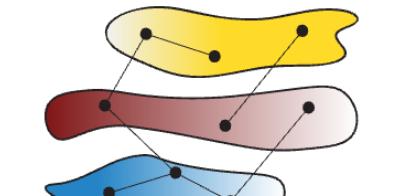




UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



CENTRE FOR
ORGANISMAL STUDIES



VELTEN GROUP



e l l i s
European Laboratory for Learning and Intelligent Systems

PROBABILISTIC FACTOR ANALYSIS FOR PERTURBATION RESPONSES IN SINGLE-CELL DATA

Francisca Gaspar Vieira
“Autumn School for Single Cell-ers”
October 2025



ABOUT ME

-
-
-
-
-



Integrated Master's in Biological Engineering, IST Lisboa

KU LEUVEN

Phage Biocontrol Research Project, KU Leuven



Master's in Computational Biology and Bioinformatics

NOVAMATH

CENTER FOR MATHEMATICS
AND APPLICATIONS

Researcher in projects: MONET, AI4MUFF and VOOmics



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386

EMBL-EBI



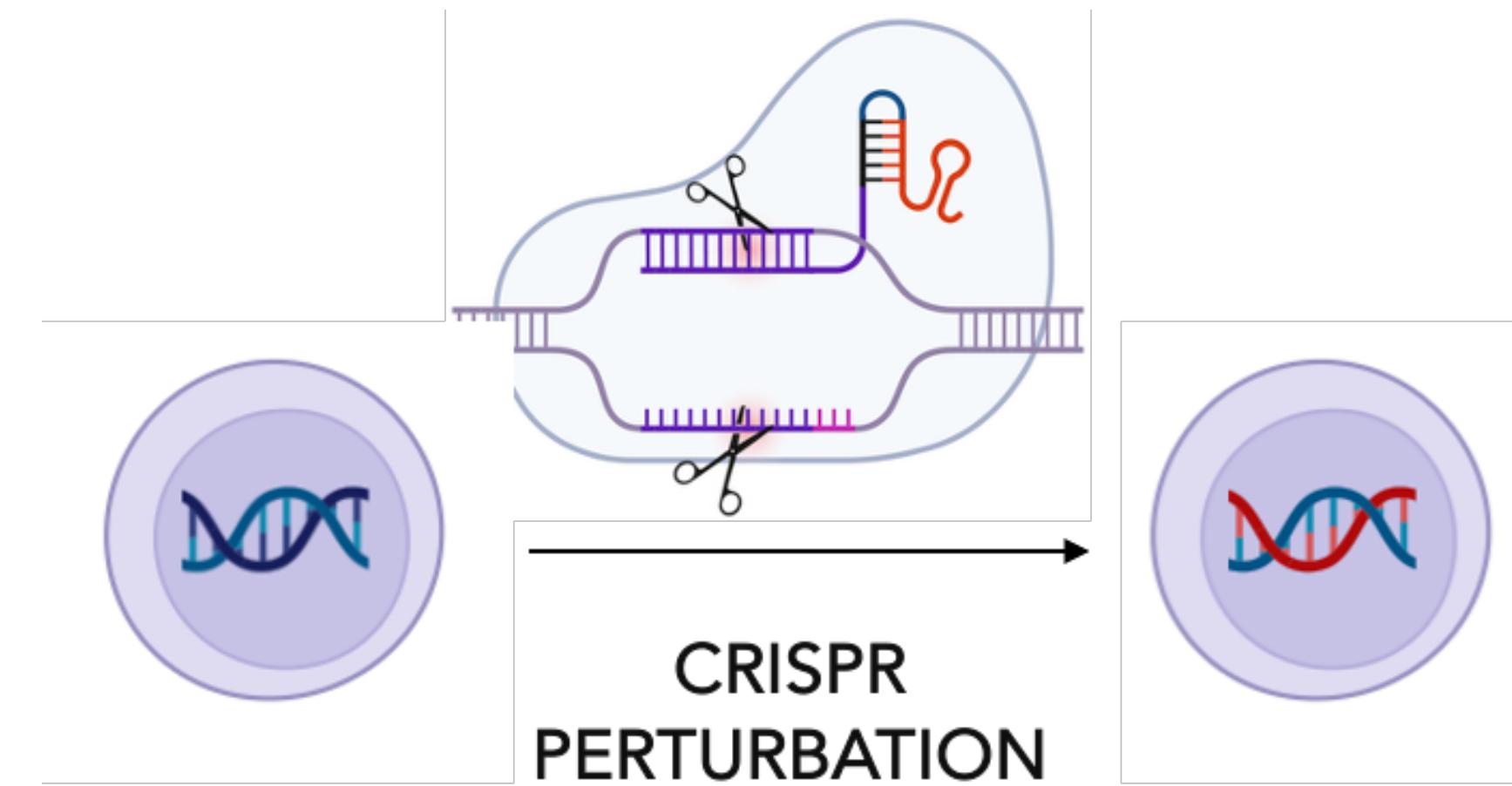
ELLIS PhD student, Heidelberg University + EMBL-EBI

Supervised by Dr. Britta Velten and Dr. Julio Saez-Rodriguez

1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

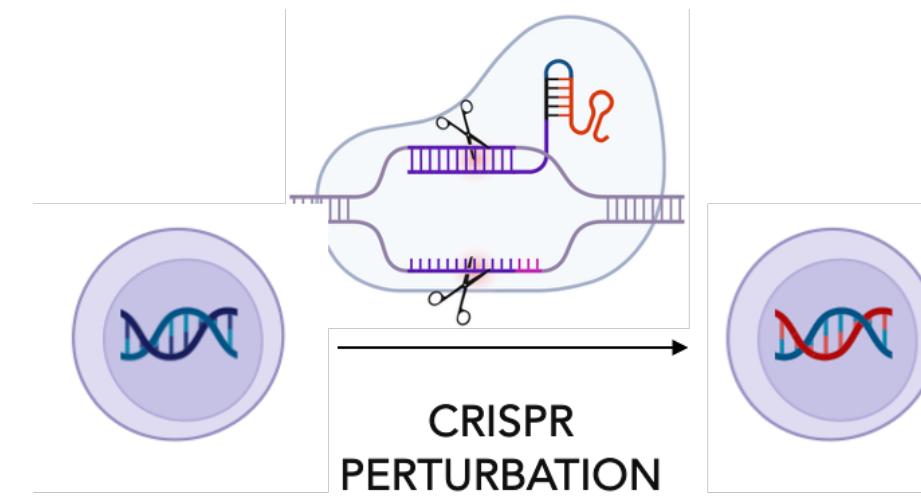
PERTURBATION RESPONSES IN SINGLE-CELL DATA

What does a certain gene actually do?

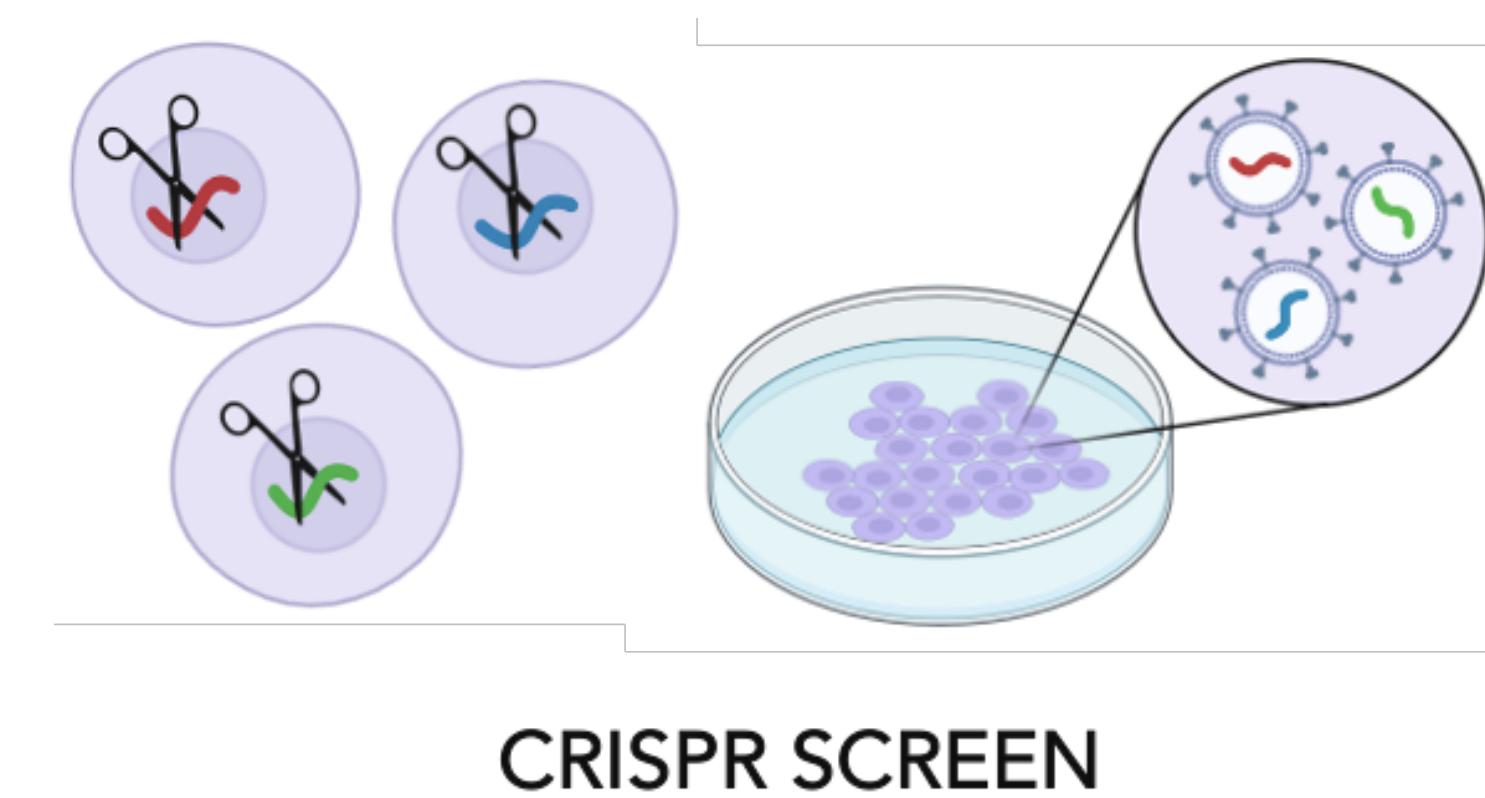


PERTURBATION RESPONSES IN SINGLE-CELL DATA

What does a certain gene actually do?

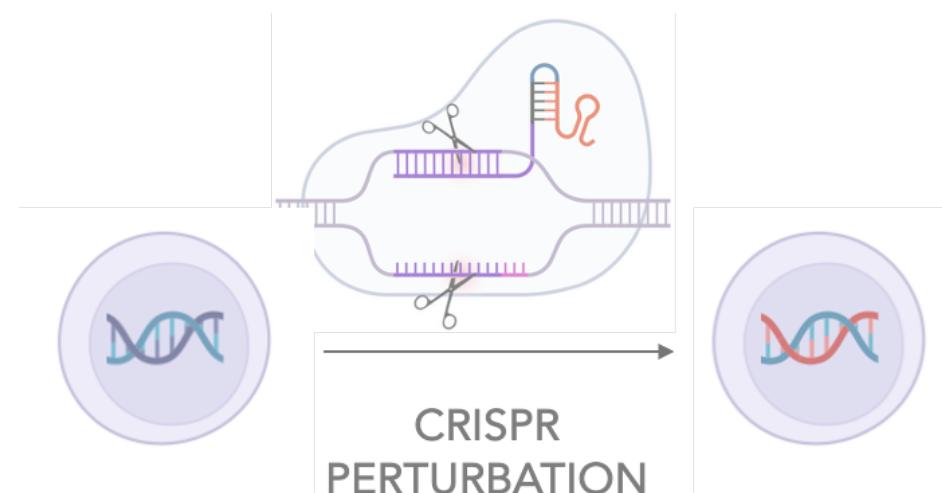


What about other genes?

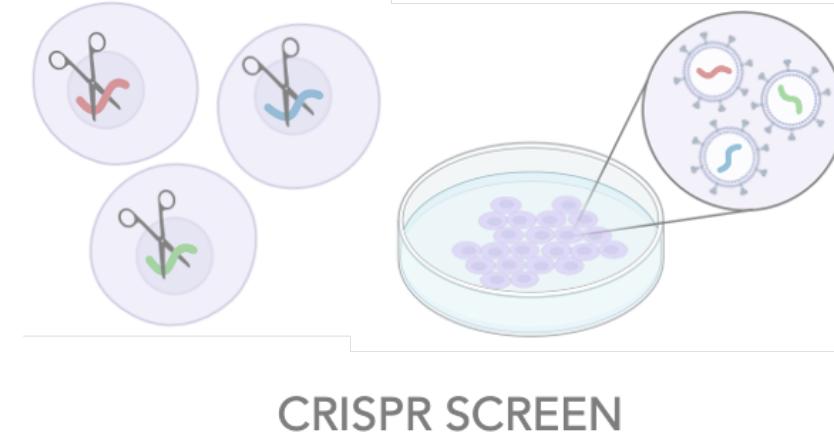


PERTURBATION RESPONSES IN SINGLE-CELL DATA

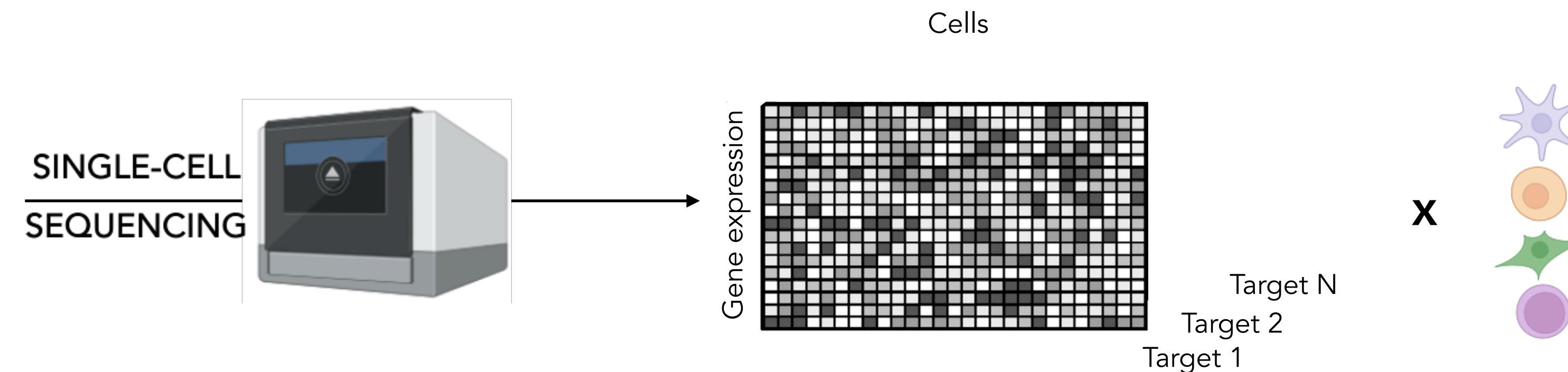
What does a certain gene actually do?



What about other genes?

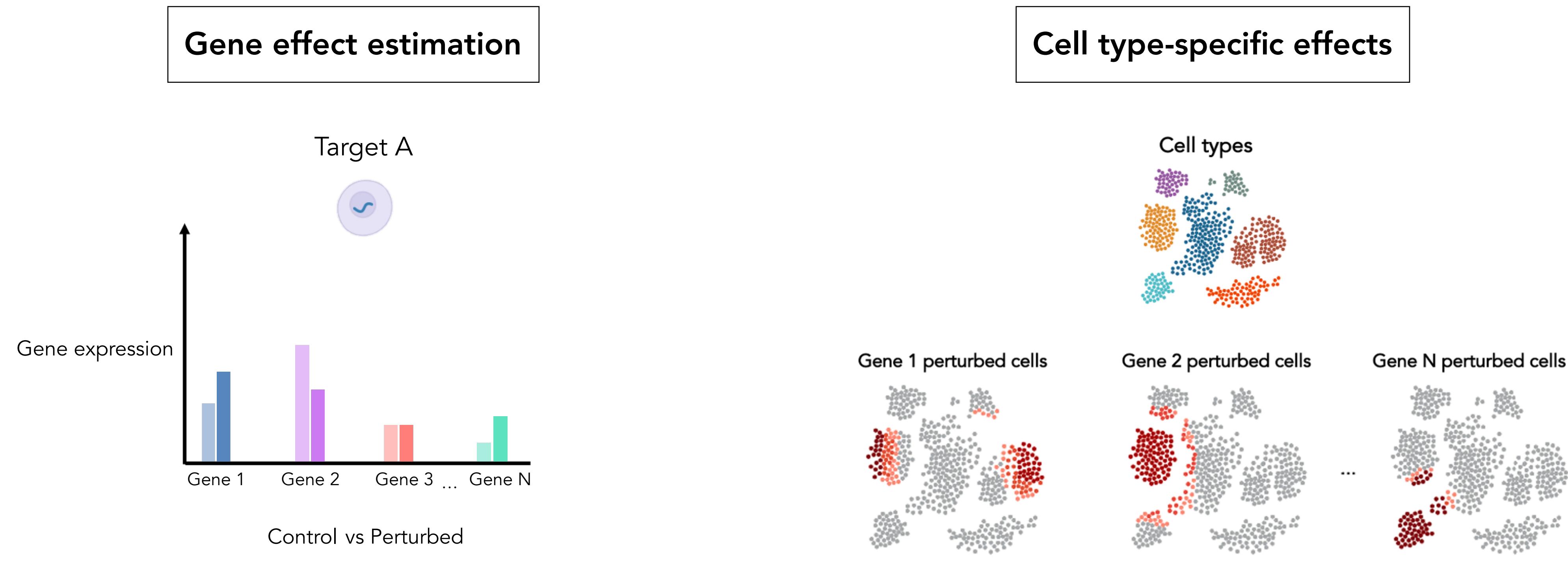


How does knocking down **each gene** reshape the transcriptional landscape of **each cell**?



Do different cell types/lines/signaling contexts respond differently to the same perturbation - how can we detect these differences?

PERTURBATION RESPONSES IN SINGLE-CELL DATA



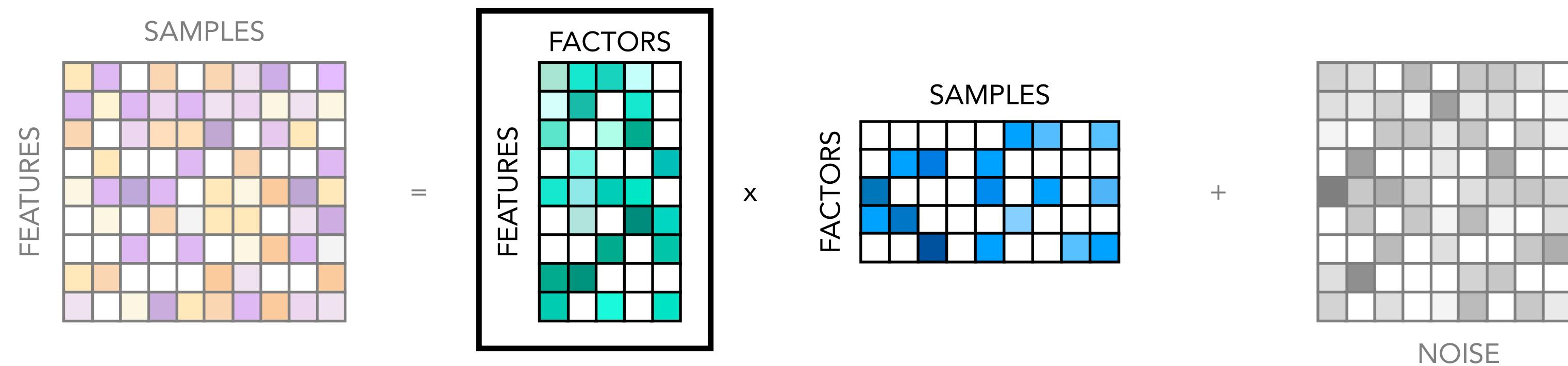
How to investigate shared effects between different biological contexts?

How to analyze all these results in a integrative way?

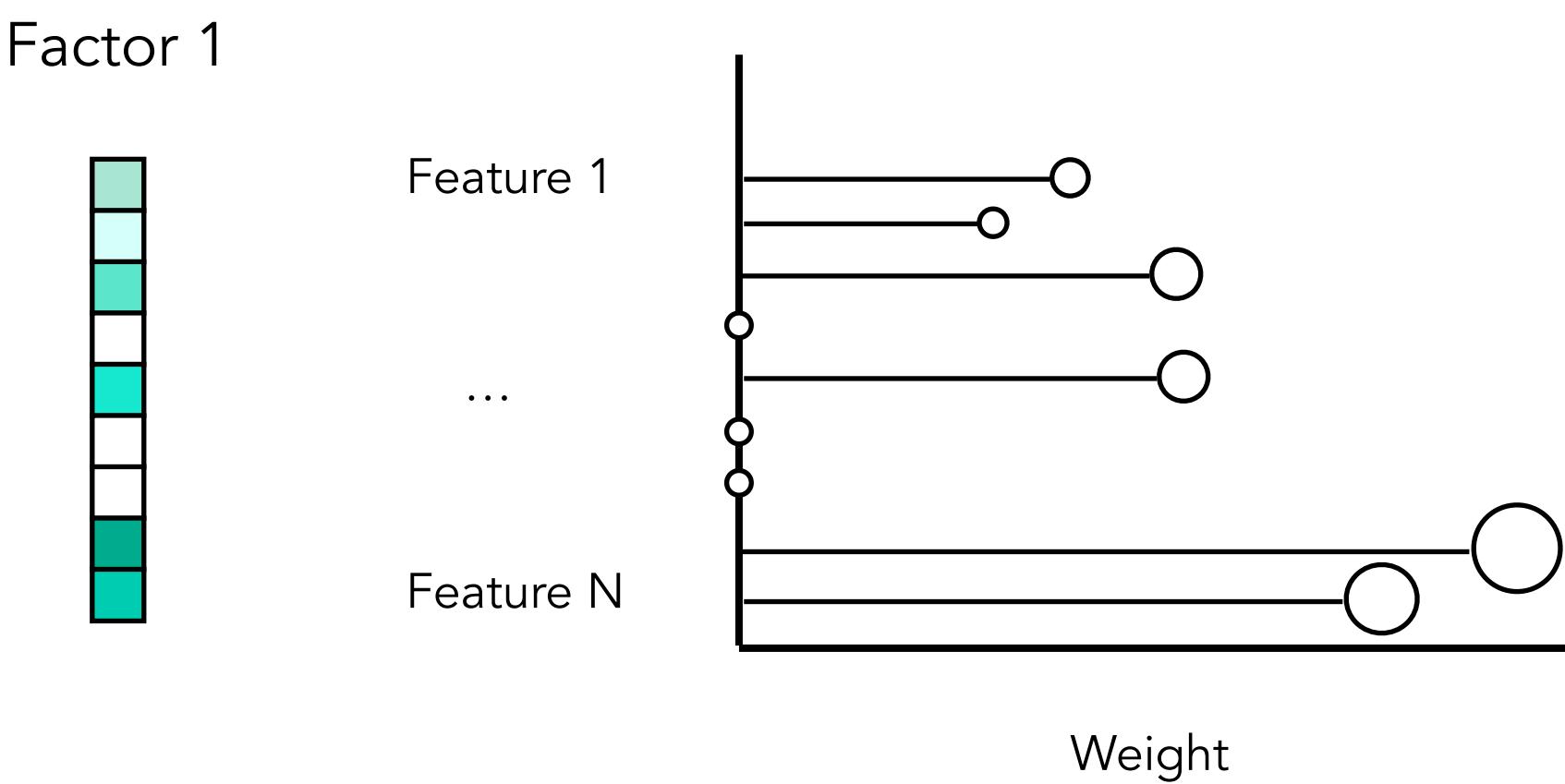
1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

APPROACH WITH FACTOR ANALYSIS

What is factor analysis?

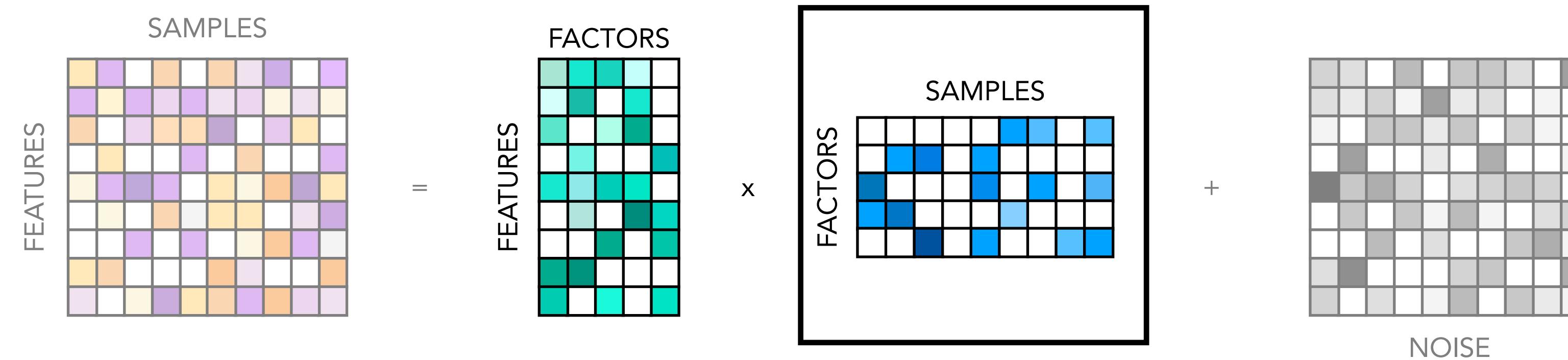


Why is it useful?

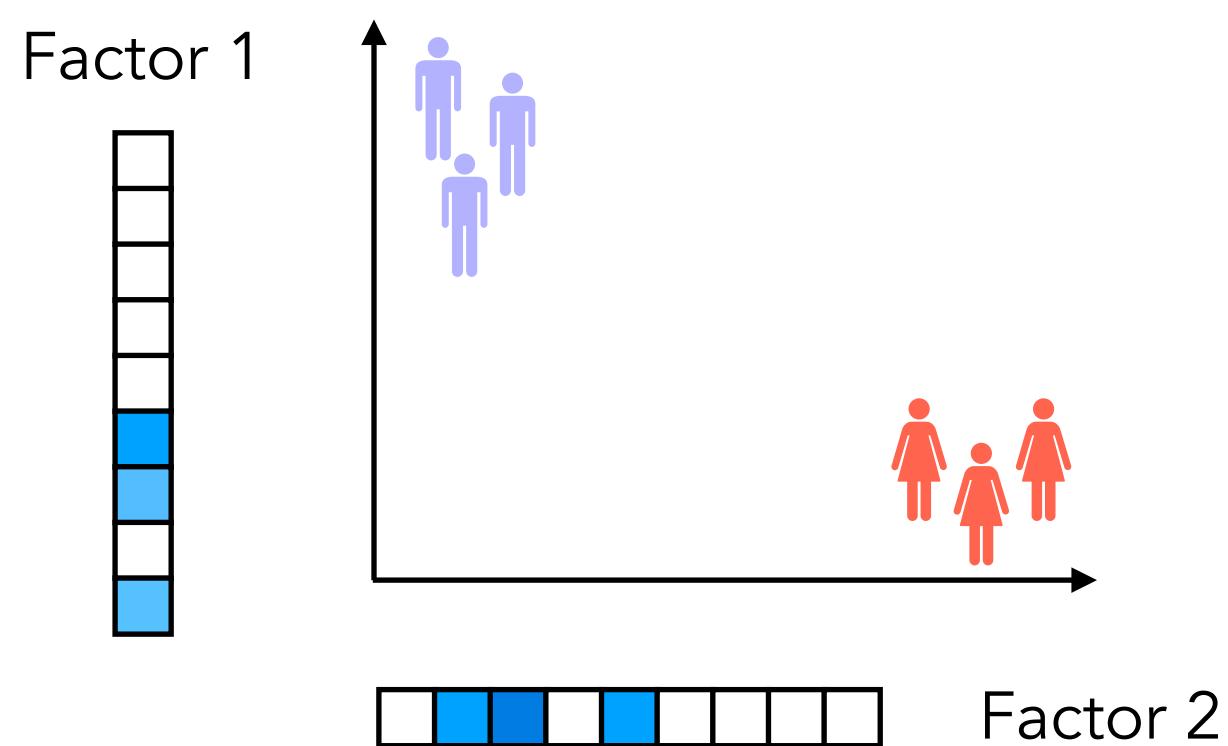


APPROACH WITH FACTOR ANALYSIS

What is factor analysis?

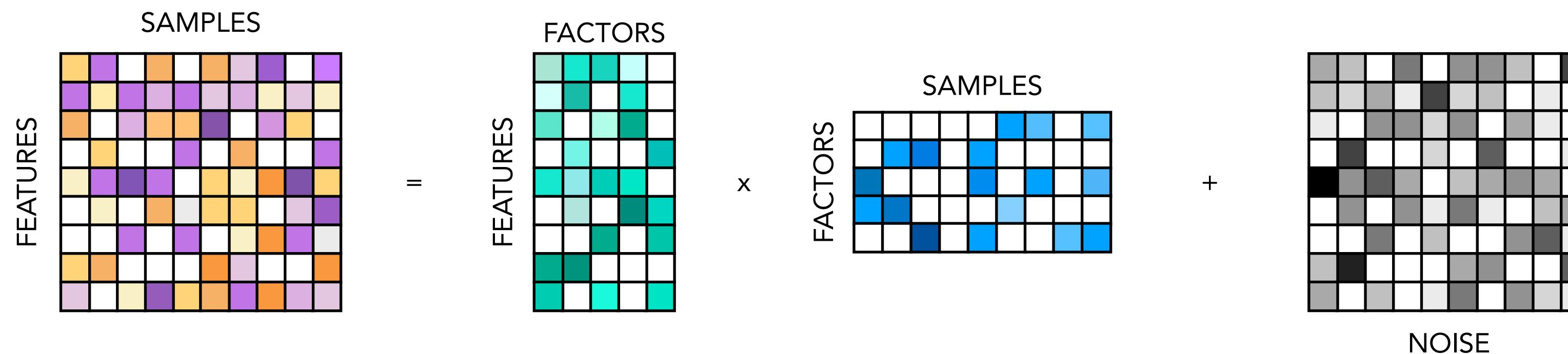


Why is it useful?



APPROACH WITH FACTOR ANALYSIS

What is behind probabilistic factor analysis?



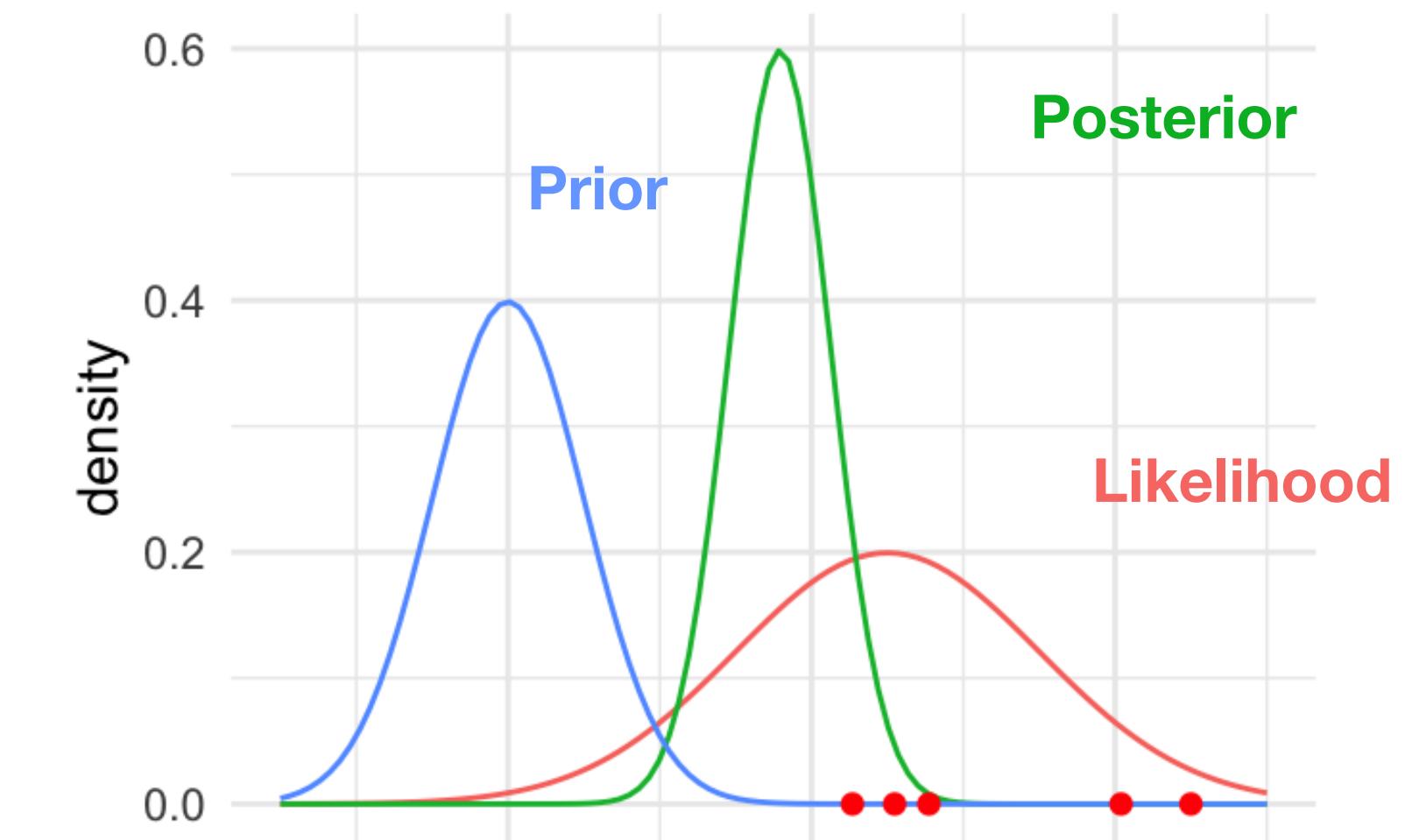
Probabilistic (Bayesian) models combine

- a *likelihood function* (statistical model for the observed data)
- a *prior* (distribution of unobserved components)

After seeing the data, we want to update our prior beliefs.

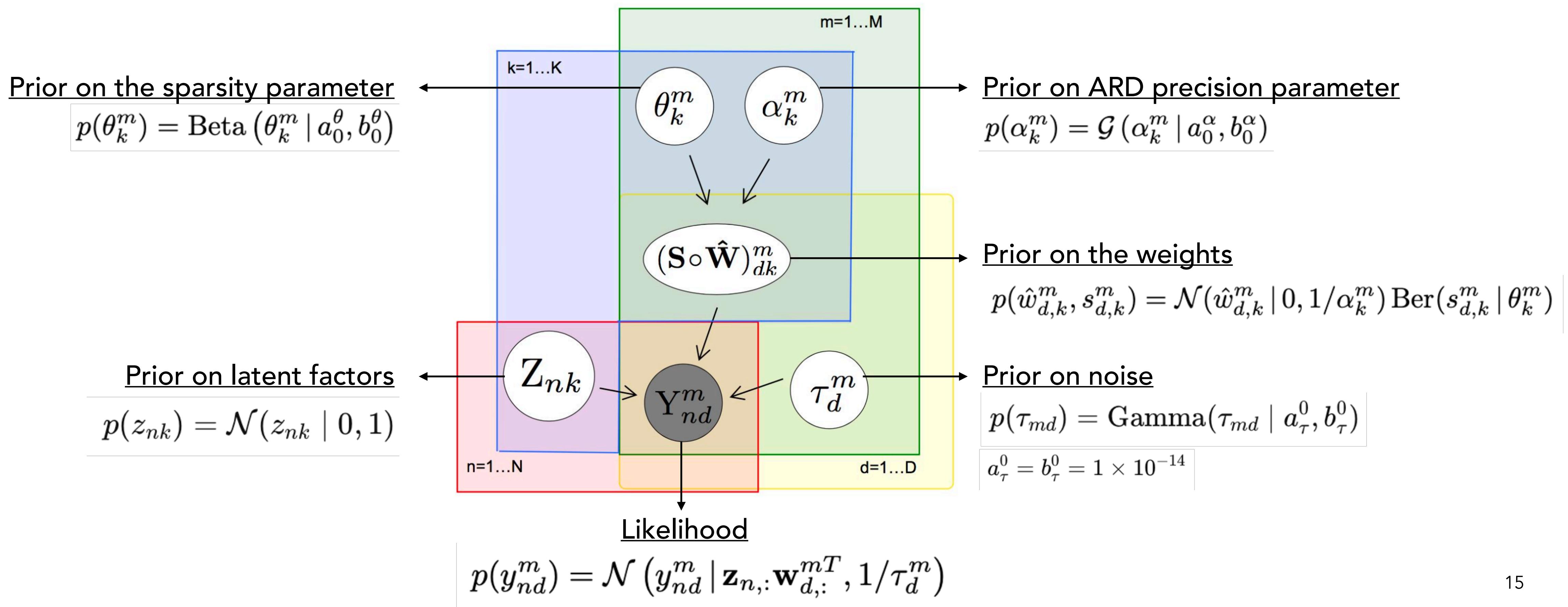
This updated belief is the **posterior**.

This can provide **regularization** & incorporation of **prior information** (prior) and an explicit model of **uncertainties**.



APPROACH WITH FACTOR ANALYSIS

$$\mathbf{Y} = \mathbf{W} \times \mathbf{Z} + \boldsymbol{\tau}$$



APPROACH WITH FACTOR ANALYSIS

Finding the posterior: An optimization problem

PROBLEM

We want to know the hidden structure (\mathbf{X}) behind our observed data (\mathbf{Y}): meaning to calculate the posterior:

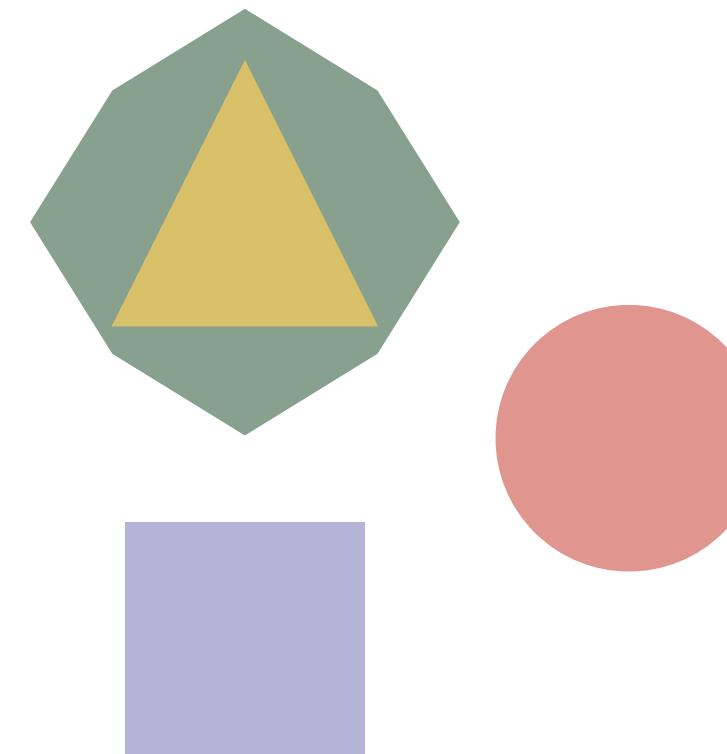
$$p(\mathbf{X} | \mathbf{Y}) = \frac{p(\mathbf{X}, \mathbf{Y})}{p(\mathbf{Y})} = \frac{p(\mathbf{X}, \mathbf{Y})}{\int p(\mathbf{X}) p(\mathbf{Y} | \mathbf{X}) d\mathbf{X}} \text{ !}$$

However, for most interesting models, the denominator becomes **intractable**.

CORE IDEA

Approximate the true posterior with a simpler and tractable distribution $q(\mathbf{X})$: turn inference into optimization.

Imagine trying to fit a complex shape with a simple one



APPROACH WITH FACTOR ANALYSIS

Finding the posterior: An optimization problem

PROBLEM

We want to know the hidden structure (\mathbf{X}) behind our observed data (\mathbf{Y}): meaning to calculate the posterior:

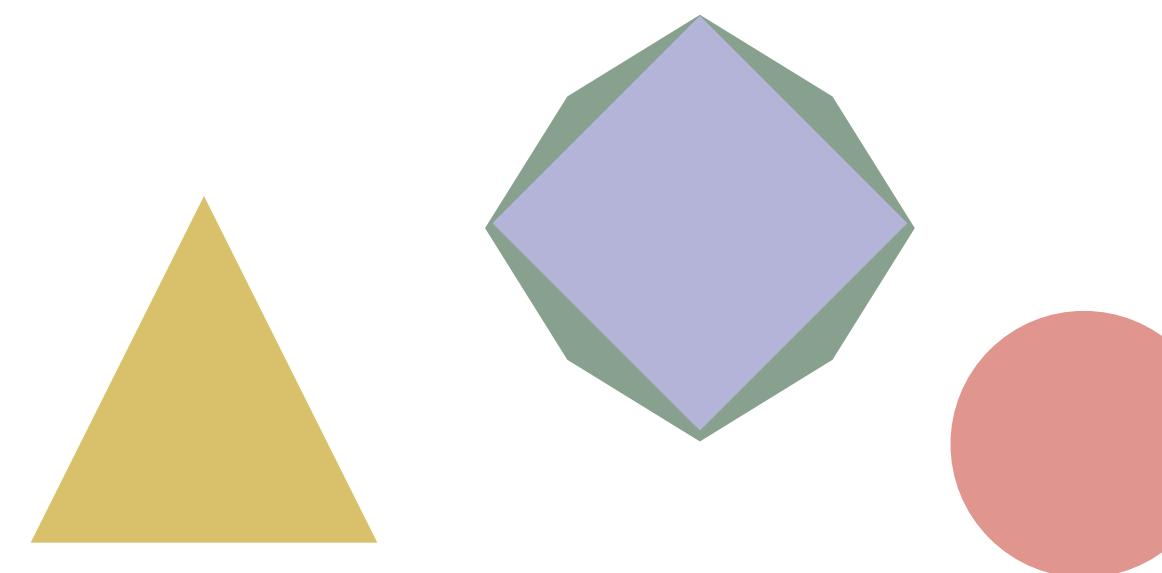
$$p(\mathbf{X} | \mathbf{Y}) = \frac{p(\mathbf{X}, \mathbf{Y})}{p(\mathbf{Y})} = \frac{p(\mathbf{X}, \mathbf{Y})}{\int p(\mathbf{X}) p(\mathbf{Y} | \mathbf{X}) d\mathbf{X}} \text{ !}$$

However, for most interesting models, the denominator becomes **intractable**.

CORE IDEA

Approximate the true posterior with a simpler and tractable distribution $q(\mathbf{X})$: turn inference into optimization.

Imagine trying to fit a complex shape with a simple one



APPROACH WITH FACTOR ANALYSIS

Finding the posterior: An optimization problem

PROBLEM

We want to know the hidden structure (\mathbf{X}) behind our observed data (\mathbf{Y}): meaning to calculate the posterior:

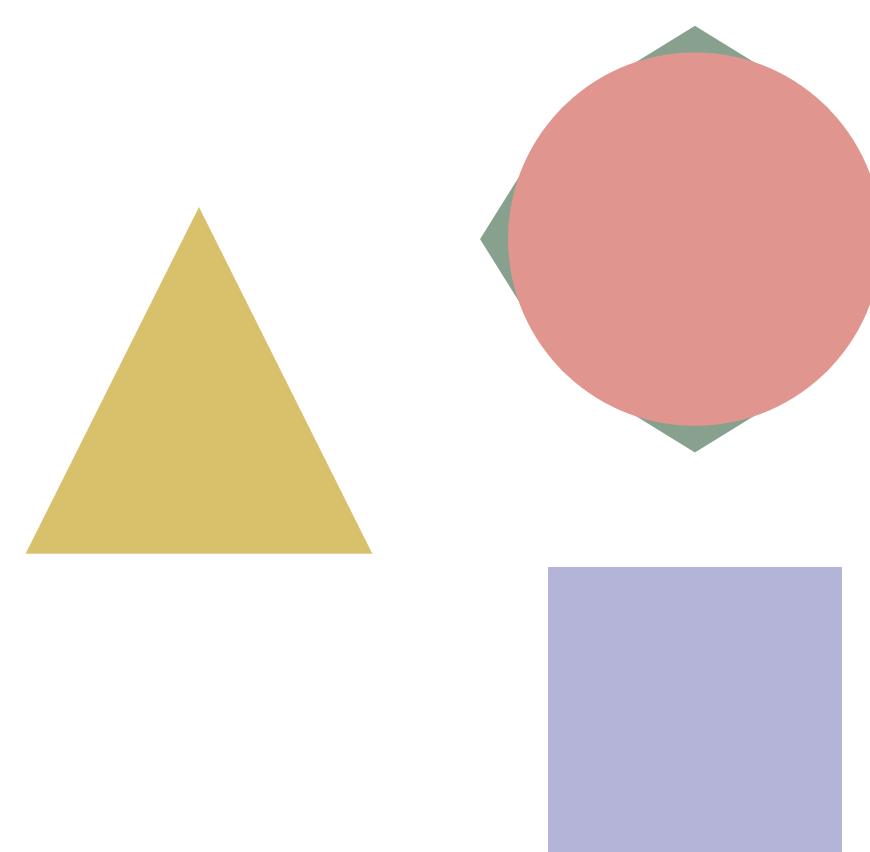
$$p(\mathbf{X} | \mathbf{Y}) = \frac{p(\mathbf{X}, \mathbf{Y})}{p(\mathbf{Y})} = \frac{p(\mathbf{X}, \mathbf{Y})}{\int p(\mathbf{X}) p(\mathbf{Y} | \mathbf{X}) d\mathbf{X}} \text{ !}$$

However, for most interesting models, the denominator becomes **intractable**.

CORE IDEA

Approximate the true posterior with a simpler and tractable distribution $q(\mathbf{X})$: turn inference into optimization.

Imagine trying to fit a complex shape with a simple one



APPROACH WITH FACTOR ANALYSIS

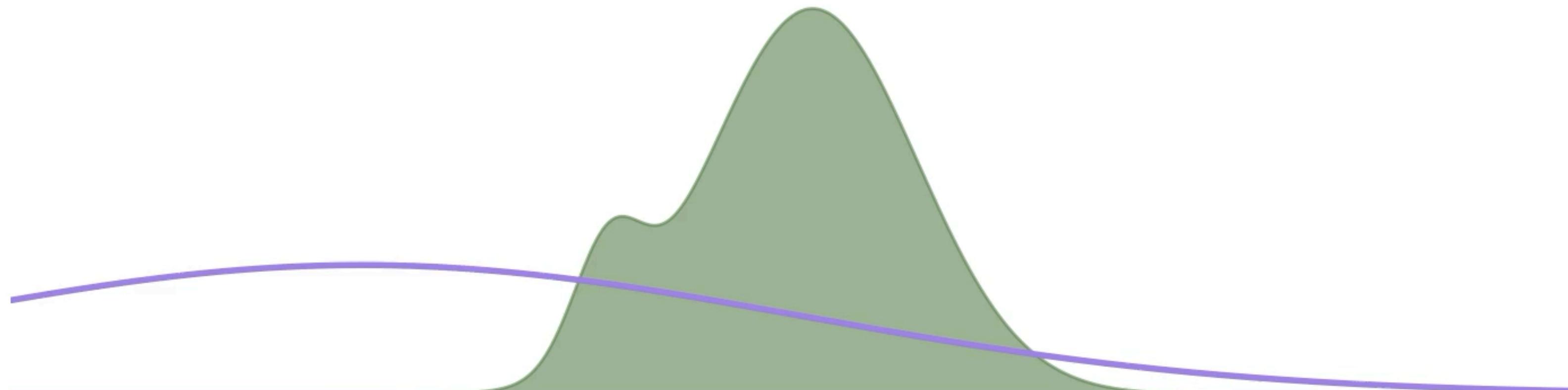
Finding the posterior: An optimization problem

1. Pick a flexible, tractable family of distributions (e.g., Gaussians)
2. Define objective/ loss function for how well $q(\mathbf{X})$ approximates $p(\mathbf{X} | \mathbf{Y})$
3. Optimize: Adjust parameters of $q(\mathbf{X})$ to minimize the mismatch

 $p(\mathbf{X} | \mathbf{Y})$ True Posterior
 $q(\mathbf{X})$ Variational Approx.

Mean (μ): -1.00
Std Dev (σ): 2.50
 $KL(p\|q)$: 62.327

Step: 1/60

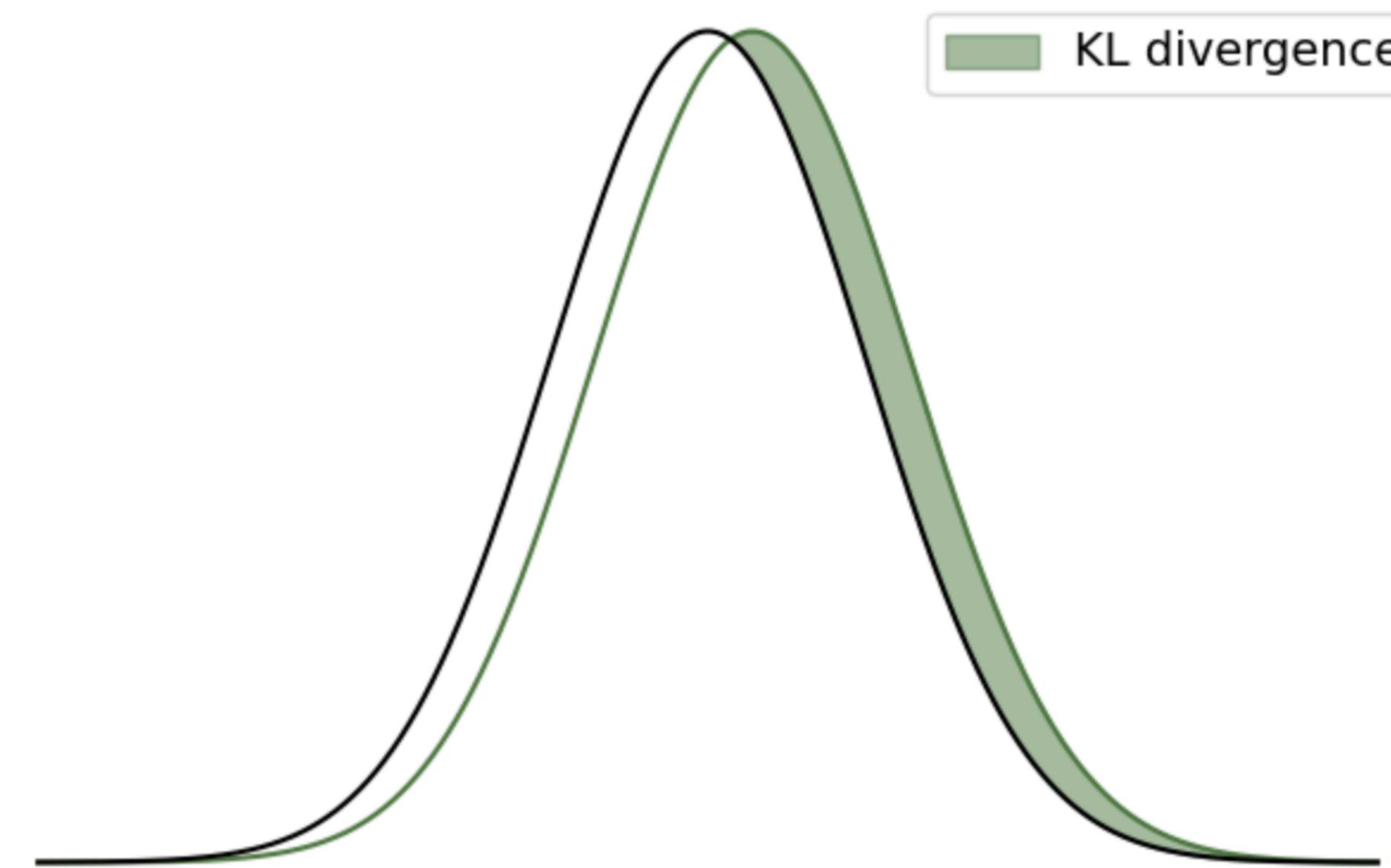


APPROACH WITH FACTOR ANALYSIS

Kullback–Leibler divergence

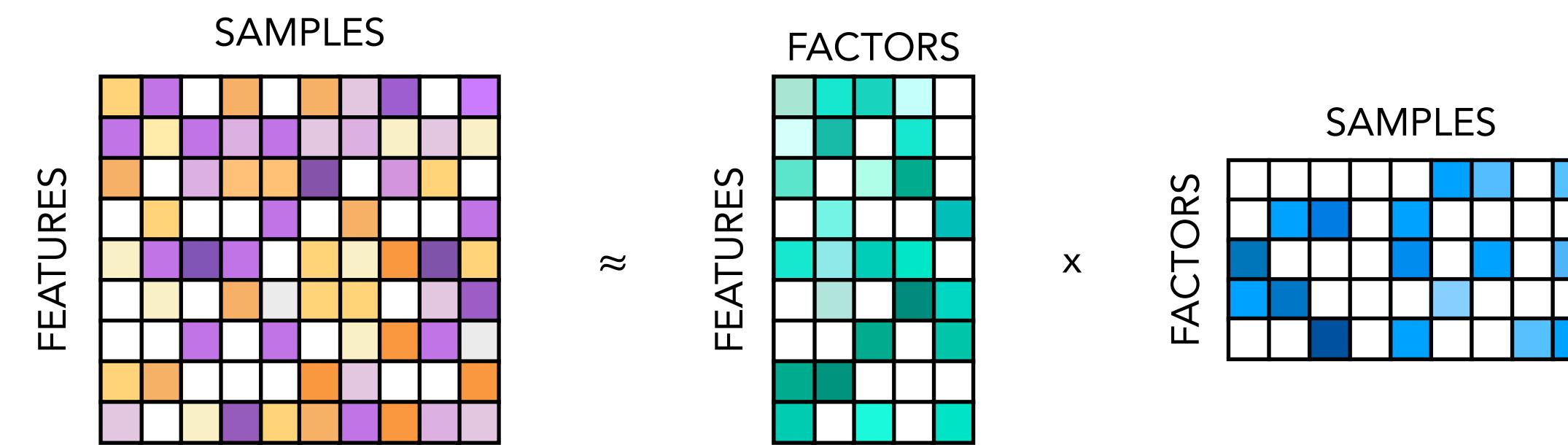
The Kullback-Leibler (KL) divergence is the standard distance measure choice when working with probability distributions:

$$\text{KL}(q(\mathbf{X}) \parallel p(\mathbf{X} \mid \mathbf{Y})) = \int q(\mathbf{X}) \log \frac{q(\mathbf{X})}{p(\mathbf{X} \mid \mathbf{Y})} d\mathbf{z}$$



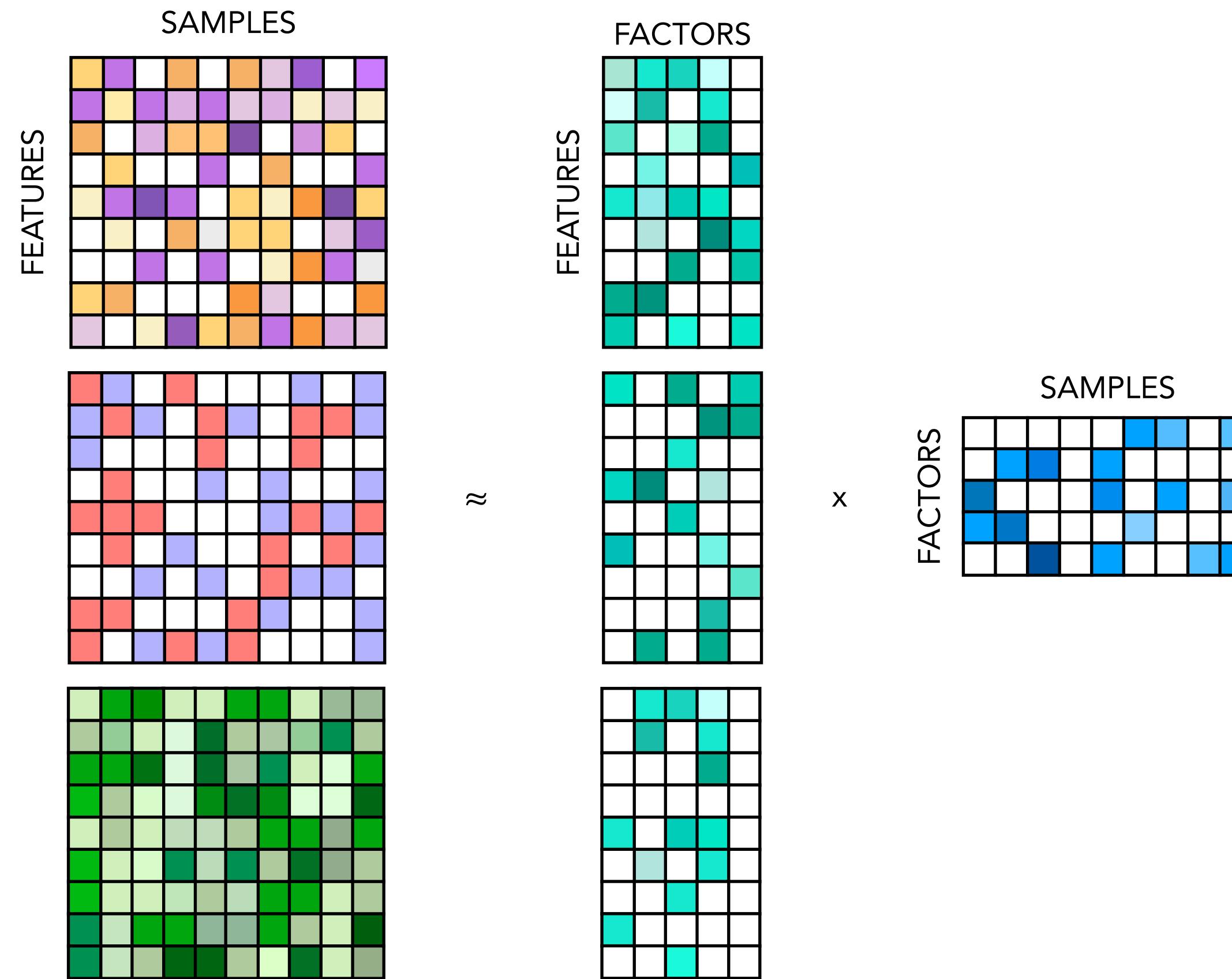
APPROACH WITH FACTOR ANALYSIS

What is factor analysis?



APPROACH WITH FACTOR ANALYSIS

What is multi view-factor analysis?



APPROACH WITH FACTOR ANALYSIS

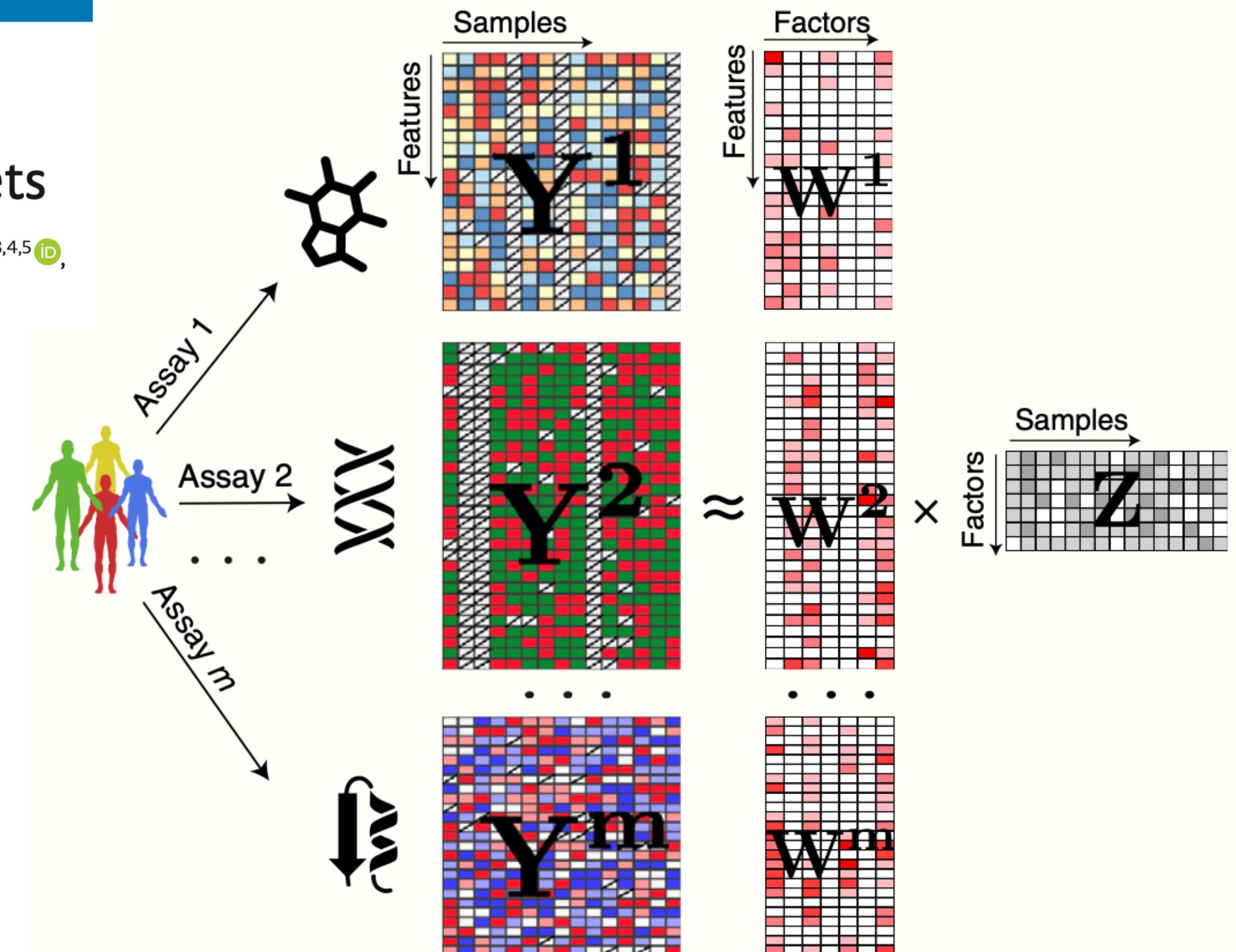
Method



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

Ricard Argelaguet^{1,†} , Britta Velten^{2,†} , Damien Arnol¹ , Sascha Dietrich³ , Thorsten Zenz^{3,4,5} , John C Marioni^{1,6,7} , Florian Buettner^{1,8,*} , Wolfgang Huber^{2,**} & Oliver Stegle^{1,2,***}

- MOFA infers an interpretable low-dimensional data representation in terms of (hidden) factors that capture major sources of variation across data modalities
- The inferred factor loadings can be sparse, thereby facilitating the linkage between the factors and the most relevant molecular features.
- **MOFA disentangles to what extent each factor is unique to a single data modality or is manifested in multiple modalities**



APPROACH WITH FACTOR ANALYSIS



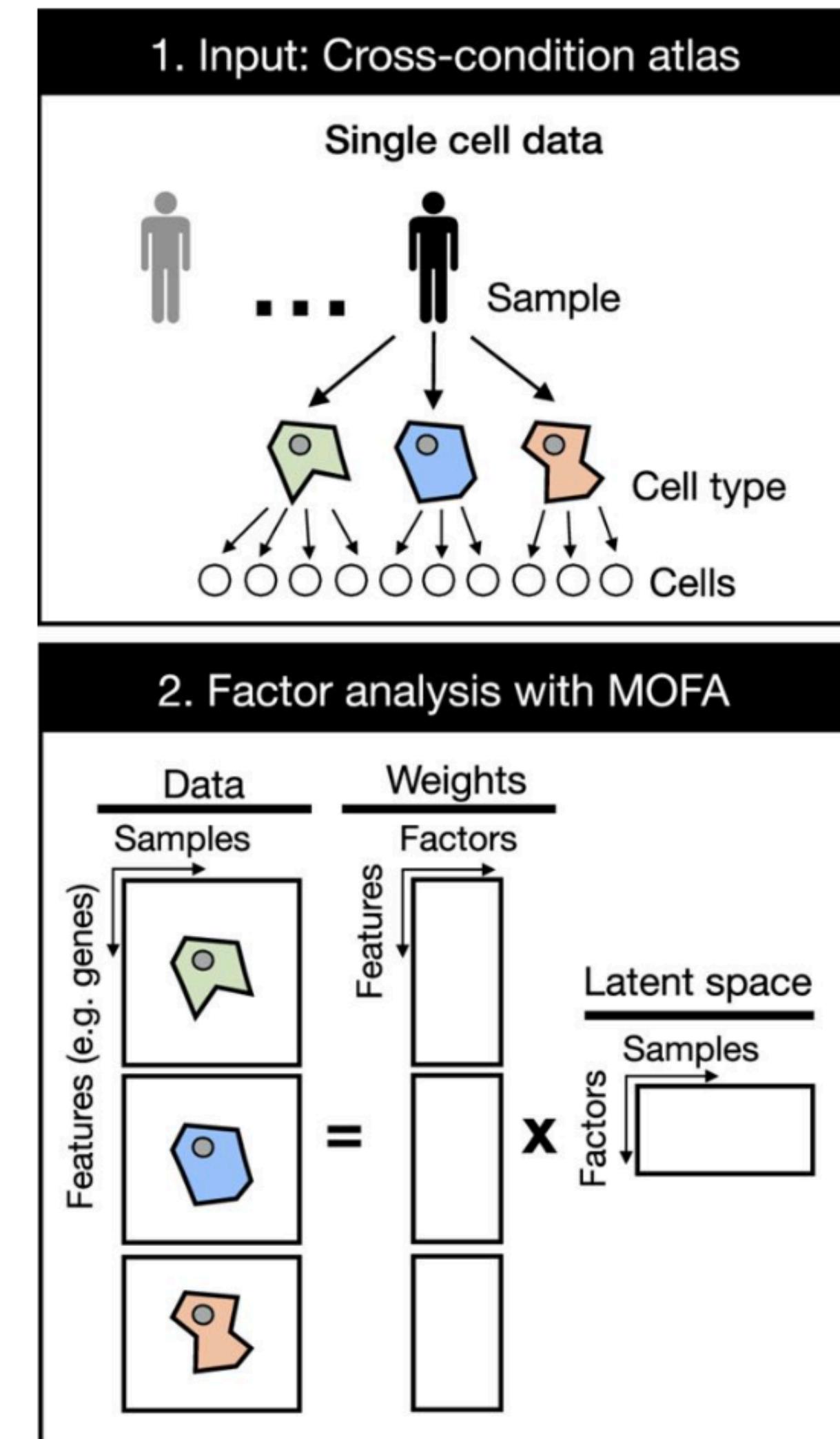
RESEARCH ARTICLE



Multicellular factor analysis of single-cell data for a tissue-centric understanding of disease

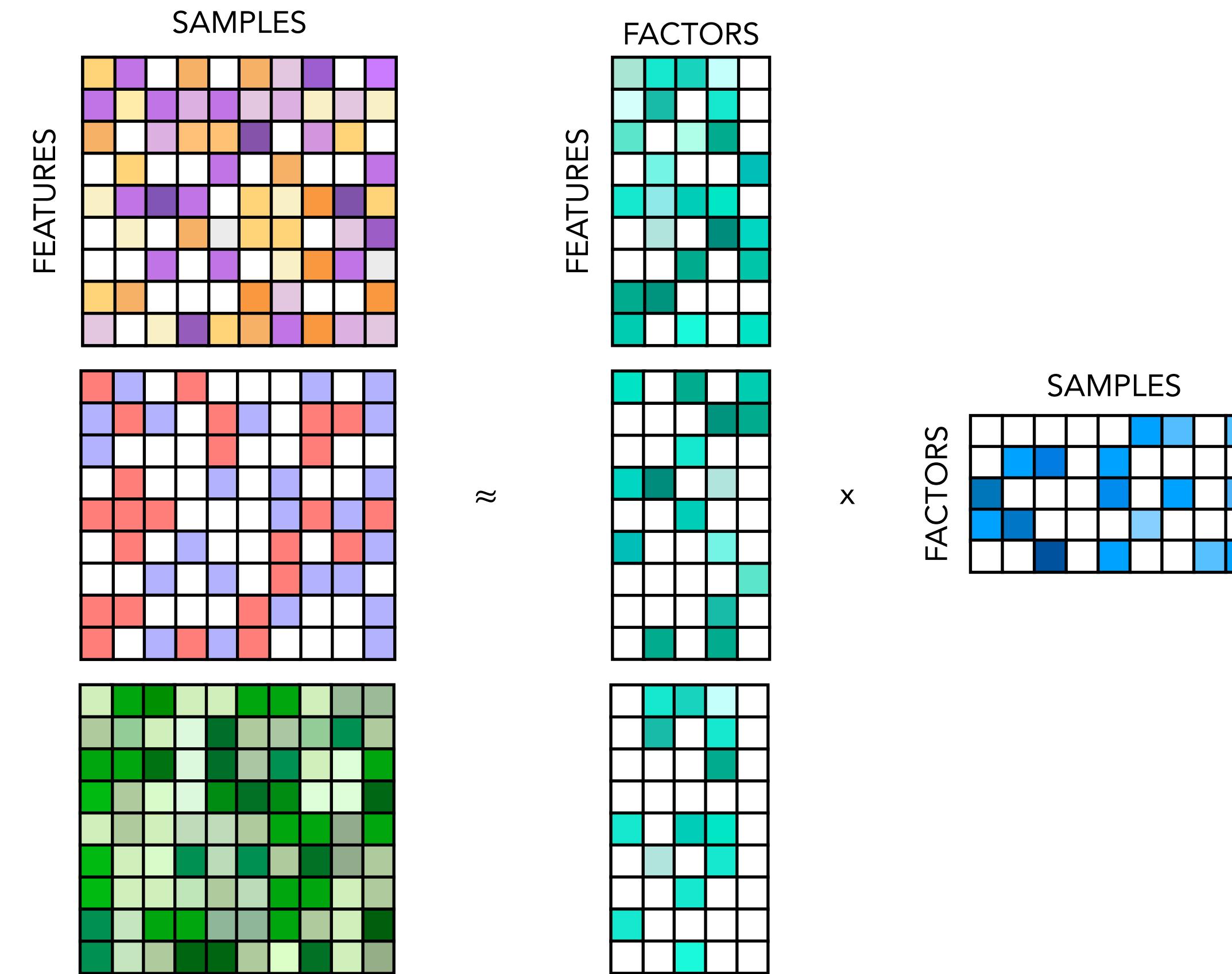
Ricardo Omar Ramirez Flores^{1*}, Jan David Lanzer¹, Daniel Dimitrov¹,
Britta Velten², Julio Saez-Rodriguez^{1*}

- Repurposes multi-omics factor analysis (MOFA) to simultaneously **decompose the variability of multiple cell types** and create a latent space that recovers multicellular transcriptional programs.
- The variables that form this latent space can be **interpreted as coordinated transcriptional changes** occurring in multiple cells - multicellular programs, providing a tissue-centric understanding of the analyzed sample.



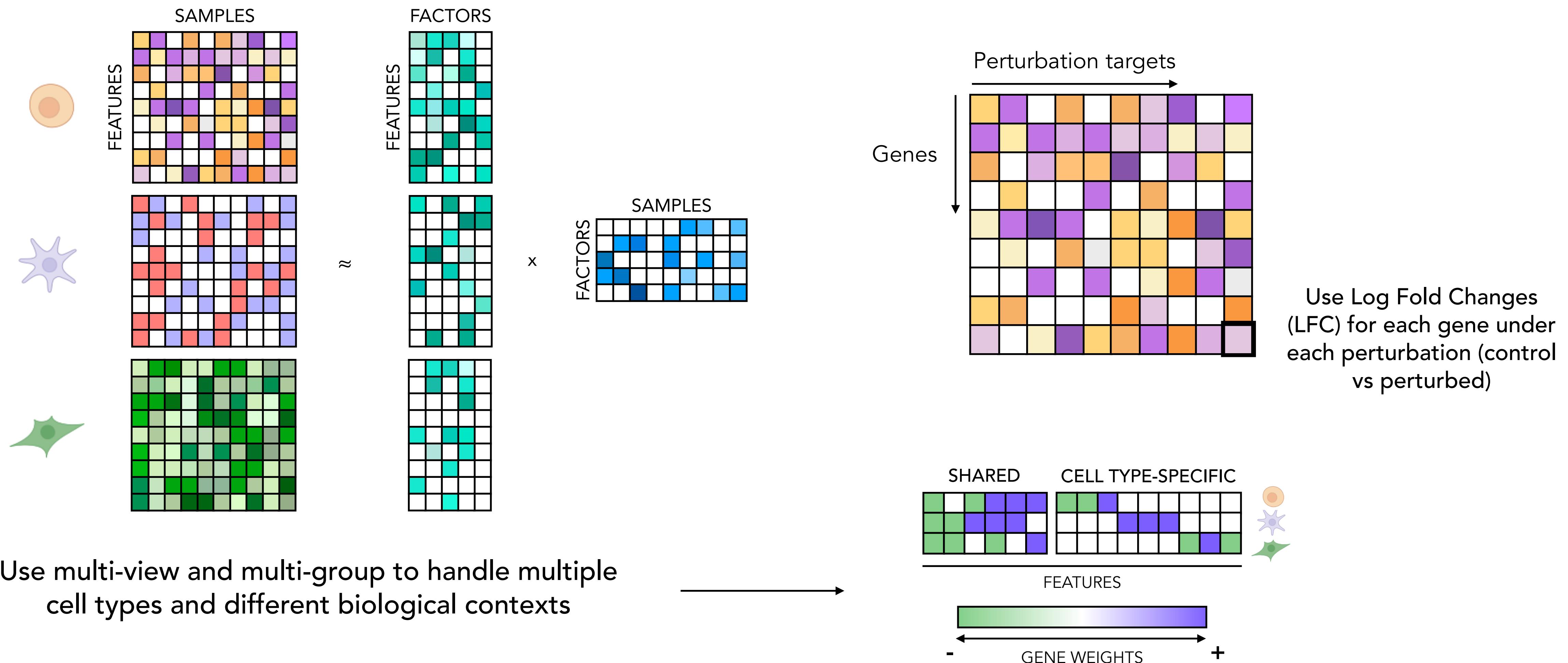
APPROACH WITH FACTOR ANALYSIS

How to use/repurpose it for perturbation responses?



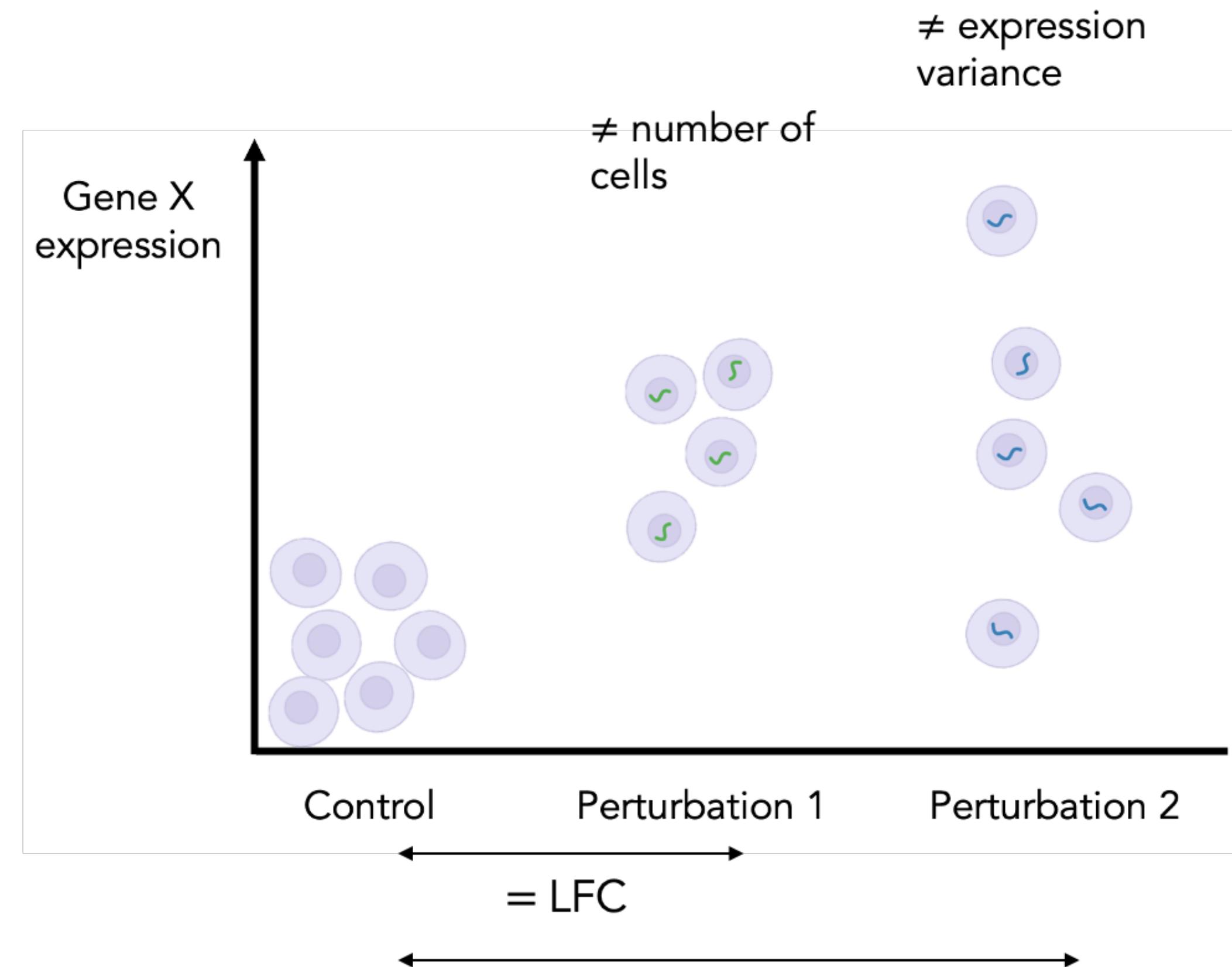
APPROACH WITH FACTOR ANALYSIS

How to use/repurpose it for perturbation responses?



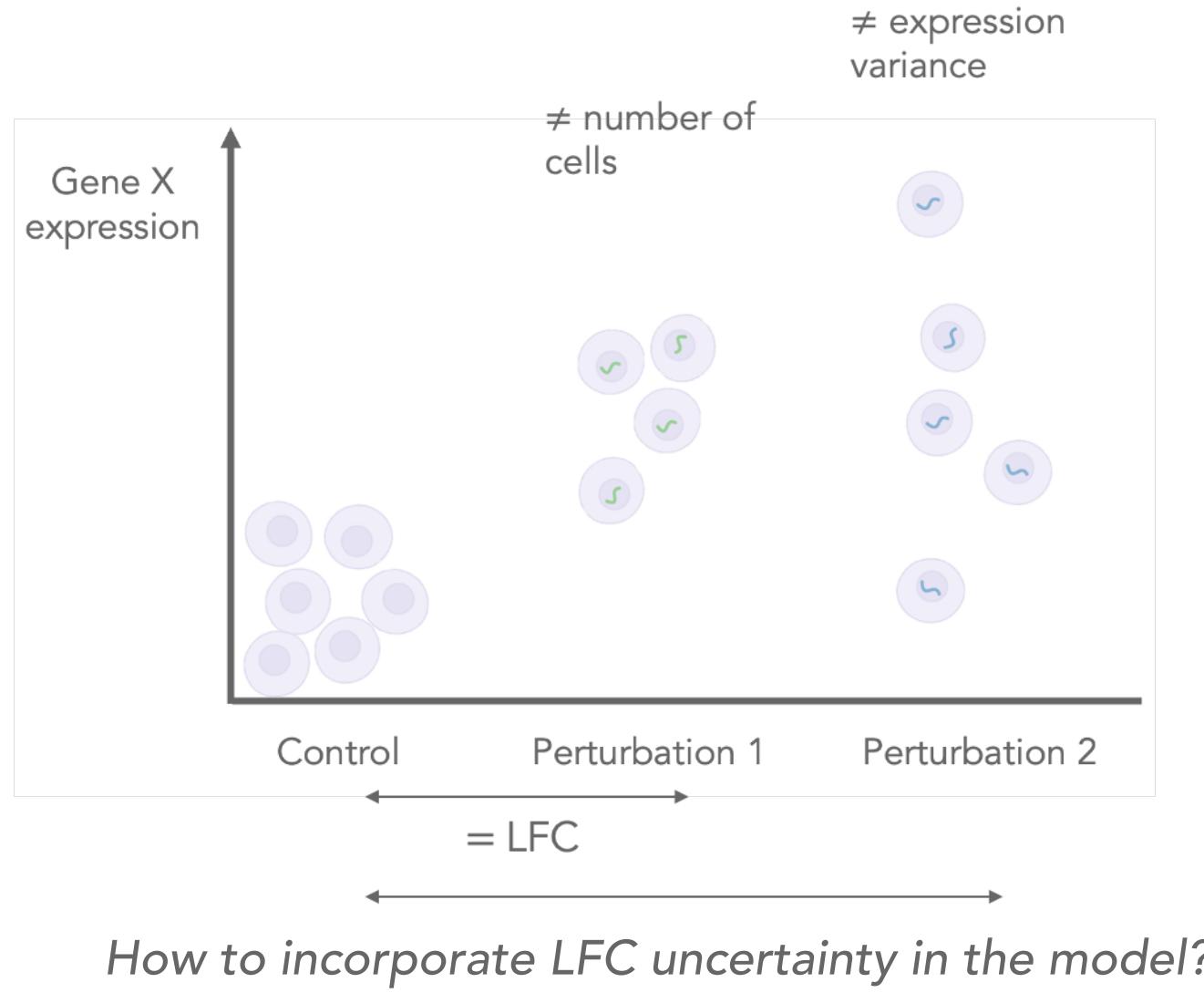
1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

INTEGRATING UNCERTAINTY



How to incorporate LFC uncertainty in the model?

INTEGRATING UNCERTAINTY



DATA AND ASSUMPTION

For each gene g and perturbation p :

We do not observe b_{gp} directly.
Instead, we observe a noisy
estimate obtained from some
inference method (e.g. DESeq2,
GLM, Seurat,...)

$$\hat{b}_{gp} \sim \mathbb{F}_{gp}, \quad \mathbb{E}[\hat{b}_{gp}] = b_{gp}$$

\hat{b}_{gp} : Estimated effect (LFC)

b_{gp} : True, unknown effect

\mathbb{F}_{gp} : Distribution with known info:

$$\left\{ \begin{array}{l} \text{Var}(\hat{b}_{gp}) = \sigma_{gp}^2 \\ \pi_{gp} = \mathbb{P}(|\hat{B}_{gp}| \geq |\hat{b}_{gp}| \mid b_{gp} = 0) \end{array} \right.$$

We assume that the matrix of perturbation effects has a **low-rank structure**:

$$b_p = W z_p + \varepsilon_p, \quad \varepsilon_p \sim \mathcal{N}(0, \Psi)$$

$$b_p \sim \mathcal{N}(0, W^T W + \Psi)$$

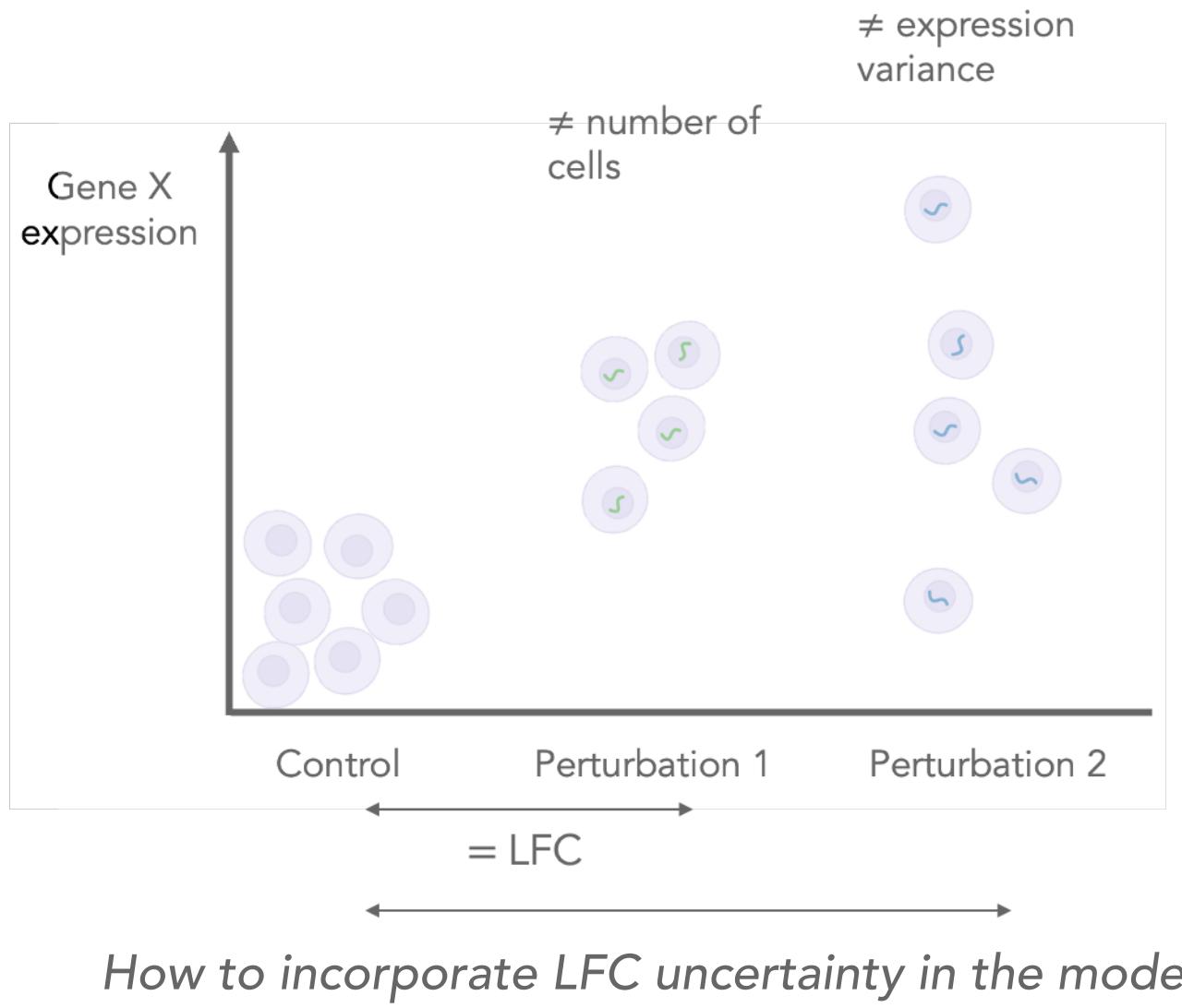
$$\Psi = \text{diag}(\phi_1^2, \dots, \phi_G^2)$$

$W \in \mathbb{R}^{G \times K}$: Gene loadings on latent factors (K)

$z_p \in \mathbb{R}^K$: Perturbation's latent factor

ε_p : Gene-specific residuals

INTEGRATING UNCERTAINTY



DATA AND ASSUMPTION

For each gene g and perturbation p : $\hat{b}_{gp} \sim \mathbb{F}_{gp}$, $\mathbb{E}[\hat{b}_{gp}] = b_{gp}$

\hat{b}_{gp} : Estimated effect (LFC)

b_{gp} : True, unknown effect

\mathbb{F}_{gp} : Distribution with known info: $\begin{cases} \text{Var}(\hat{b}_{gp}) = \sigma_{gp}^2 \\ \pi_{gp} = \mathbb{P}(|\hat{B}_{gp}| \geq |\hat{b}_{gp}| \mid b_{gp} = 0) \end{cases}$

We assume that the matrix of perturbation effects has a **low-rank structure**:

$$b_p = Wz_p + \varepsilon_p, \quad \varepsilon_p \sim \mathcal{N}(0, \Psi)$$

$$b_p \sim \mathcal{N}(0, W^T W + \Psi)$$

$$\Psi = \text{diag}(\phi_1^2, \dots, \phi_G^2)$$

$W \in \mathbb{R}^{G \times K}$: Gene loadings on latent factors (K)

$z_p \in \mathbb{R}^K$: Perturbation's latent factor

ε_p : Gene-specific residuals

APPROACHES

1. Direct naïve application

$$\hat{b}_{gp} = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

2. Variance weighting (z-score)

$$\frac{\hat{b}_{gp}}{\sigma_{gp}} = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

3. P-value weighting

$$\hat{b}_{gp} \cdot (-\log_{10} \pi_{gp}) = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

4. Decomposition of noise

$$\hat{b}_{gp} = \sum_k w_{gk} z_{kp} + \varepsilon_g + \zeta_{gp}$$

INTEGRATING UNCERTAINTY

APPROACHES

1. Direct naïve application

$$\hat{b}_{gp} = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

2. Variance weighting (z-score)

$$\frac{\hat{b}_{gp}}{\sigma_{gp}} = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

3. P-value weighting

$$\hat{b}_{gp} \cdot (-\log_{10} \pi_{gp}) = \sum_k w_{gk} z_{kp} + \varepsilon_g$$

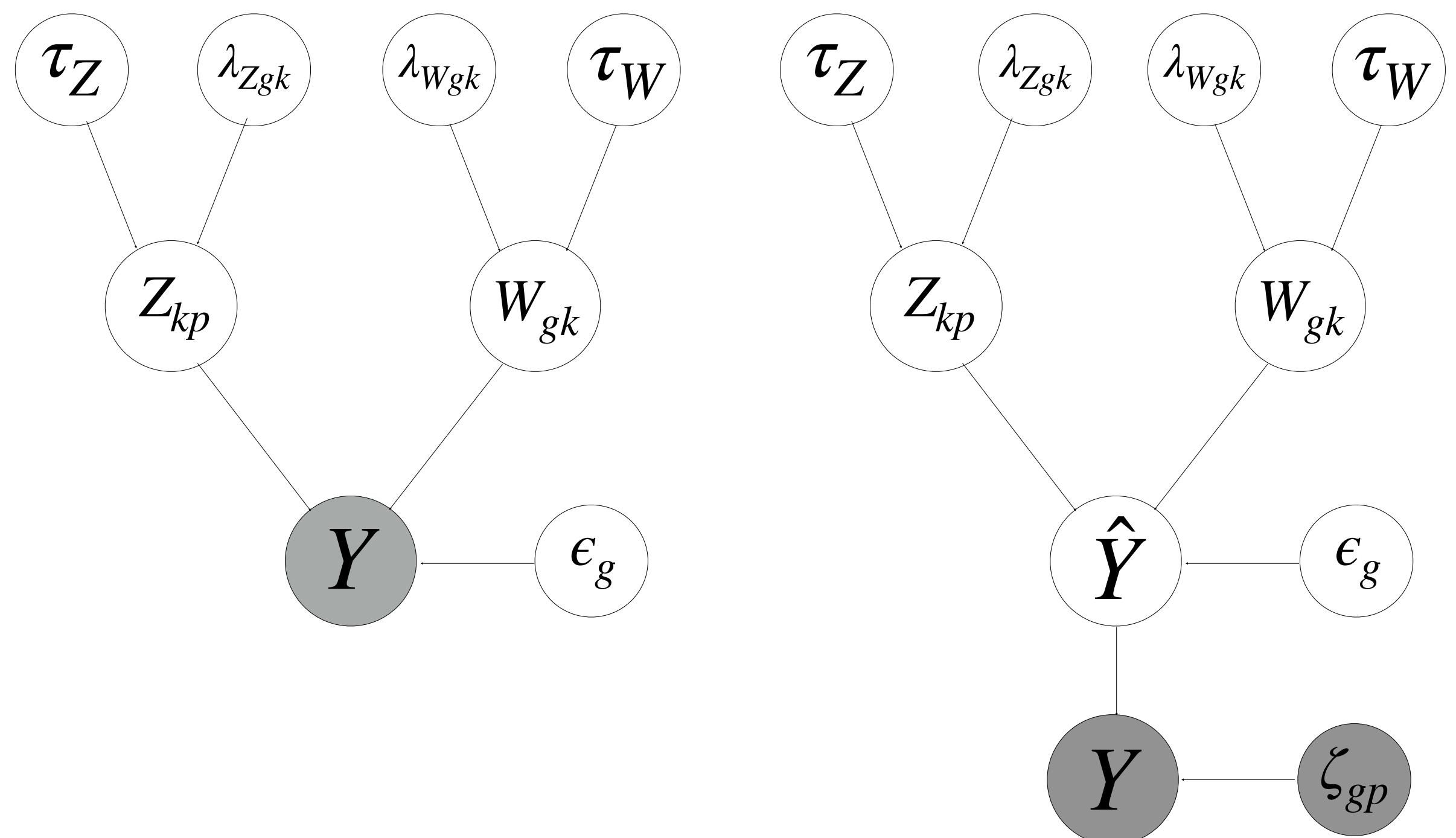
4. Decomposition of noise

$$\hat{b}_{gp} = \sum_k w_{gk} z_{kp} + \varepsilon_g + \zeta_{gp}$$

IMPLEMENTATION

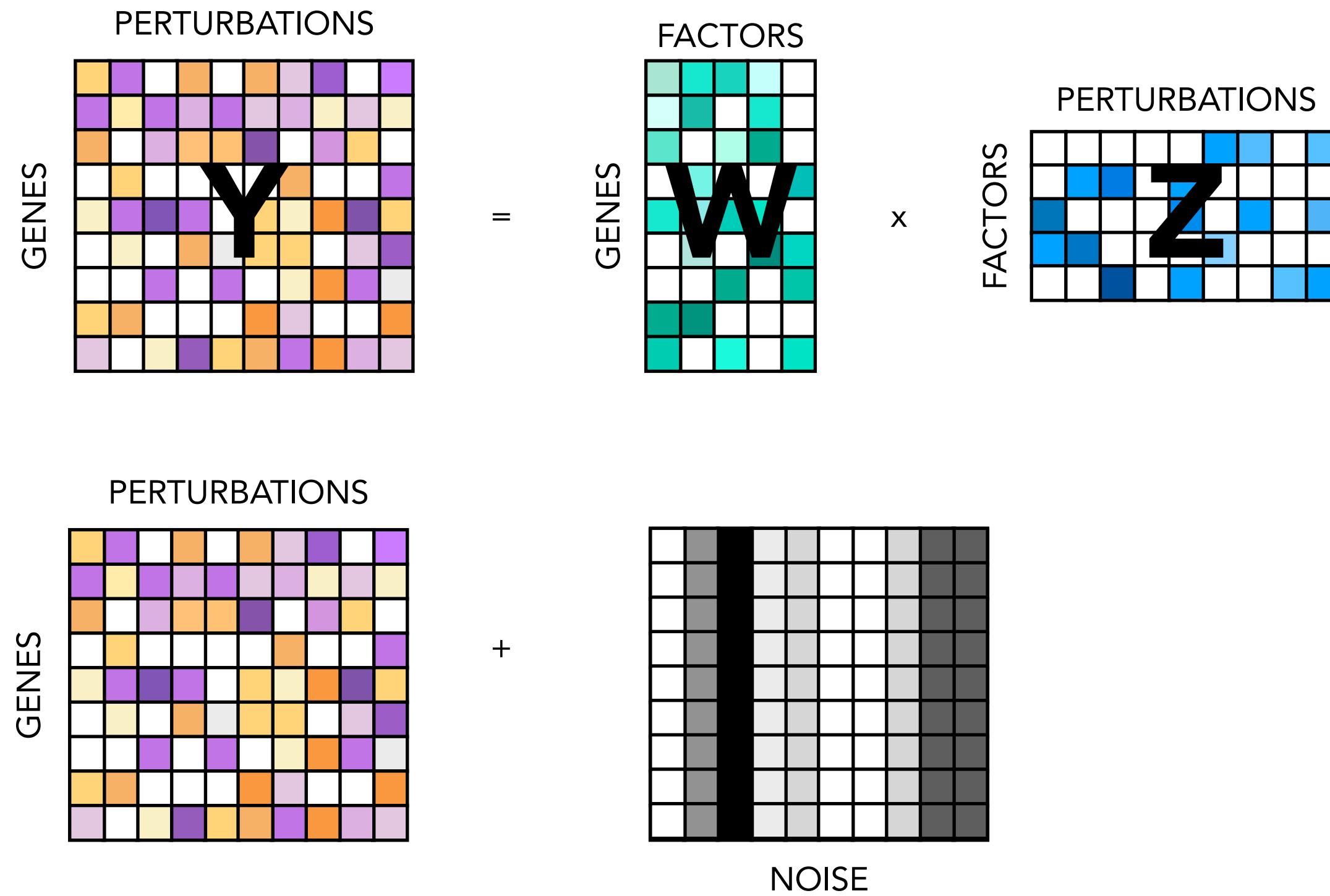


Implementation of two different models in Pyro:



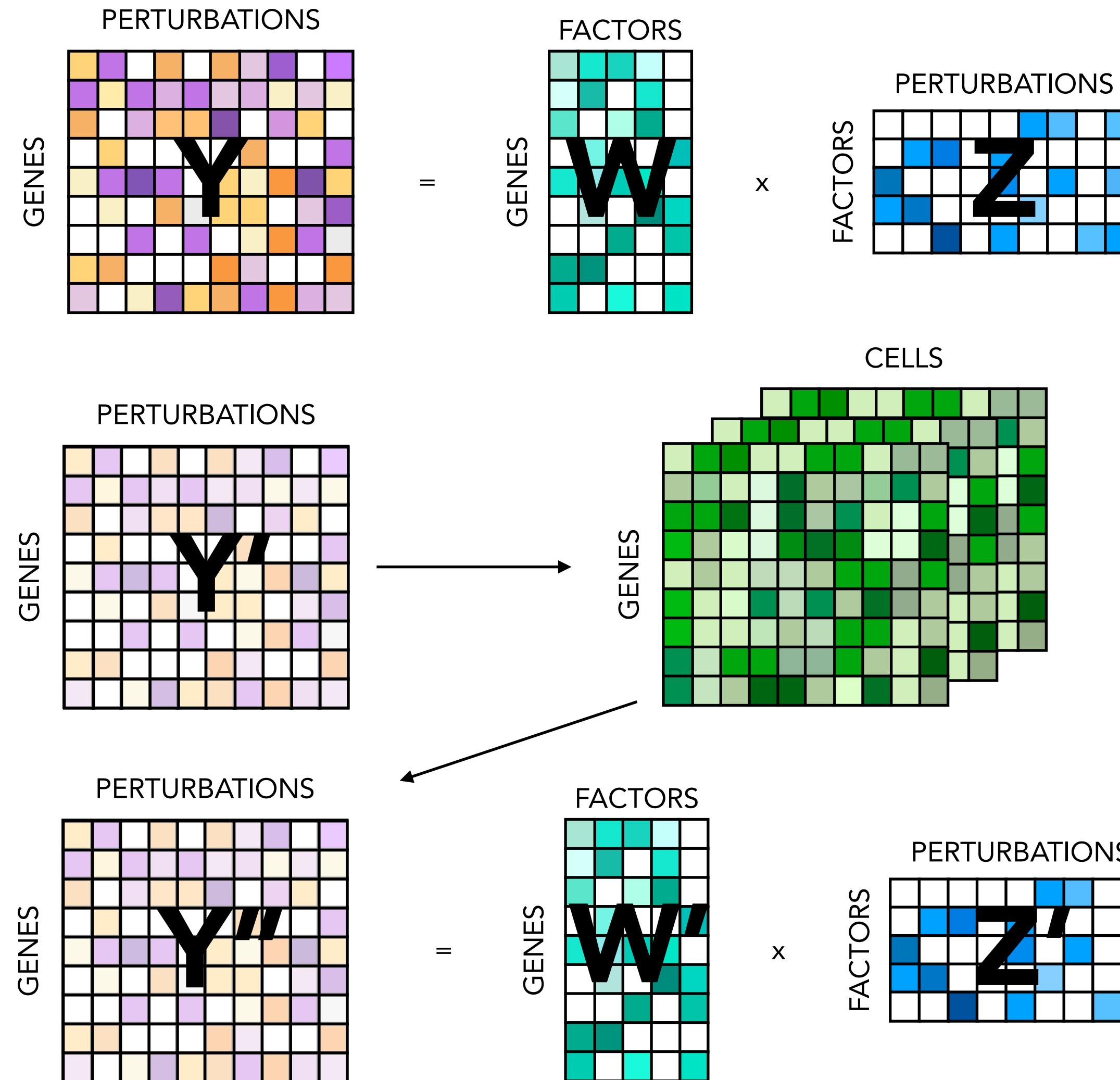
1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

DATA SIMULATION



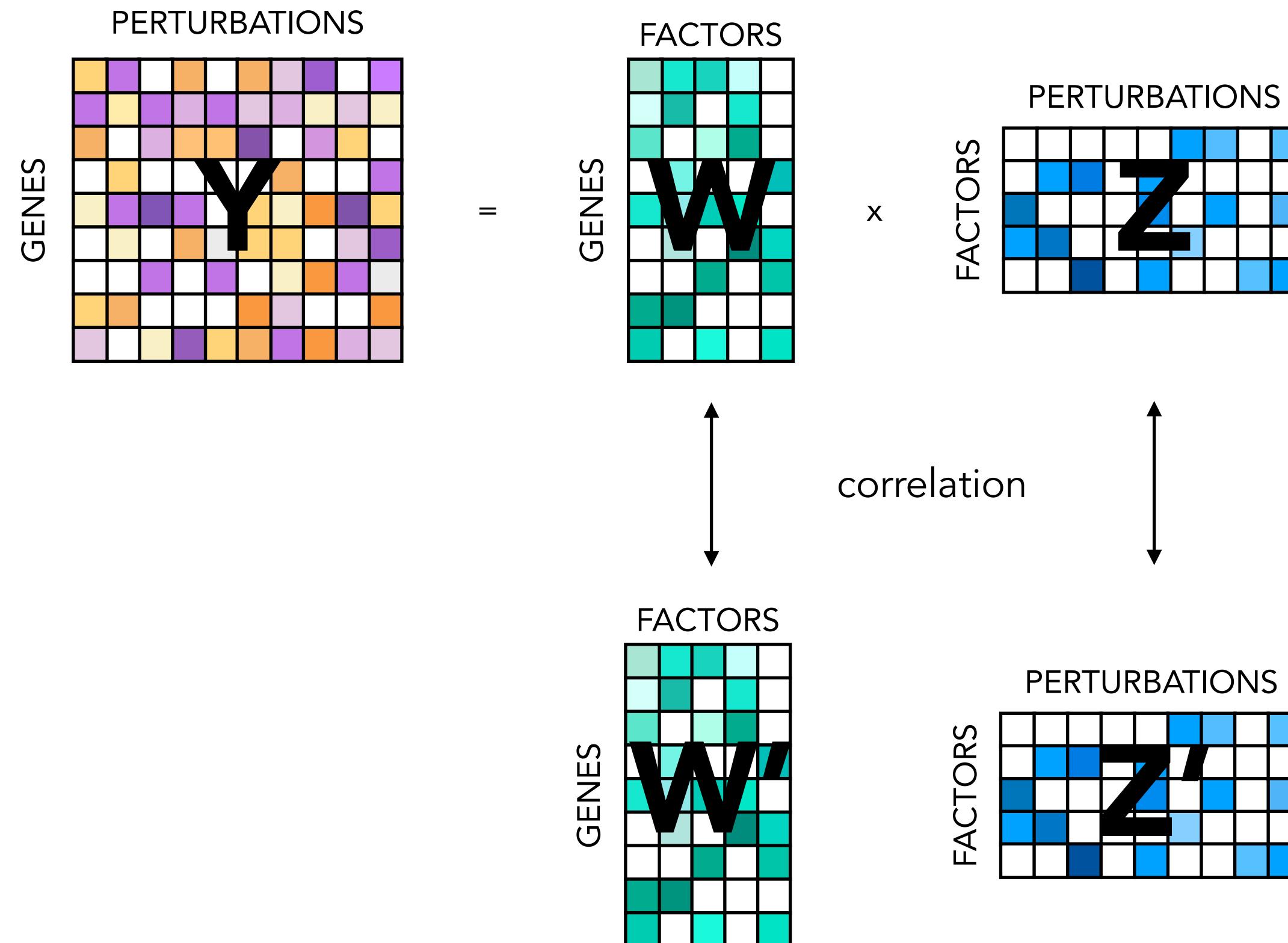
1. Simulate sparse W and Z ;
2. Compute Y ;
3. Add noise:
 - independent
 - gene-specific
 - perturbation-specific

DATA SIMULATION



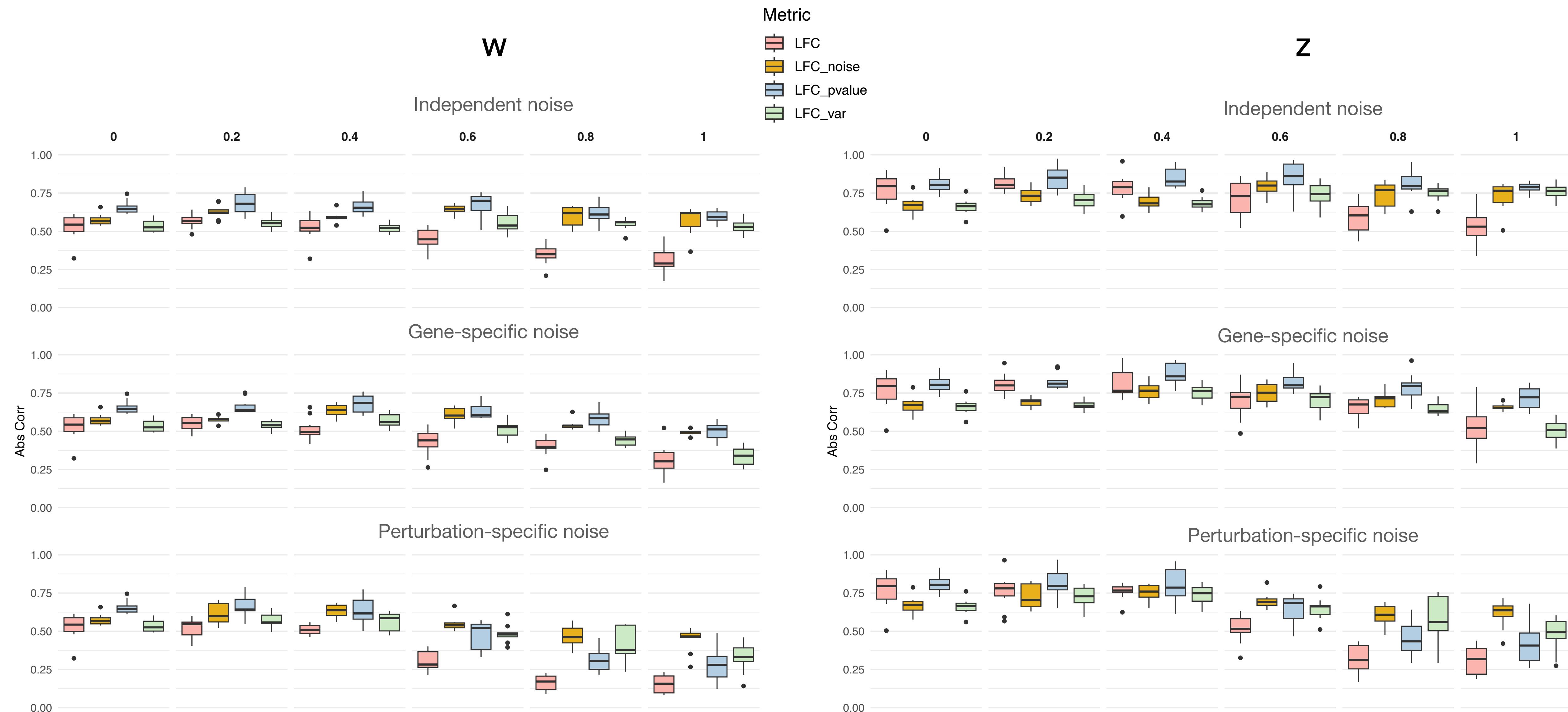
1. Simulate sparse W and Z ;
2. Compute Y ;
3. Add noise:
 - independent
 - gene-specific
 - perturbation-specific
4. Generate control and perturbed single-cell expression based on Y (with adaptation from splatter);
5. Calculate new LFCs and weighted LFCs;
6. Use FA models to get new estimated W' and Z' ;

DATA SIMULATION



1. Simulate sparse W and Z ;
2. Compute Y ;
3. Add noise:
 - independent
 - gene-specific
 - perturbation-specific
4. Generate control and perturbed single-cell expression based on Y (with adaptation from splatter);
5. Calculate new LFCs and weighted LFCs;
6. Use FA models to get new estimated W' and Z' ;
7. Compare estimated with simulated matrices.

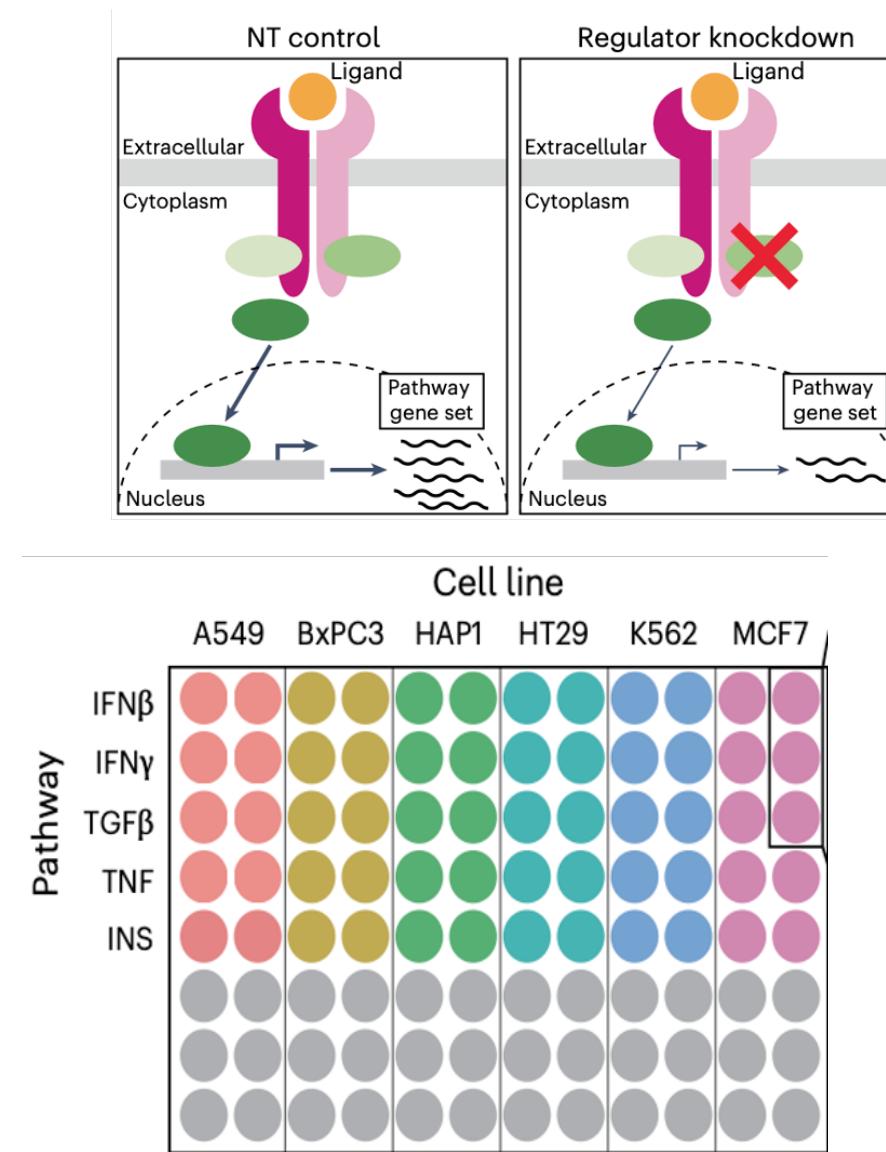
DATA SIMULATION



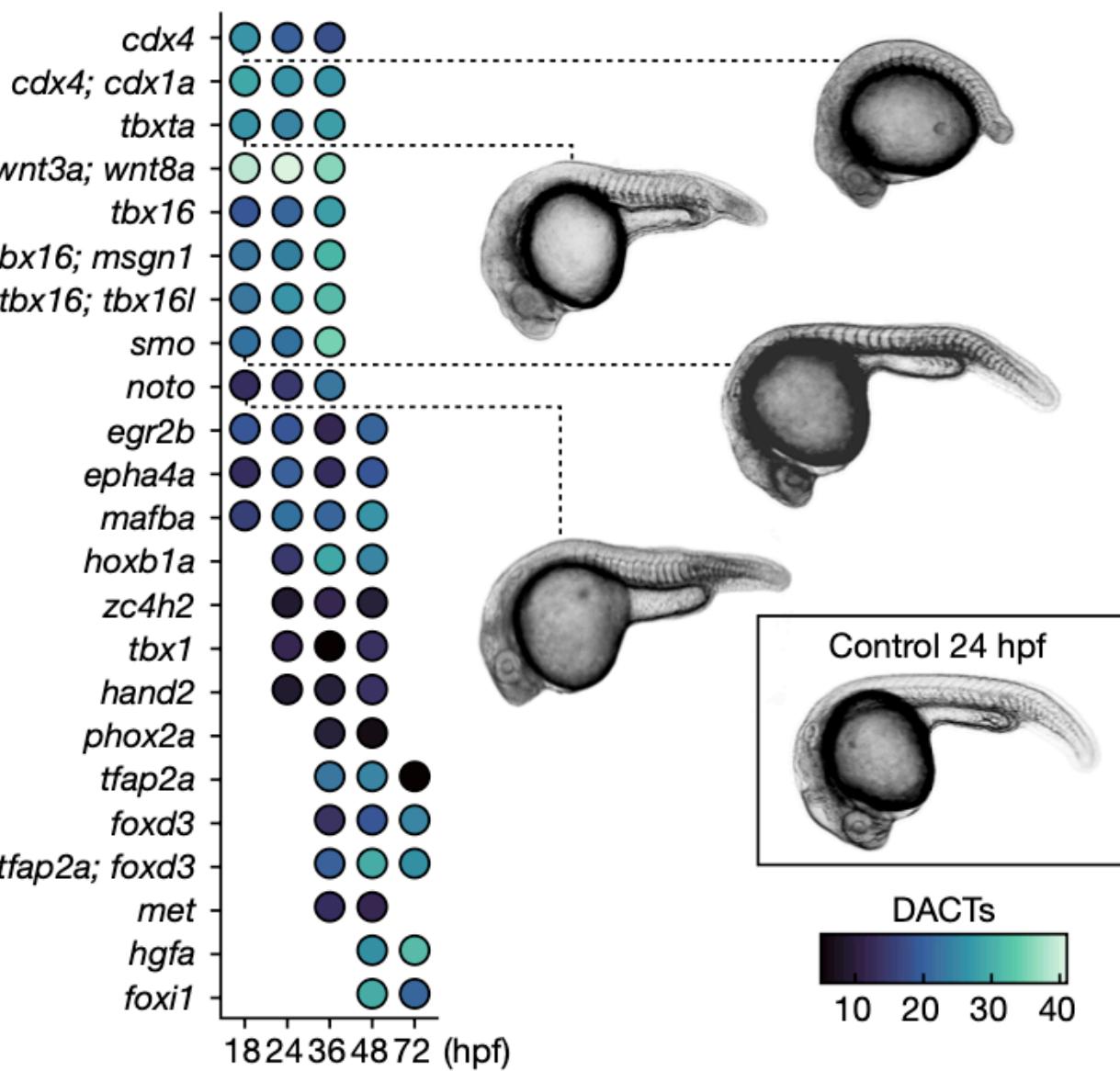
1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

REAL DATASETS

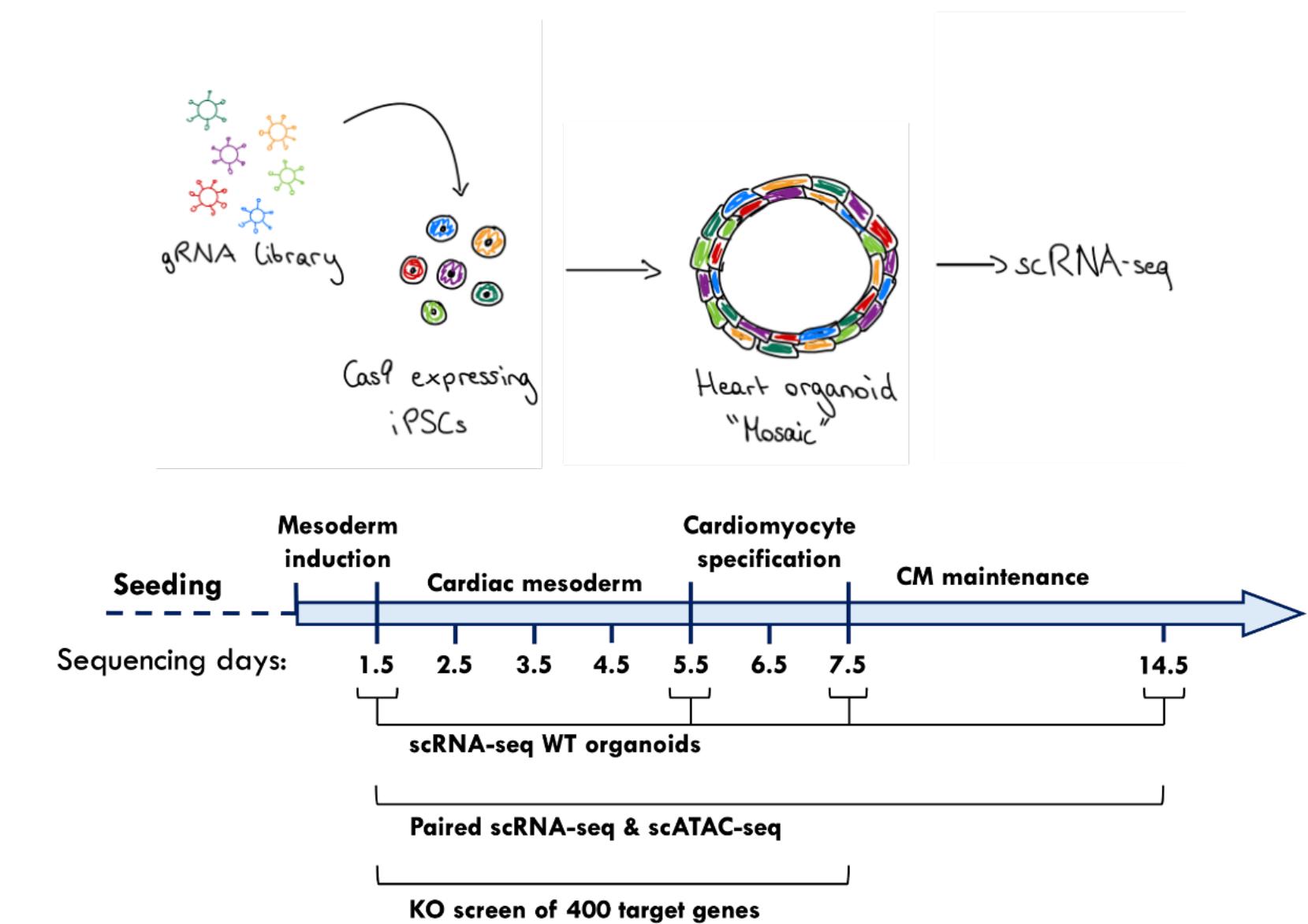
Jiang et al. 2025



Saunders et al. 2023



Heart Organoid dataset, Stegle Group, DKFZ
(still to be published)



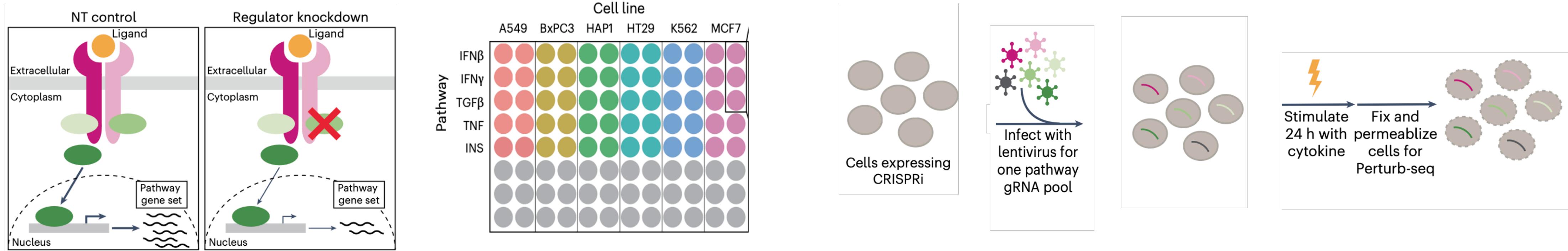
- 6 human cancer cell lines:
lung, breast, colon, bone marrow, pancreas
- 5 cytokine stimuli for major pathways:
IFN β , IFN γ , TGF β , TNF, INS
- 44–61 literature-curated regulators per pathway
- ~ 2.6 million cells profiled
- ~30K genes

- 1812 zebrafish embryos
- 33 major tissues, 99 broad cell types
- 23 genetic perturbations
- 3.2 million cells
- >8 embryos/condition
- 19 timepoints

- 400 genetic perturbations
- 15 cell types
- ~1 million cells
- ~20K genes
- 7 timepoints

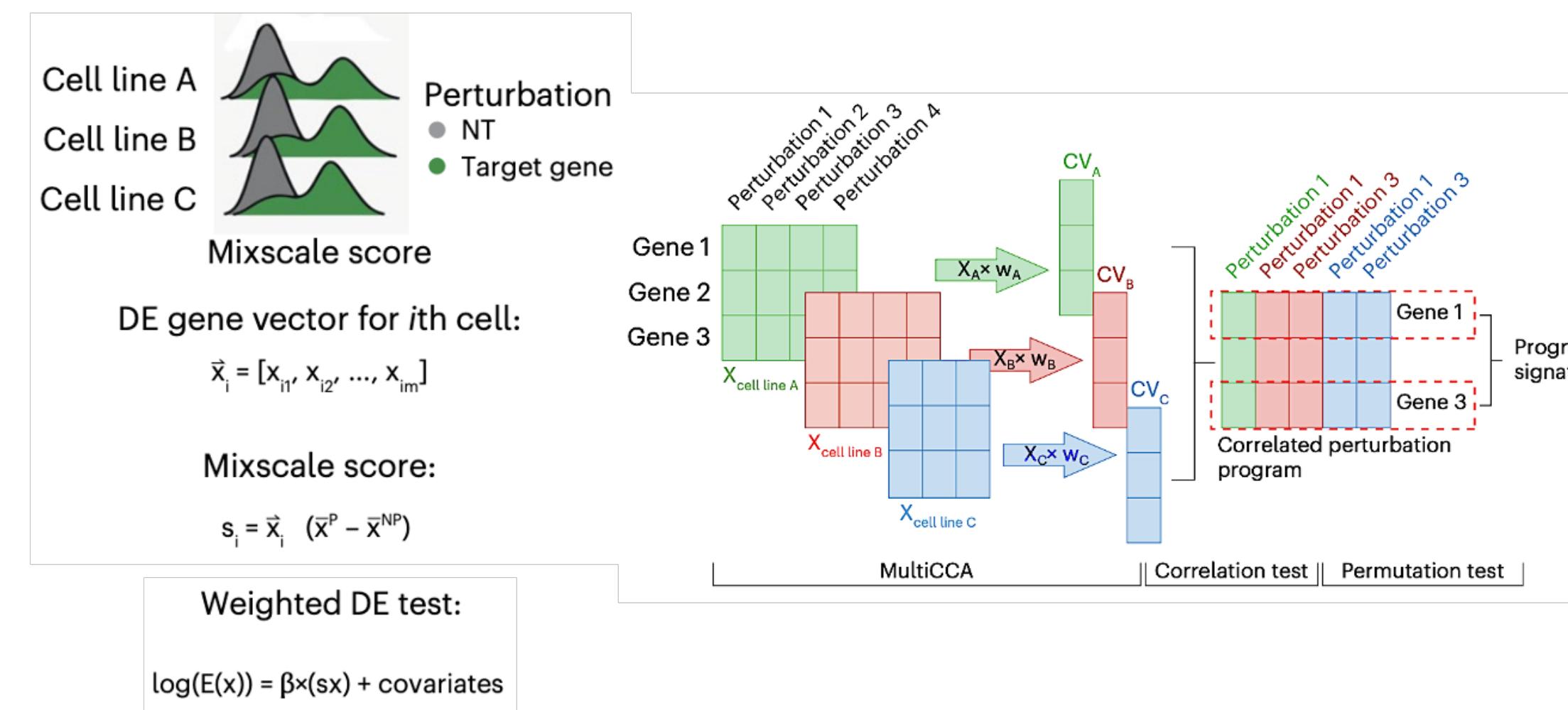
REAL DATASETS

Jiang et al. 2025



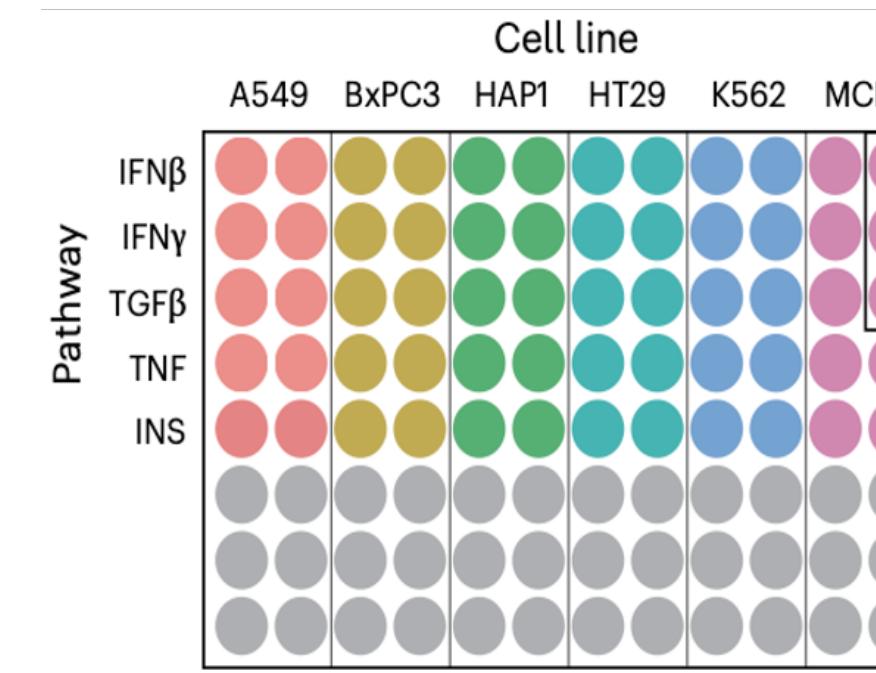
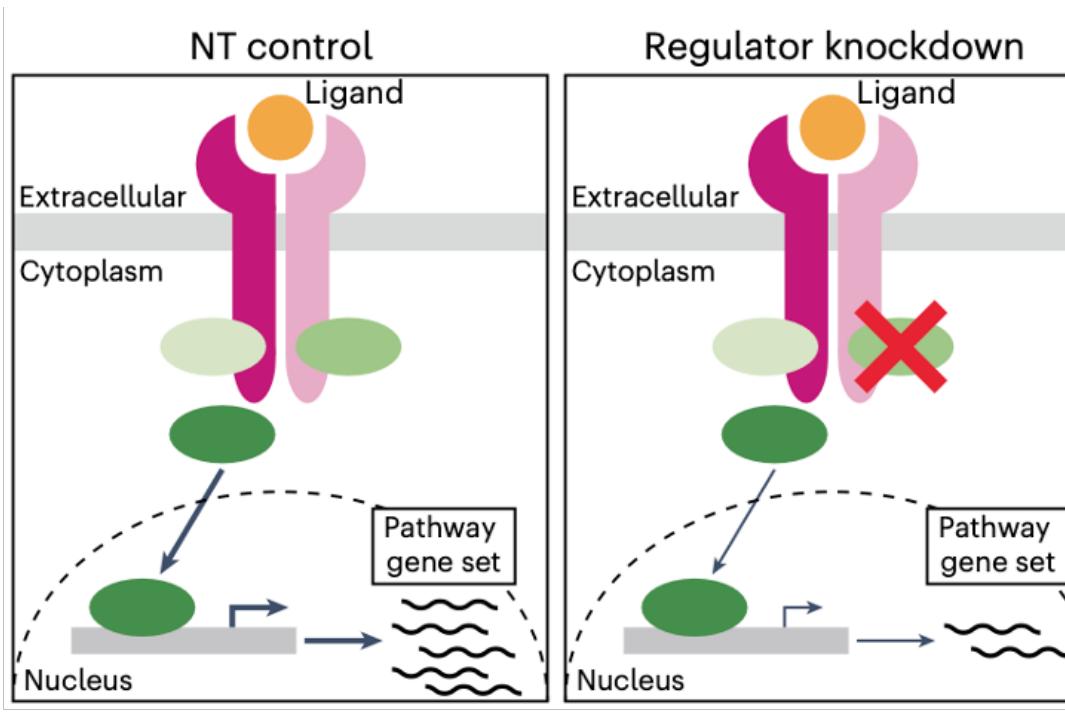
Authors strategy:

Weighted differential expression test + MultiCCA decomposition method

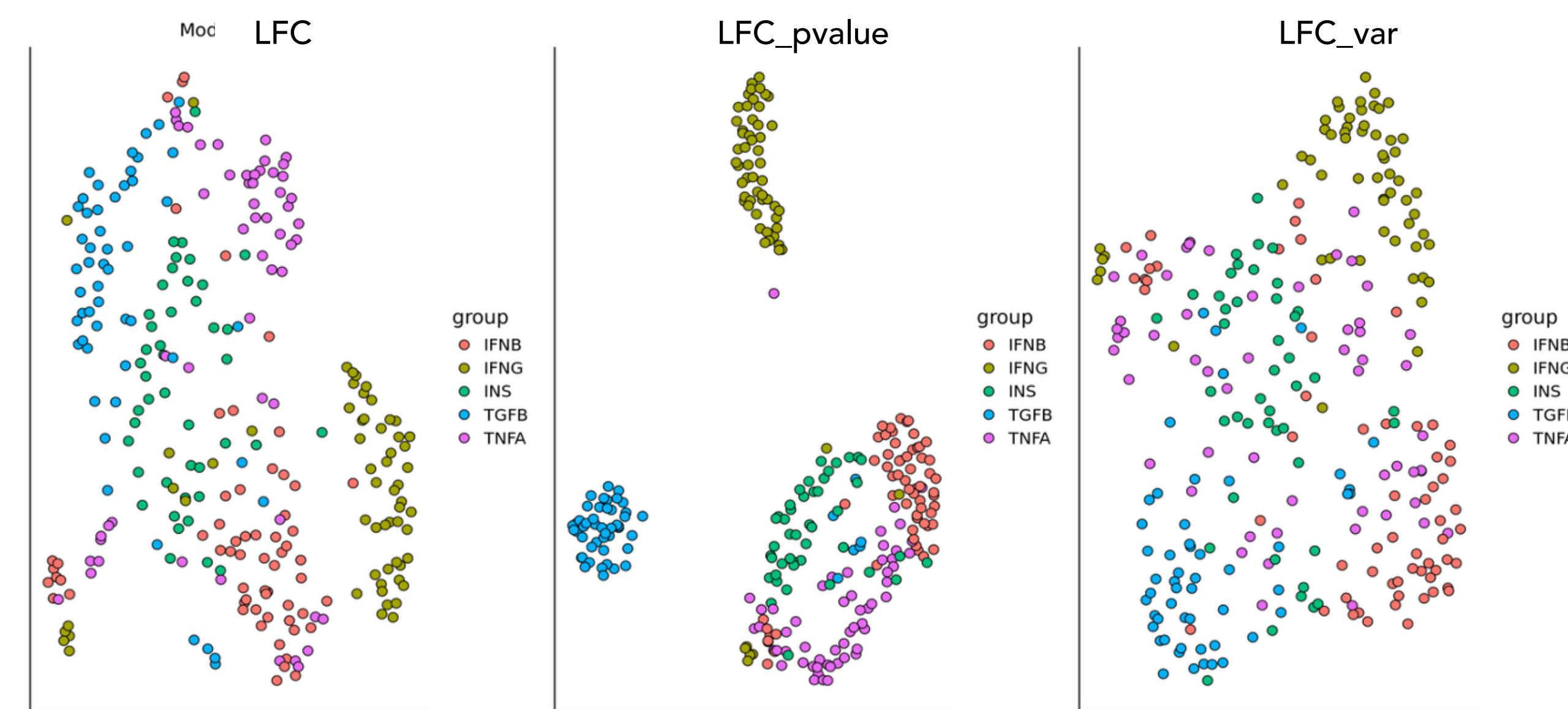


REAL DATASETS

Jiang et al. 2025



Application of MOFA:

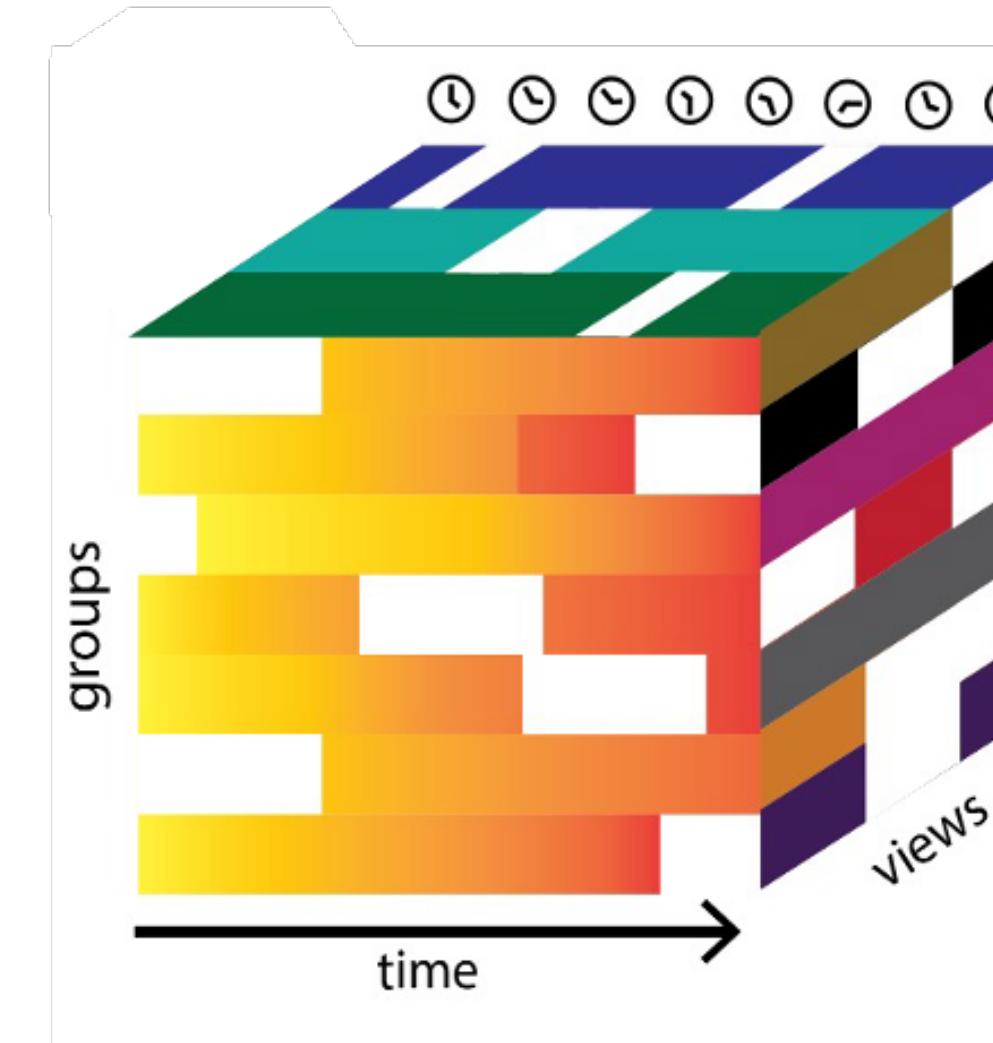
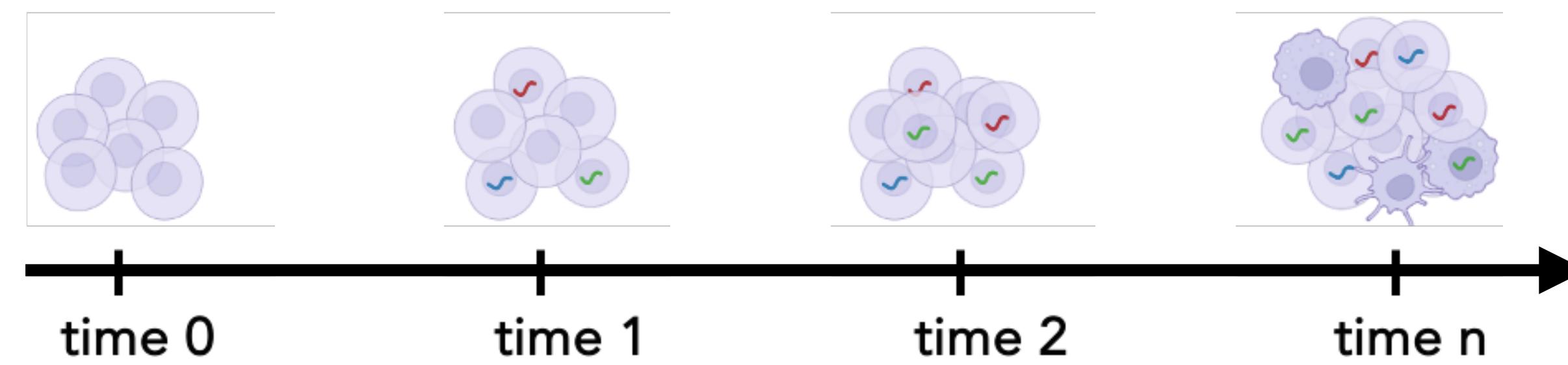


1. Context
2. Approach using Factor Analysis
3. Integrating Uncertainty
4. Data Simulation
5. Application in Real Datasets
6. Conclusion and Future Perspectives

CONCLUSION

- Single-cell CRISPR screens offer rich, context-specific insights—but are still challenging to analyze and interpret;
- Multi-view/group factor models (like MOFA) might allow a more comprehensive analysis of perturbation responses across several biological contexts;
- Integrating uncertainty into the models can improve the interpretability and reliability of inferred responses:
 - Initial simulations suggest that uncertainty-informed models recover more relevant and structured latent factors.
 - Real-data applications are underway - early results look promising.

FUTURE STEPS



ACKNOWLEDGMENTS

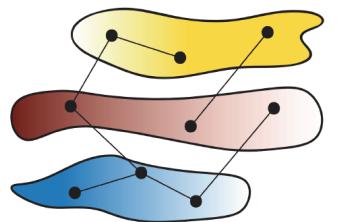
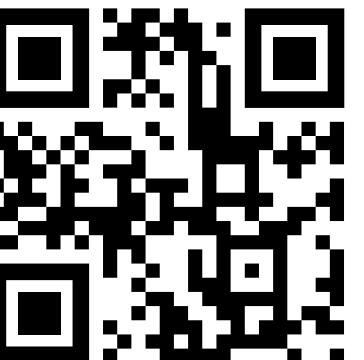
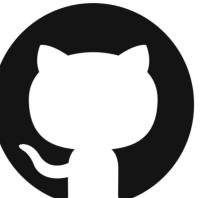
AG Velten

JunProf Dr. Britta Velten
Dr. Nikolai Köhler
Jana Braunger
Zehua Zhang
Purusharth Saxena
Wangjun Hu
Dr. Stijn Hawinkel
Jan Sprengel
Fabian Linsenmeier
Thomas Greulich
Dr. Yen-Hsi Beyer

PhD advisors and TAC Members

JunProf Dr. Britta Velten
Prof Dr. Julio Saez-Rodriguez
JunProf Dr. Lauren Saunders
Dr. Ricardo Ramirez-flores

PAGES



UNIVERSITÄT
HEIDELBERG
ZUKUNFT
SEIT 1386



CENTRE FOR
ORGANISMAL STUDIES



Interdisciplinary Center
for Scientific Computing



e l l i s
European Laboratory for Learning and Intelligent Systems

APPENDIX

ADAPTED SPLATTER

1. Simulate Gene Mean Expression

- Sample gene means λ_i from a **Gamma distribution**:

$$\lambda_i \sim \text{Gamma}(\alpha, \beta)$$

2. Introduce High Expression Outliers

- With probability π_O , a gene is an outlier.

- Replace its mean with an inflated value:

$$\lambda'_i = \text{median}(\lambda) \times F$$

- Inflation factor F sampled from a **log-normal distribution**:

$$F \sim \text{LogNormal}(\mu_O, \sigma_O)$$

3. Adjust for Library Size

- Simulate library sizes L_j from a **log-normal distribution**:

$$L_j \sim \text{LogNormal}(\mu_L, \sigma_L)$$

- Adjust mean expression per cell:

$$\lambda_{i,i} = \lambda_i \times \frac{L_j}{\sum L_j}$$

↓

4. Enforce Mean-Variance Trend

- Sample the **biological coefficient of variation (BCV)** from a **scaled inverse chi-squared distribution**:

$$\text{BCV}_i \sim \text{Scaled-Inv-}\chi^2(\nu, s^2(\lambda_i))$$

- Adjust gene means using a **Gamma distribution**:

$$\lambda_{i,j} \sim \text{Gamma} \left(\frac{1}{\text{BCV}_i^2}, \frac{\lambda_i}{\text{BCV}_i^2} \right)$$

5. Generate Count Matrix

- Sample counts from a **Poisson distribution**:

$$Y_{i,j} \sim \text{Poisson}(\lambda_{i,j})$$

6. Introduce Dropout (Zero Inflation)

- Define dropout probability using a **logistic function**:

$$P(\text{zero}|\lambda_{i,j}) = \frac{1}{1 + e^{-k(\lambda_{i,j} - x_0)}}$$

- Replace values with zero based on **Bernoulli sampling**:

$$Z_{i,j} \sim \text{Bernoulli}(P(\text{zero}|\lambda_{i,j}))$$

STEPS

1. Estimate α and β from control cells in real dataset
2. Generate control dataset (I am not using outliers nor dropout), with random number of cells.
3. Use simulated Y values to get gene means for each perturbation:
`control_means * 2^logFC`
4. Generate perturbation dataset with random number of cells for each perturbation:
 - Input gene means instead of sampling from gamma

APPENDIX

What is BBVI and why do we need it?

CORE IDEA

BBVI is a flexible alternative that uses **sampling and gradients** to estimate and optimize the ELBO, without **model-specific derivations**.

GOAL

We want to compute the gradient: $\nabla_{\phi} \mathcal{L}(\phi) = \nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z)} \left[\log \frac{p(x, z)}{q_{\phi}(z)} \right]$

PROBLEM

Exact integration is intractable, we need to **approximate the expectation using sampling**
But z **is sampled from $q_{\phi}(z)$ which depends on ϕ** , so we can't move the gradient inside naively

$$\nabla_{\phi} \mathbb{E}_{z \sim q_{\phi}(z)} [f(z)] \neq \mathbb{E}_{z \sim q_{\phi}(z)} [\nabla_{\phi} f(z)]$$

APPROACH

Rewrite the expression as an **expectation of the gradient** using tricks

Score-Function Estimator (REINFORCE)

Reparameterization Trick

APPENDIX

BBVI in practice

Pseudocode:

1. Initialize variational parameters λ

2. Repeat until convergence:

- Sample $z^{(1)}, \dots, z^{(S)} \sim q_\lambda(z)$

→ use reparameterization if available, otherwise sample directly

- For each sample $z^{(s)}$, compute:

- $\log p(x, z^{(s)})$

- $\log q_\lambda(z^{(s)})$

- If using score-function:

$$\nabla_\lambda \log q_\lambda(z^{(s)})$$

- If using reparameterization:

$$\nabla_\lambda \left(\log p(x, z^{(s)}) - \log q_\lambda(z^{(s)}) \right)$$

- Estimate gradient of ELBO:

- Score-function:

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S \left(\log p(x, z^{(s)}) - \log q_\lambda(z^{(s)}) \right) \cdot \nabla_\lambda \log q_\lambda(z^{(s)})$$

- Reparameterization:

$$\nabla_\lambda \mathcal{L} \approx \frac{1}{S} \sum_{s=1}^S \nabla_\lambda \left(\log p(x, z^{(s)}) - \log q_\lambda(z^{(s)}) \right)$$

- Update λ using gradient ascent (e.g. Adam)



You write only the model and the guide.

Pyro takes care of ELBO gradients, chooses the right estimator, and runs the optimizer in just a few lines of code.

Model $p(x, z)$

```
def model(x): ...
    pyro.sample(...)
```

Guide $q_\phi(z)$

```
def guide(x): ...
    pyro.sample(...)
```

ELBO Objective

```
loss=Trace_ELBO()
```

Gradient Ascent on ELBO

```
svi.step(x)
```

✗ No ELBO derivation

✗ No gradient estimator coding

✗ No optimizer steps written by hand