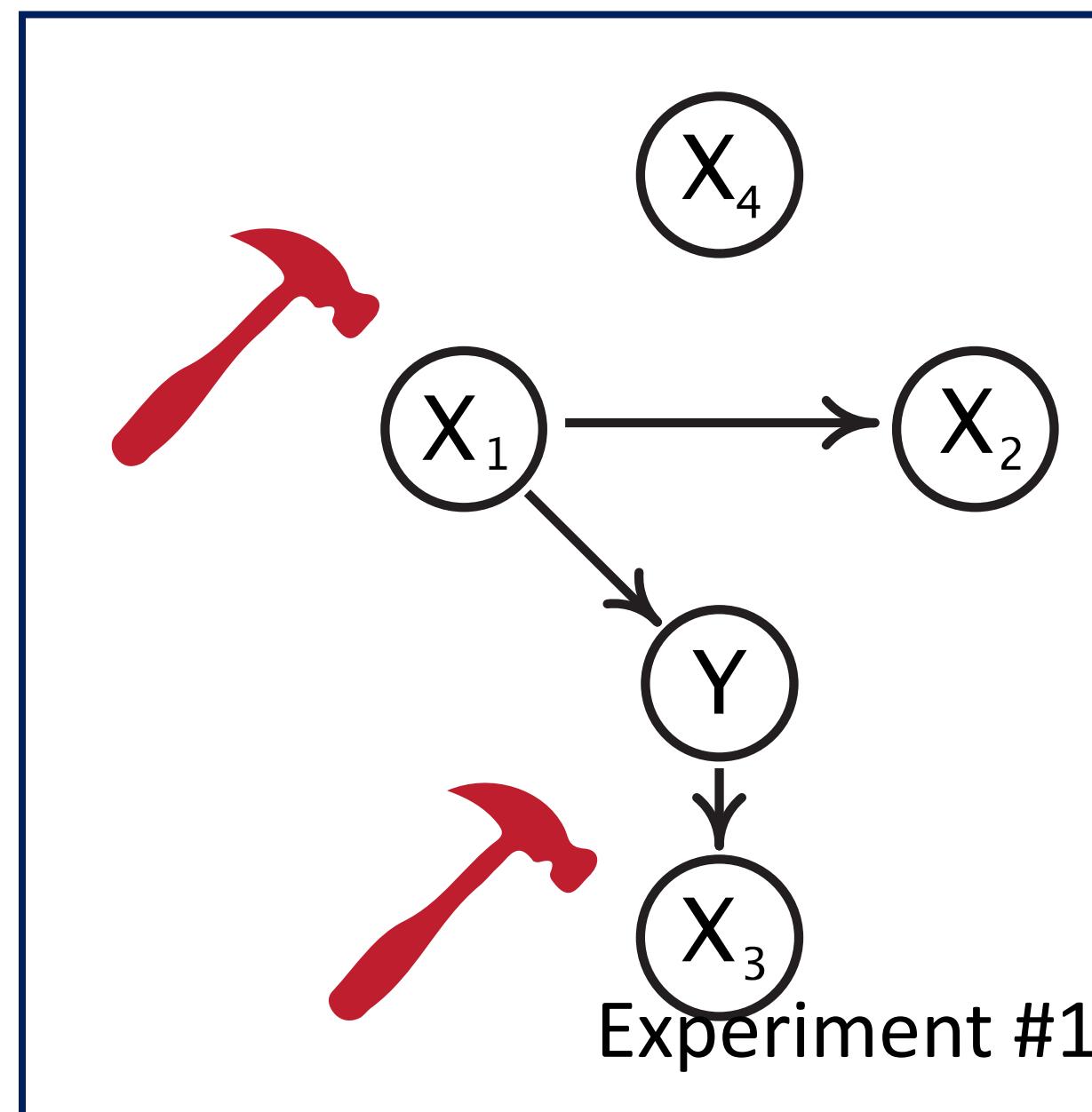
Cannot reject any  
experiments

# What if we specify a wrong predictor?

Q. Does this model hold the invariance condition?

$$Y = X_1 + \textcolor{red}{X}_3 + \varepsilon$$

do( $X=0$ )



Observation:

$$Y = \varepsilon$$

$$X_1 = 0, X_3 = 0$$

$$Y = X_1 + \varepsilon$$

$$X_1 = \varepsilon_1, X_3 = 0$$

$$Y = X_1 + \varepsilon$$

$$X_1 = \varepsilon_1, X_3 = Y + \varepsilon_3$$

# Does this model make sense in all the experiments?

Observation:

$$\#1 \quad Y = \varepsilon_y \quad X_3 = 0$$

$$\#2 \quad Y = X_1 + \varepsilon_y \quad X_3 = 0$$

$$\#3 \quad Y = X_1 + \varepsilon_y \quad X_3 = \varepsilon_3$$

Regression Residuals:

$$\boxed{\varepsilon_y} \sim 0$$

$$\rightarrow \cancel{X_1} + \boxed{\varepsilon_y} \sim \cancel{X_1} + 0$$

$$\cancel{X_1} + \cancel{\varepsilon_y} \sim X_1 + Y + \varepsilon_3$$

$$\sim \cancel{X_1} + (X_1 + \cancel{\varepsilon_y}) + \varepsilon_3$$

#1

#2

#3

**Invariance violated**

Reject the experiment #3

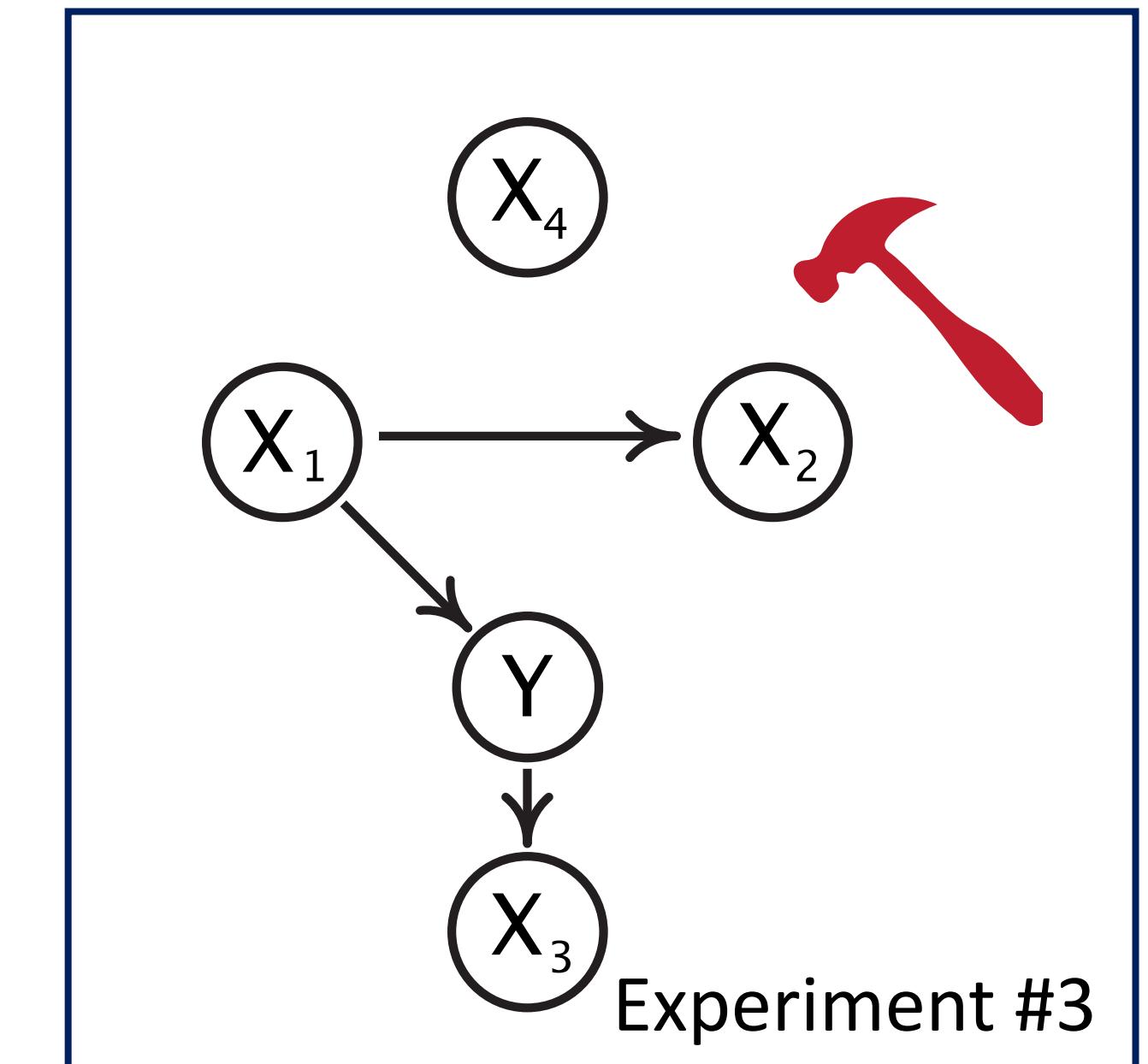
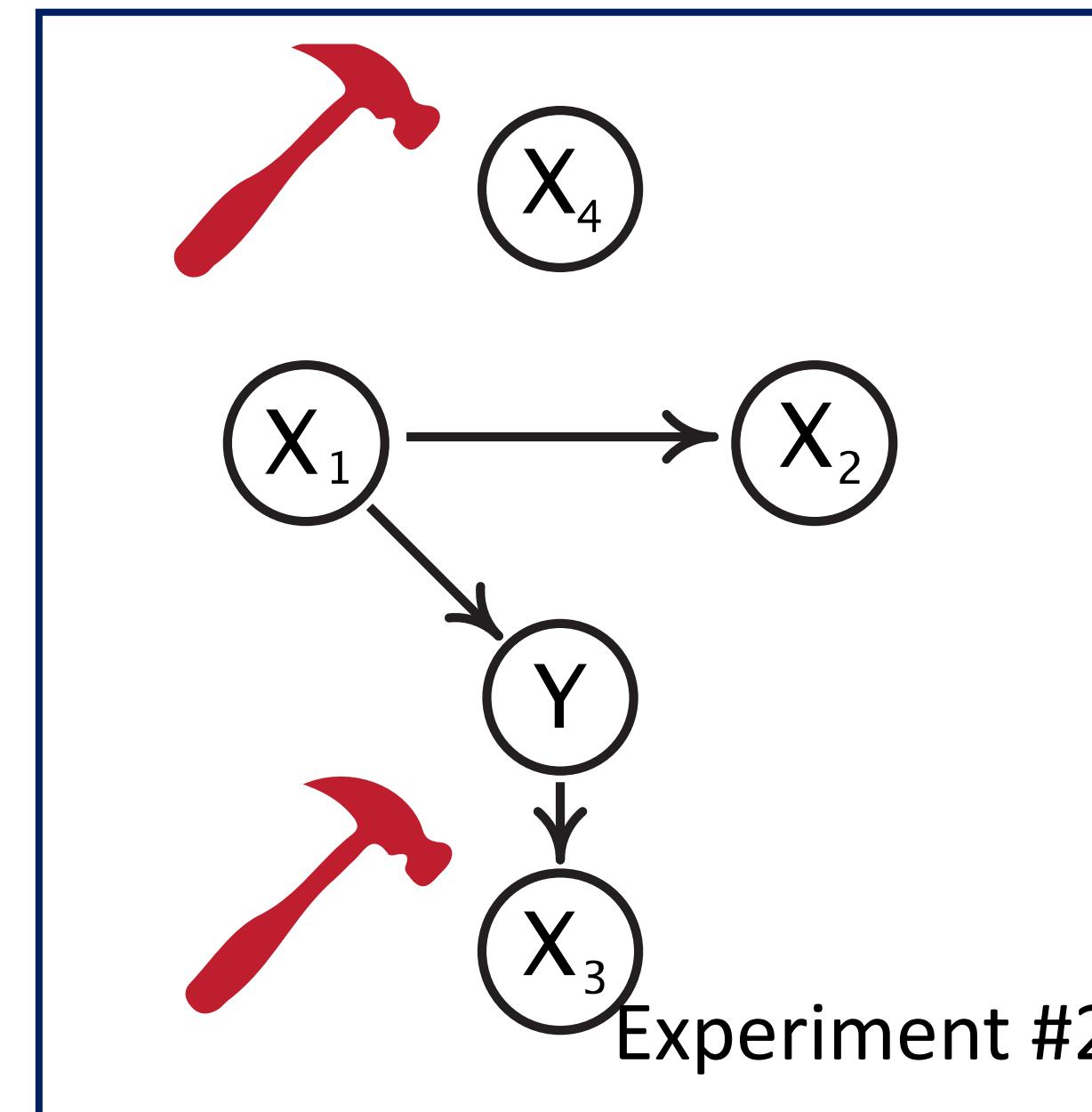
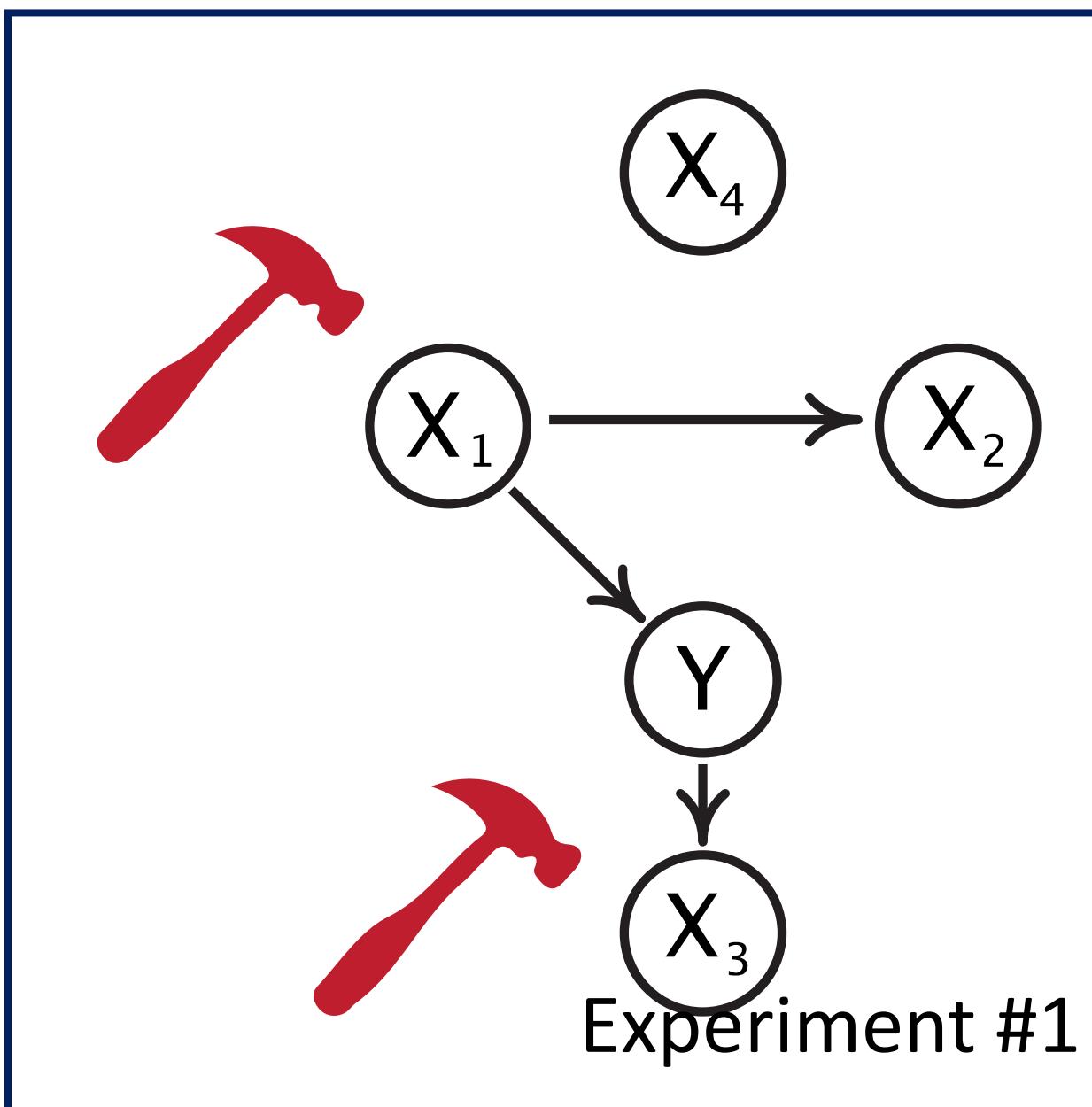
→ Reject the invariance of  
 $X_3$

# What if we specify a wrong predictor?

do( $X=0$ )

Q. Does this model hold the invariance condition?

$$Y = X_2 + \varepsilon_y$$



Observation:

$$Y = \varepsilon$$

$$X_2 = 0$$

$$Y = X_1 + \varepsilon$$

$$X_2 = X_1 + \varepsilon_2$$

$$Y = X_1 + \varepsilon$$

$$X_2 = 0$$

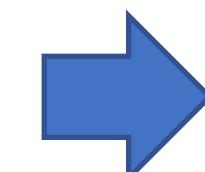
# Does this model make sense in all the experiments?

Observation:

$$\#1 \quad Y = \varepsilon_y \quad X_2 = \varepsilon_2$$

$$\#2 \quad Y = X_1 + \varepsilon_y \quad X_2 = X_1 + \varepsilon_2$$

$$\#3 \quad Y = X_1 + \varepsilon_y \quad X_2 = 0$$



Regression Residuals:

$$\begin{aligned} & \cancel{X}_1 + \varepsilon_y \sim \cancel{X}_1 + \varepsilon_2 \\ & \boxed{X}_1 + \varepsilon_y \sim \varepsilon_1 \end{aligned}$$

#1

#2

#3

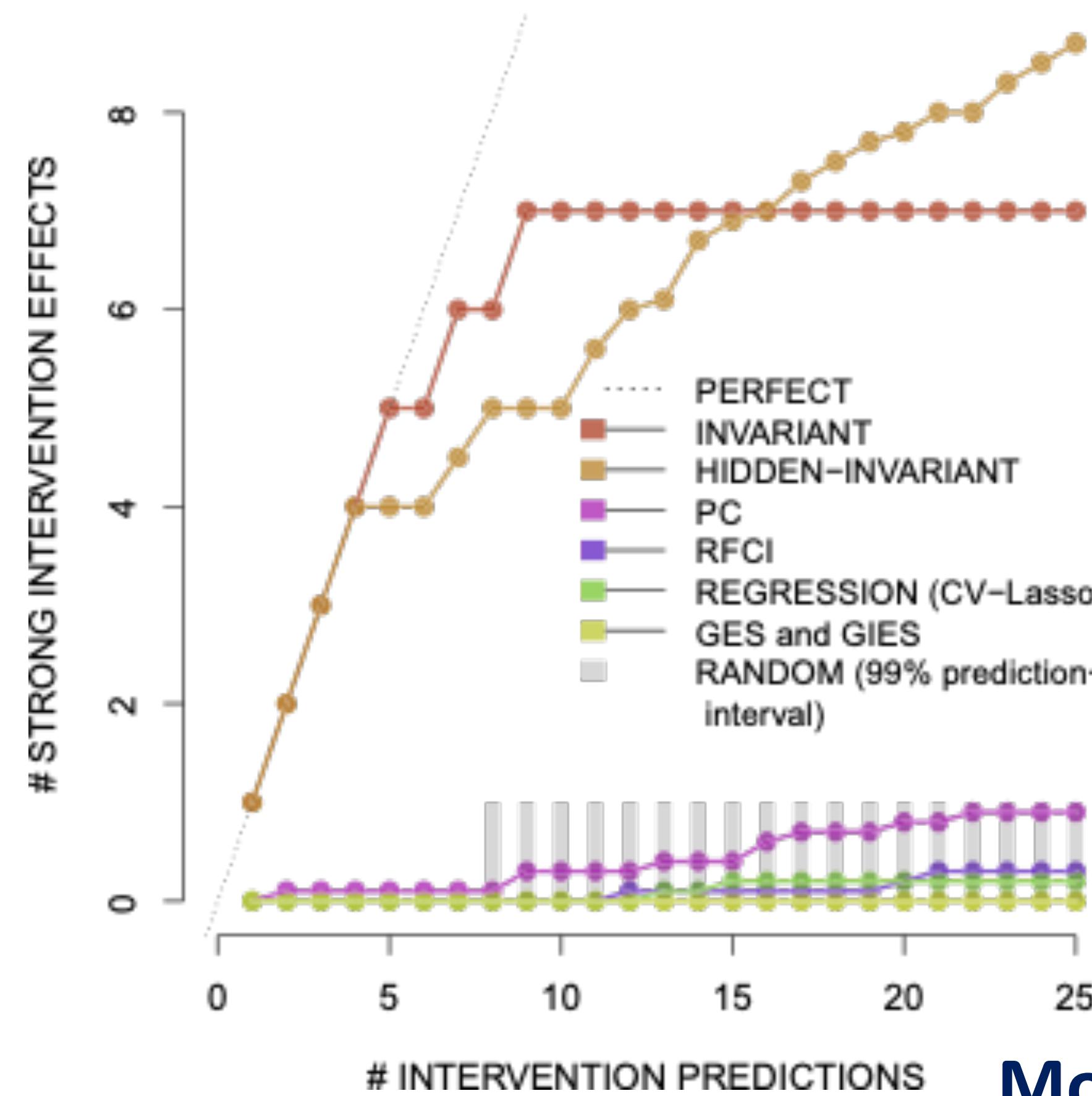
Invariance violated

Reject the experiment #3

→ Reject the invariance of  
 $X_2$

# Causal Model Discovery by Invariance outperforms regression-based methods

Evaluation  
of the actual  
Knock-out  
experiments



Causal regression

$$\mathbb{E}[y | do(X_j = x)]$$

Regression on observed data

$$\mathbb{E}[y | X = x]$$

More perturbation Exp.

# Causal triangulation to gain confidence in scientific explanations

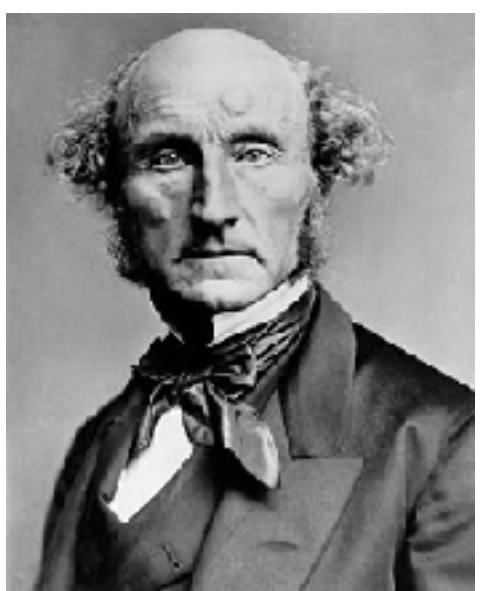
Multiple lines of “orthogonal evidence”

## Observation

1. The occurrence of  $X_1$  and  $X_2$  with  $Y_1$  and  $Y_2$
2. The occurrence of  $X_2$  with  $Y_2$

## Induction (conclusion)

$$X_1 \rightarrow Y_1 \quad \text{or} \quad X_1 \leftarrow Y_1$$

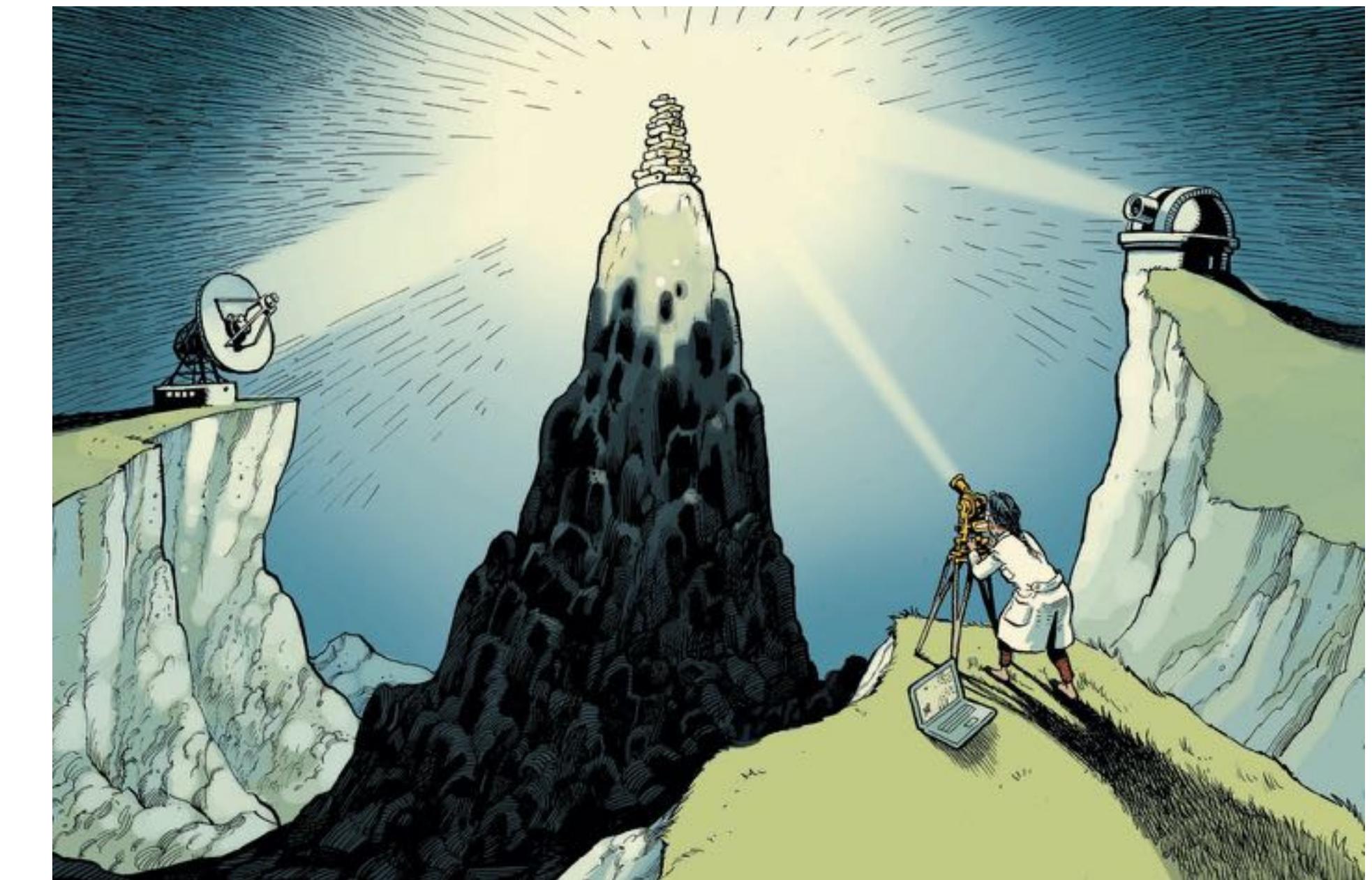


John Stuart Mill

### SECOND CANON.

*If an instance in which the phenomenon under investigation occurs, and an instance in which it does not occur, have every circumstance save one in common, that one occurring only in the former; the circumstance in which alone the two instances differ, is the effect, or cause, or a necessary part of the cause, of the phenomenon.*

JS Mill, *A system of Logic* (1843)



Munafo & Davey Smith, "Repeating experiments is not enough", Nature (2018)



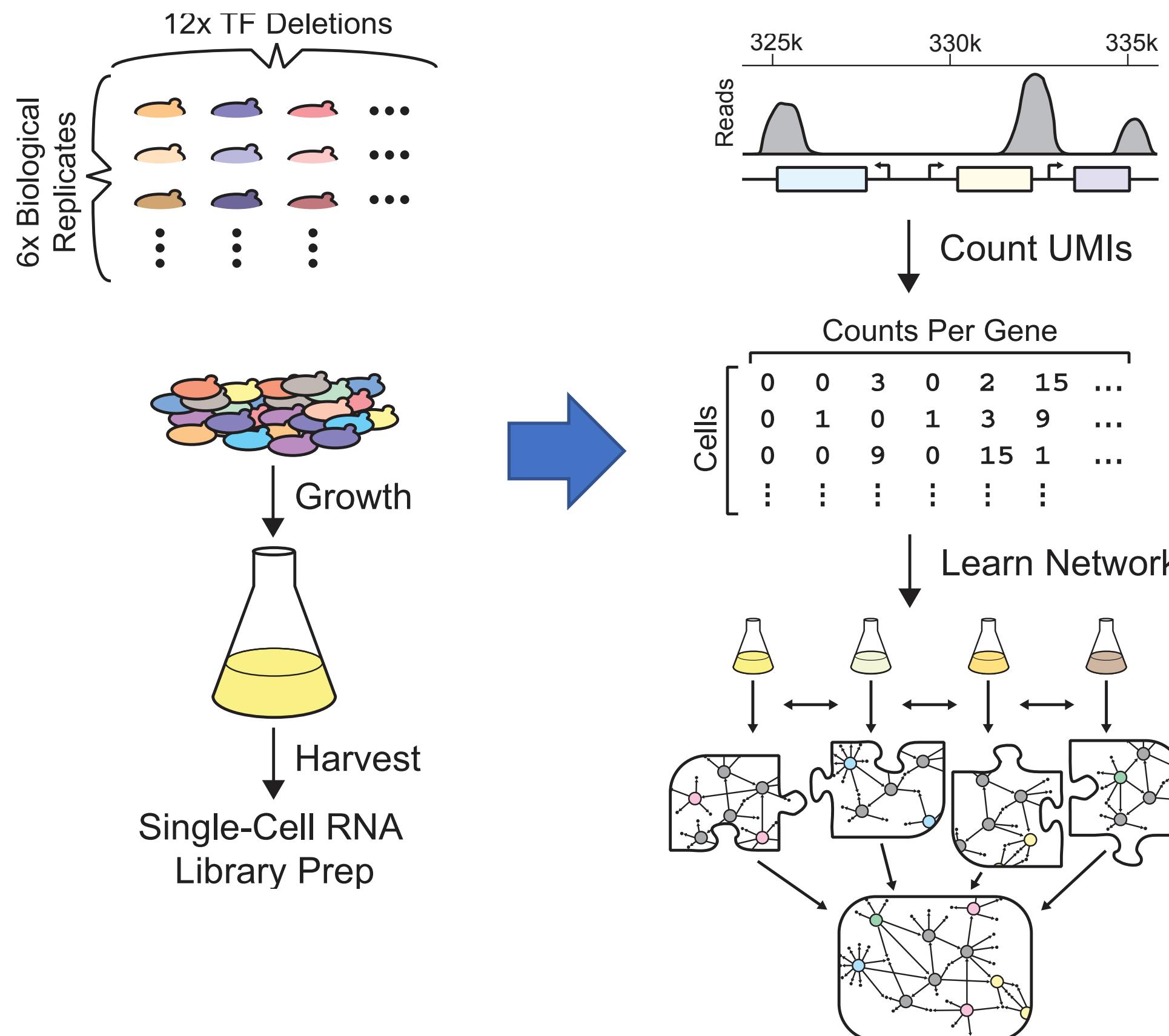
Peter Lipton

Contrastive  
Explanation &  
causal  
triangulation,  
Philosophy of  
Science (1991)



George  
Davey Smith

# Can we borrow the awesome power of genetics?



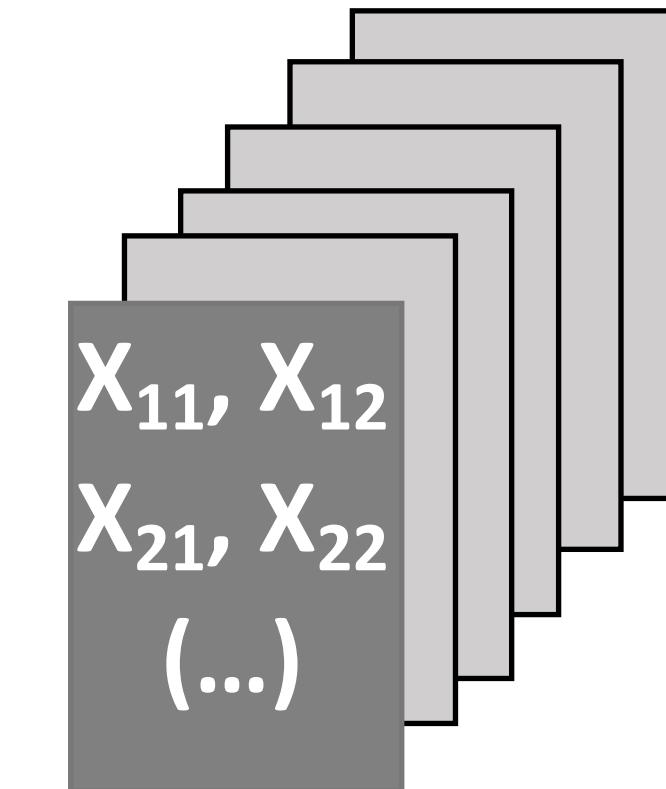
**Early 2000's**

$$\begin{matrix} X_{11}, X_{12}, X_{13} \\ X_{21}, X_{22}, X_{23} \\ \dots \end{matrix}$$

A single or only a handful of conditions

*Many computational biologists gave up on this NP-hard problem... learning BN from data*

**After perturb-seq**

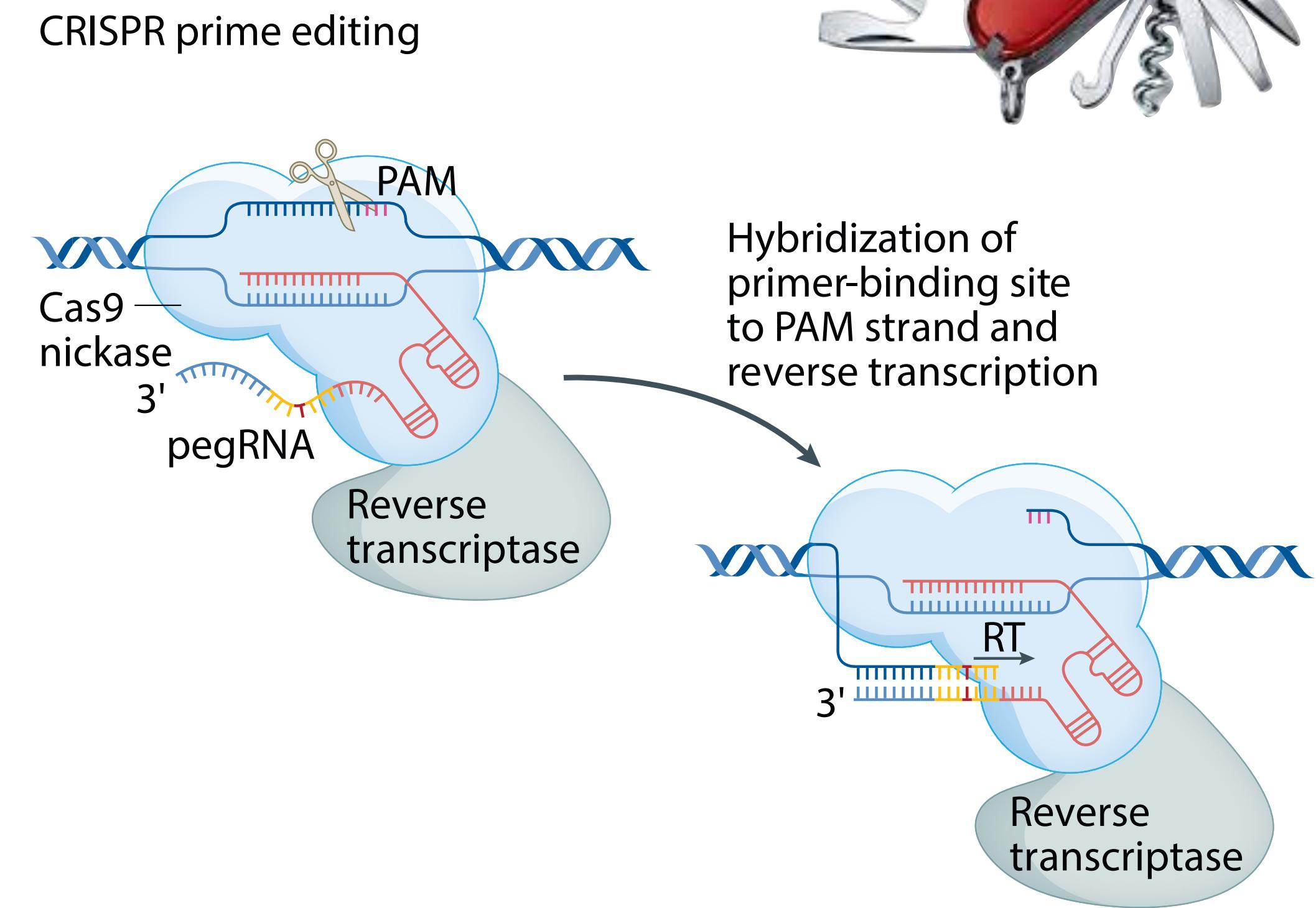
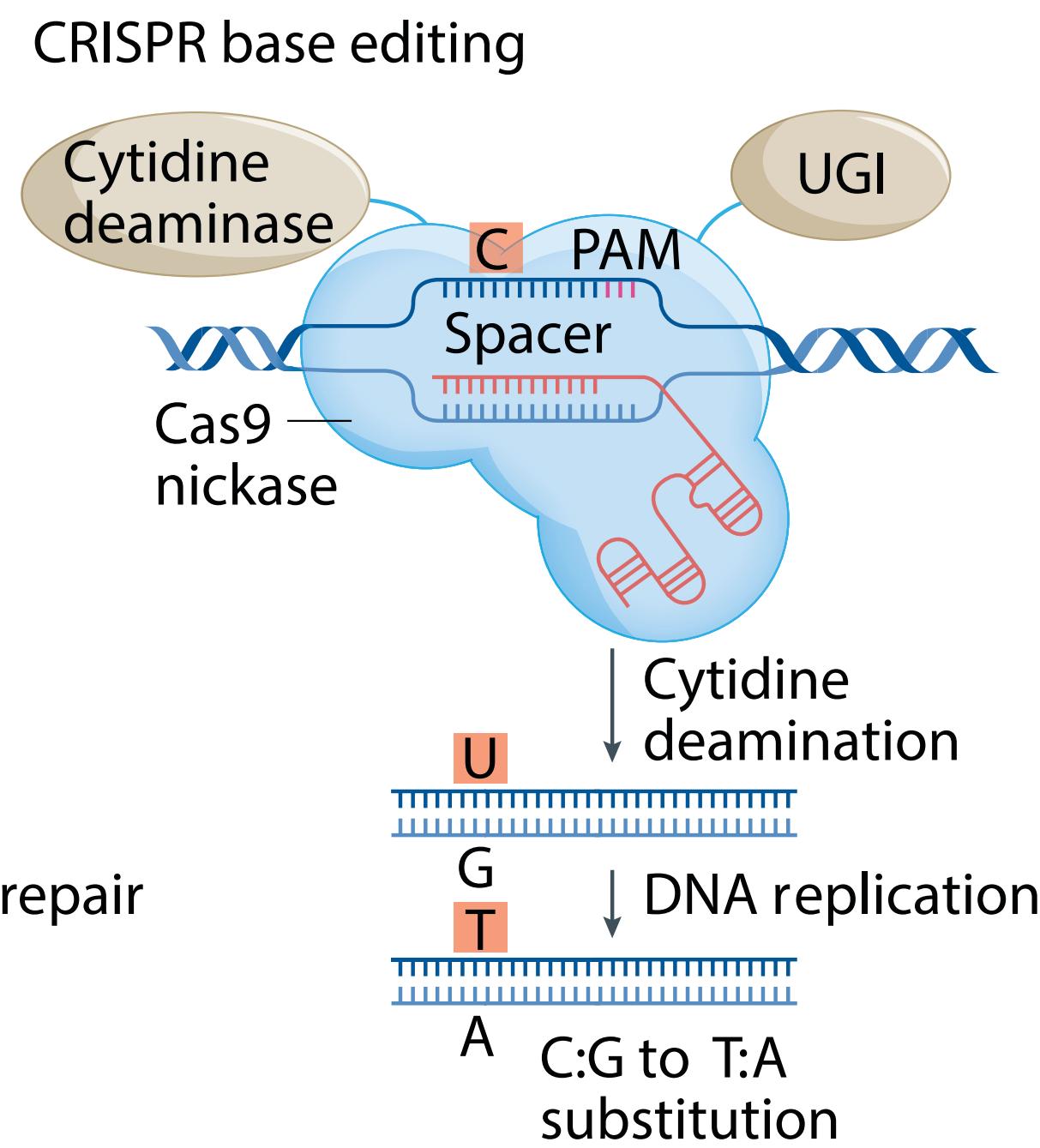
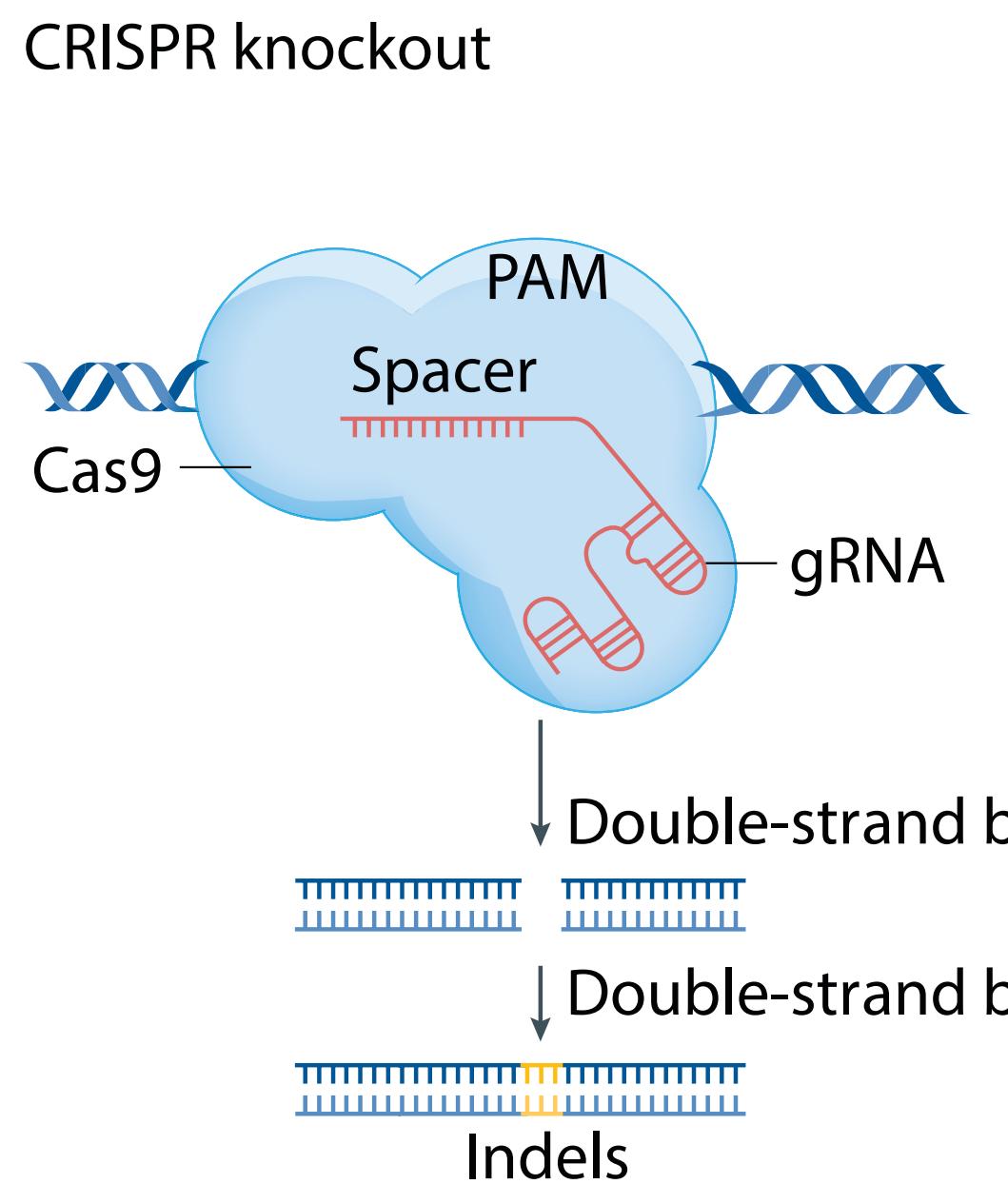


Massive (random) perturbation assays → hundreds of data matrices

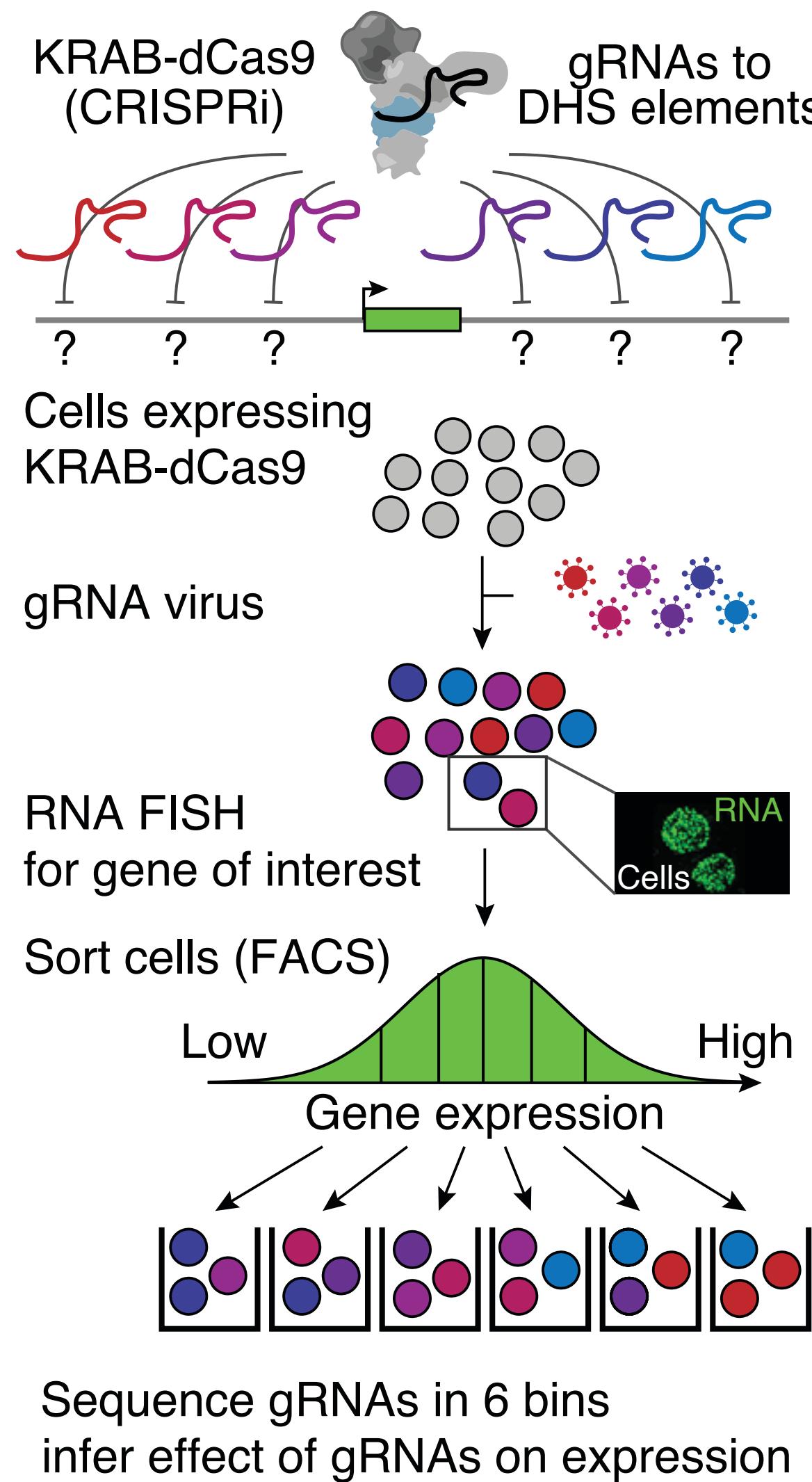
# CRISPR-dCas9 system made genetic perturbation studies so easy and scalable



# Molecular scissors for genome editing



# CRISPRi to disrupt gene regulatory elements of a target gene to interrogate connections between them



DE–G connections

Tested genes/  
regulatory DEs

Other genes

Distal elements

gRNA effect on  
gene expression (%)

GATA1 FlowFISH

–log<sub>10</sub> P

37

17

242

55

100 kb

22 kb

DHS H3K27ac

GATA1

HDAC6

Cells

RNA

+

0

-100

-200

-300

-400

-500

-600

-700

-800

-900

-1000

-1100

-1200

-1300

-1400

-1500

-1600

-1700

-1800

-1900

-2000

-2100

-2200

-2300

-2400

-2500

-2600

-2700

-2800

-2900

-3000

-3100

-3200

-3300

-3400

-3500

-3600

-3700

-3800

-3900

-4000

-4100

-4200

-4300

-4400

-4500

-4600

-4700

-4800

-4900

-5000

-5100

-5200

-5300

-5400

-5500

-5600

-5700

-5800

-5900

-6000

-6100

-6200

-6300

-6400

-6500

-6600

-6700

-6800

-6900

-7000

-7100

-7200

-7300

-7400

-7500

-7600

-7700

-7800

-7900

-8000

-8100

-8200

-8300

-8400

-8500

-8600

-8700

-8800

-8900

-9000

-9100

-9200

-9300

-9400

-9500

-9600

-9700

-9800

-9900

-10000

-10100

-10200

-10300

-10400

-10500

-10600

-10700

-10800

-10900

-11000

-11100

-11200

-11300

-11400

-11500

-11600

-11700

-11800

-11900

-12000

-12100

-12200

-12300

-12400

-12500

-12600

-12700

-12800

-12900

-13000

-13100

-13200

-13300

-13400

-13500

-13600

-13700

-13800

-13900

-14000

-14100

-14200

-14300

-14400

-14500

-14600

-14700

-14800

-14900

-15000

-15100

-15200

-15300

-15400

-15500

-15600

-15700

-15800

-15900

-16000

-16100

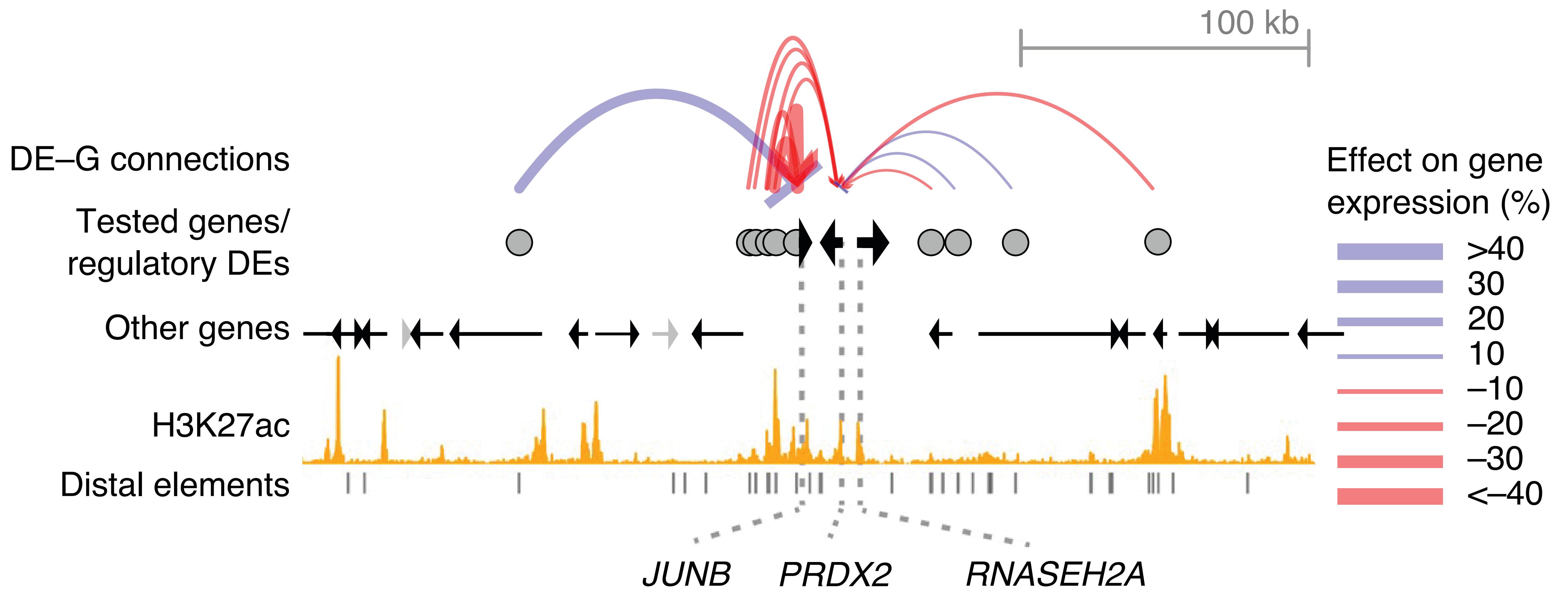
-16200

-16300

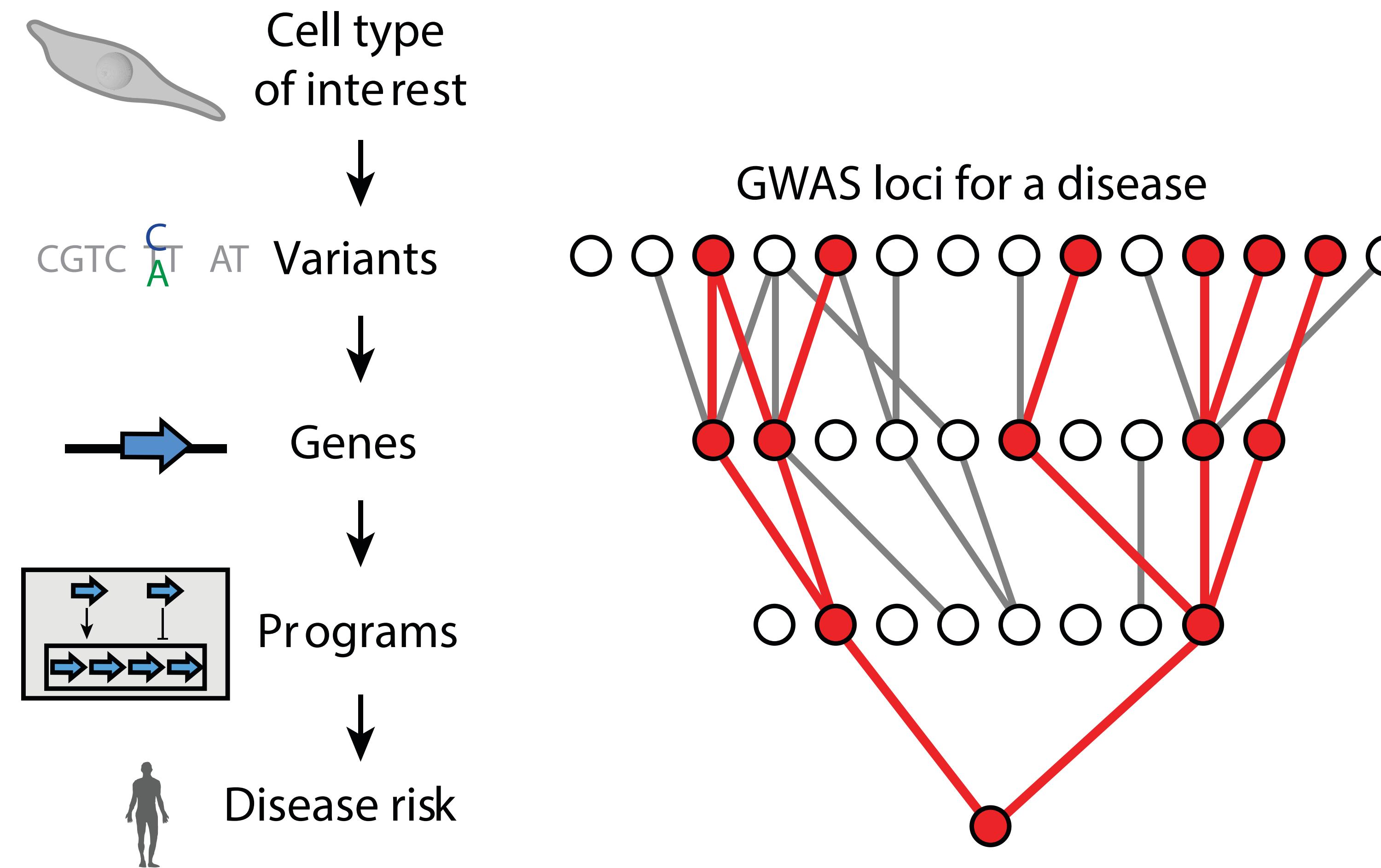
-16400

-16500

# We can systematically link epigenetic regions to target genes by experiments

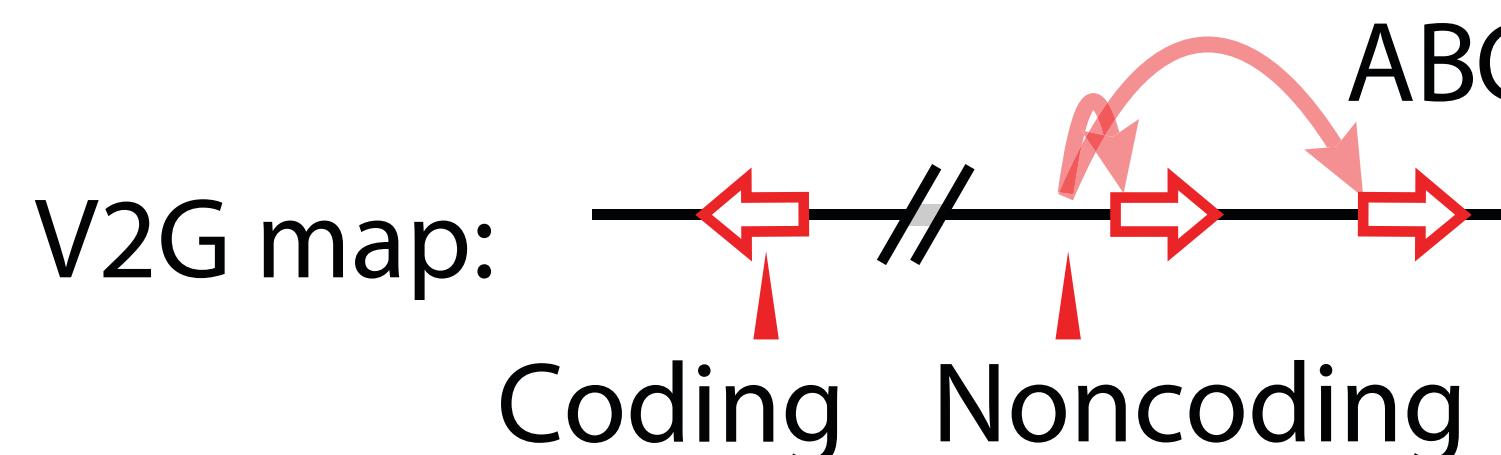


# GWAS variant to target gene

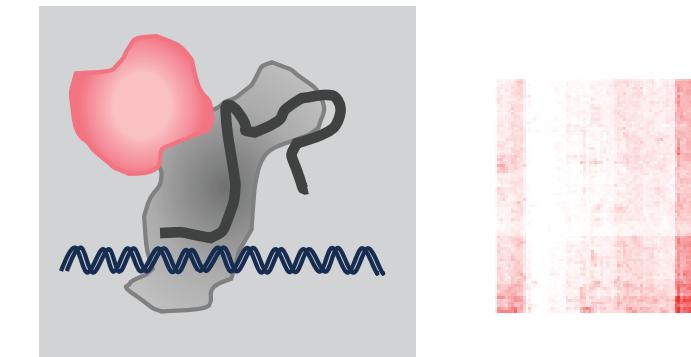


# GWAS variant to gene to pathway

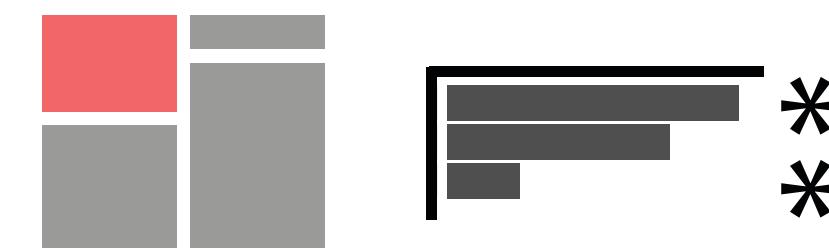
V2G2P enrichment test:



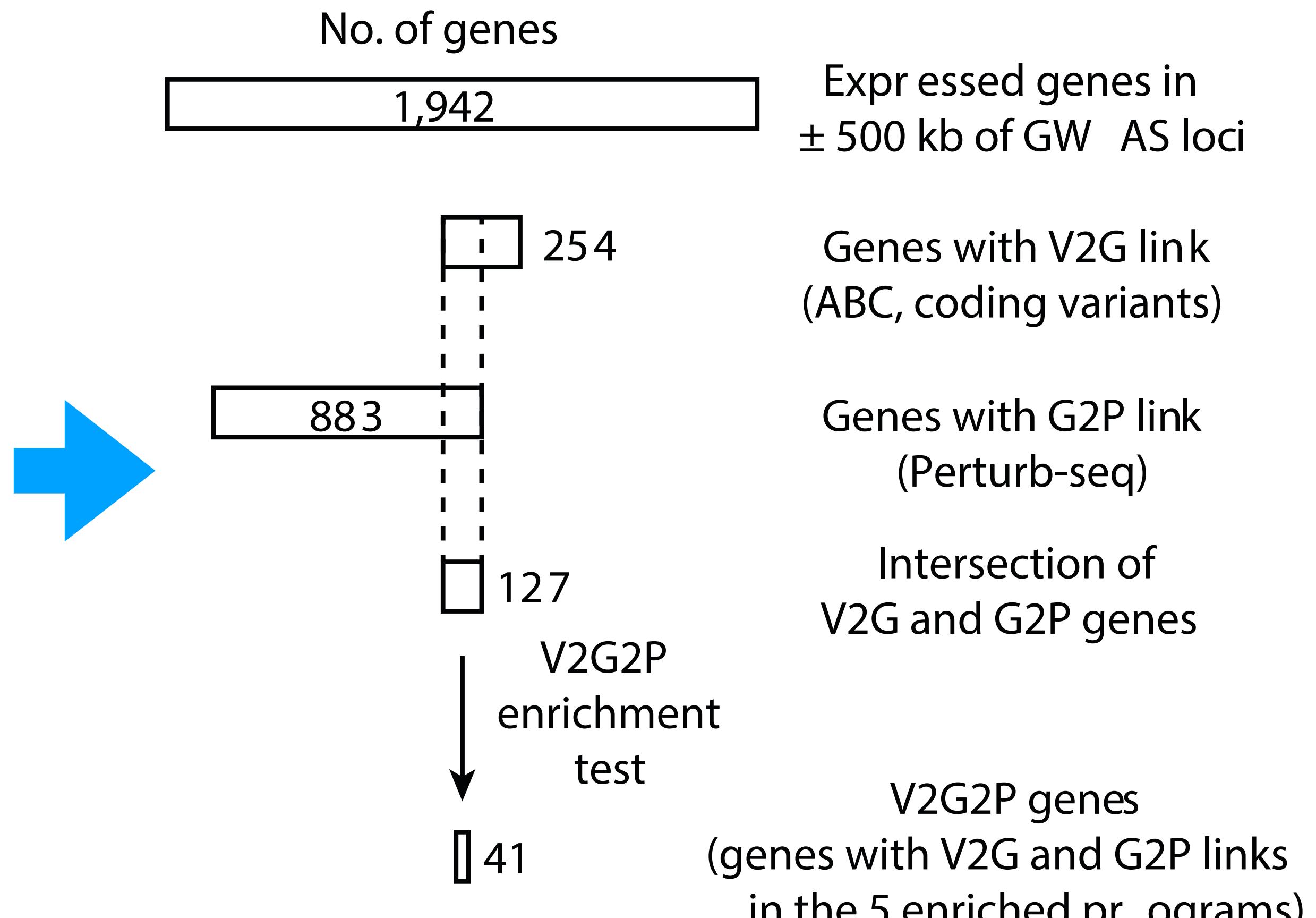
G2P map:



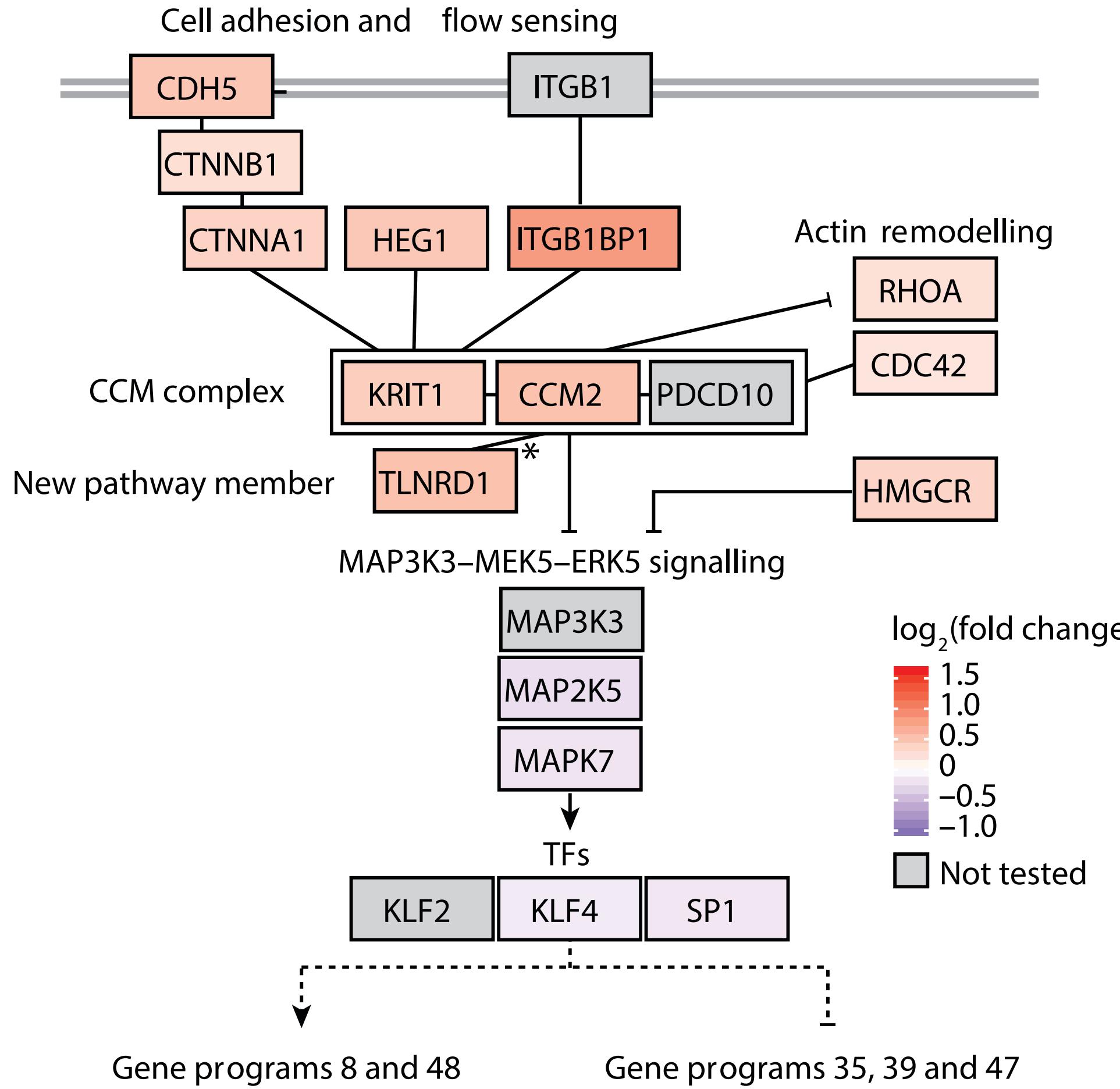
CRISPRi-Perturb-seq



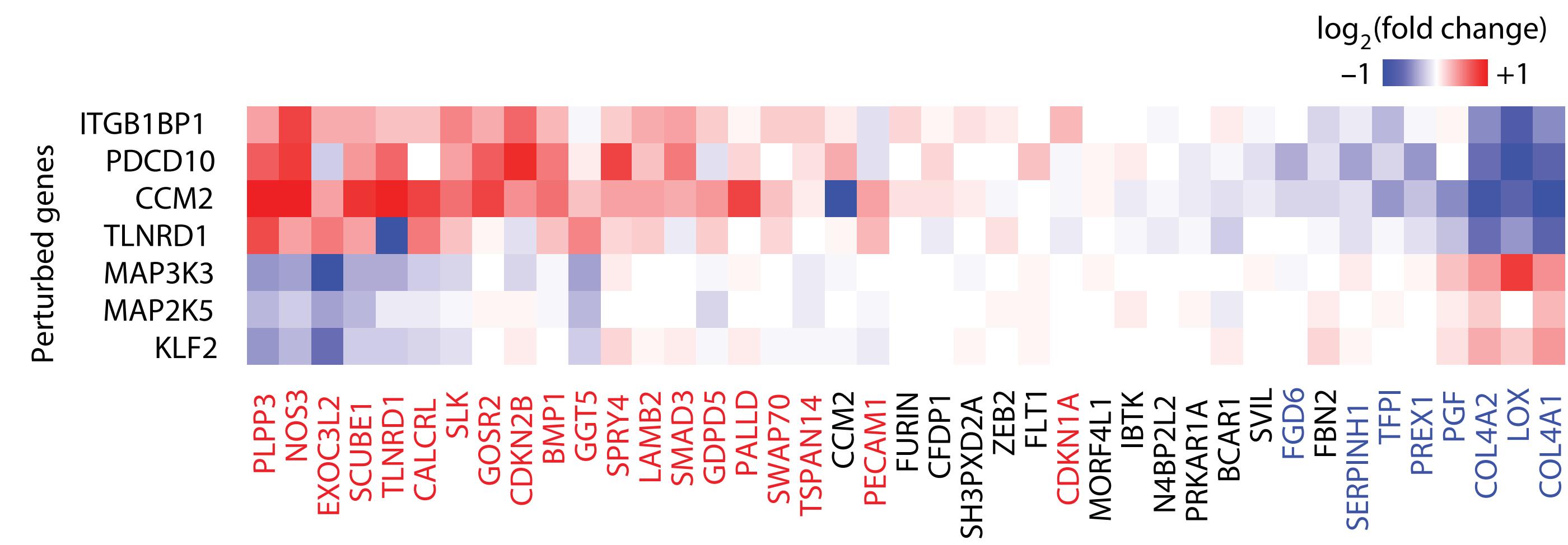
Prioritize programs



# How can we prove Variant to Gene to Pathway?



# A series of Knock-down experiments



# Multiple CRISPR screening can be combined

---

nature biotechnology



Article

<https://doi.org/10.1038/s41587-023-01964-9>

## Scalable genetic screening for regulatory circuits using compressed Perturb-seq

---

Received: 5 January 2023

---

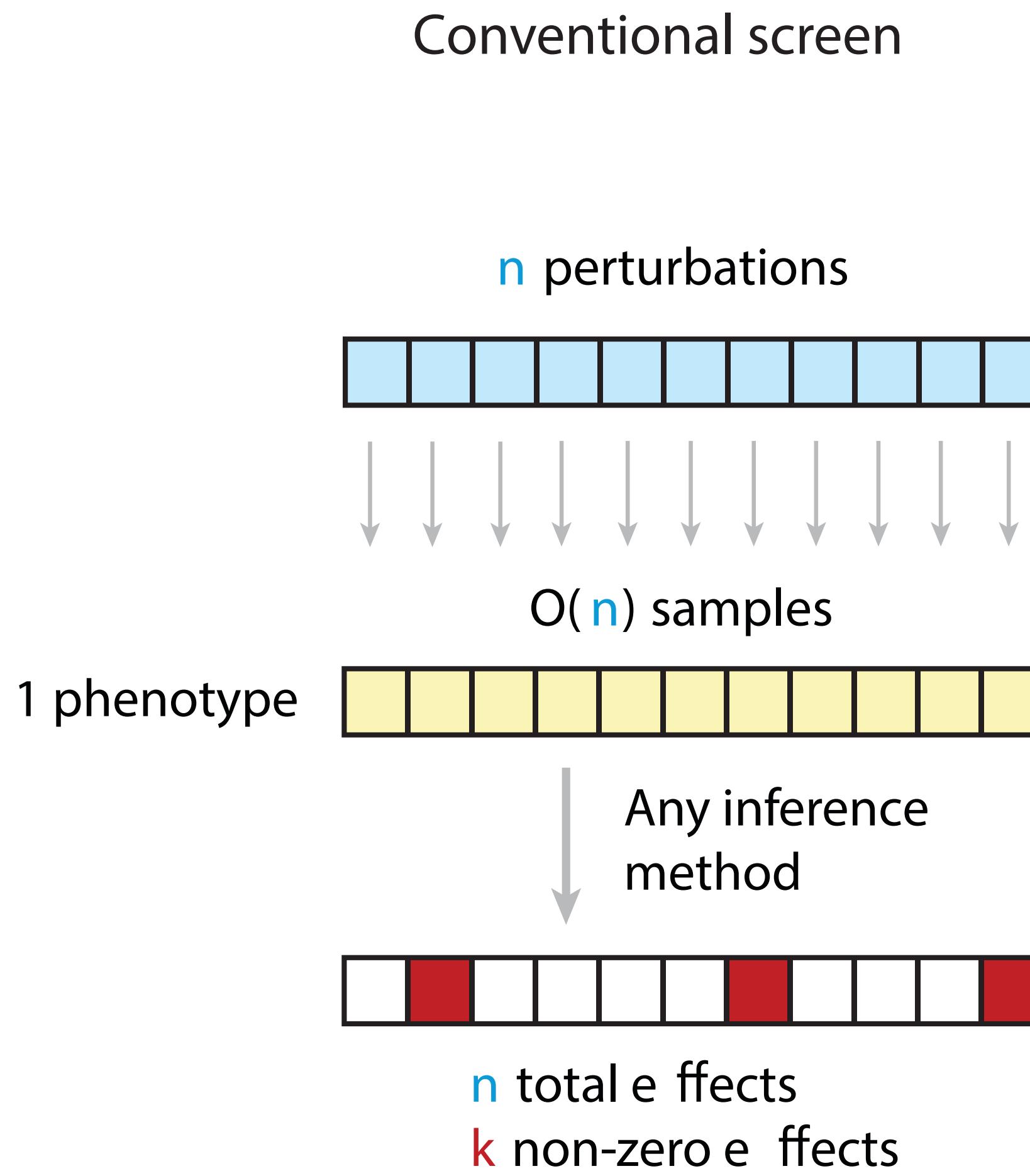
Accepted: 22 August 2023

---

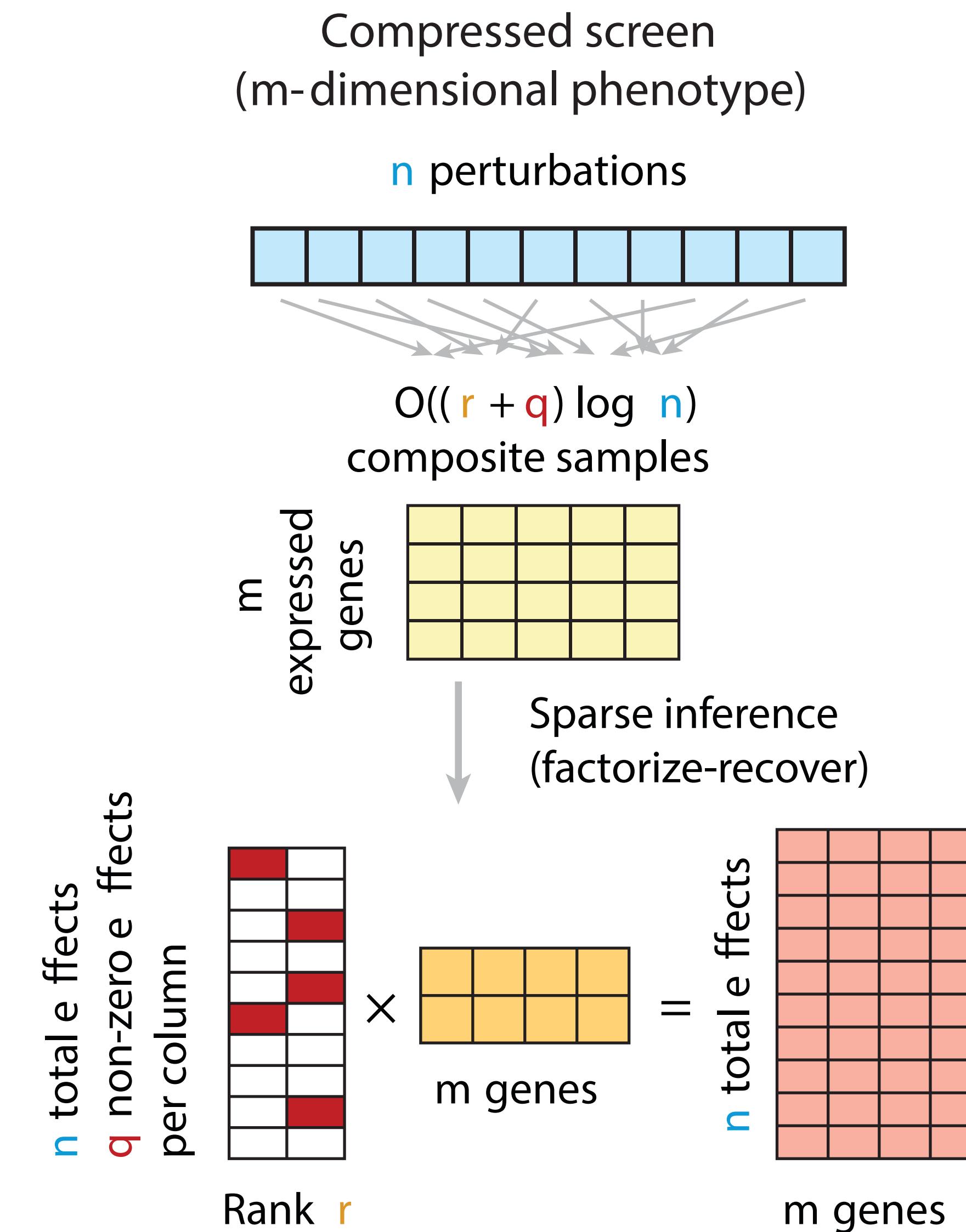
Published online: 23 October 2023

Douglas Yao<sup>1</sup>, Loic Binan<sup>2</sup>, Jon Bezney<sup>2,13</sup>, Brooke Simonton<sup>2</sup>, Jahanara Freedman<sup>2</sup>, Chris J. Frangieh<sup>2,3</sup>, Kushal Dey<sup>14</sup>, Kathryn Geiger-Schuller<sup>5</sup>, Basak Eraslan<sup>5</sup>, Alexander Gusev<sup>12,6,7,16</sup>, Aviv Regev<sup>2,14,15</sup> & Brian Cleary<sup>8,9,10,11,12,16</sup>✉

# Multiple CRISPR screening can be combined

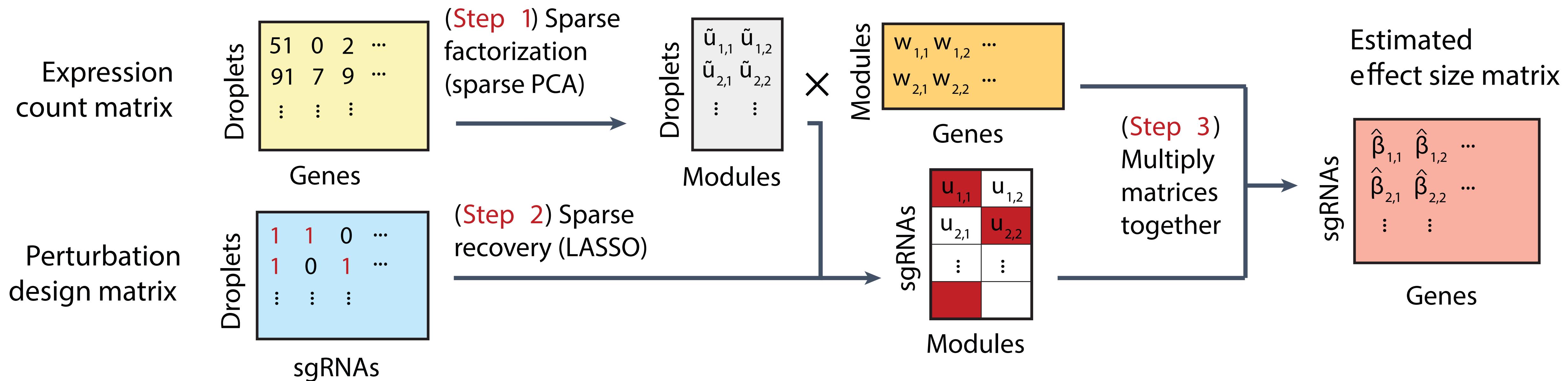


vs.



# Compressive sensing (sparse modelling) improves CRISPR experimental design

## Inference with FR-Perturb



# Today's lecture: Multiomics data integration

- **Why do we do multiomics data integration?**
  - view #1: borrowing information across modalities
  - view #2: efforts to provide mechanistic explanations
- **Global, unsupervised multiomics data integration**
  - Multiomics Factorization (and variants)
  - Network-based data integration
- **Local, linking between layers to understand mechanisms**
  - Deep dive into mechanisms of gene regulatory mechanisms