



THE UNIVERSITY  
*of* EDINBURGH

# Higher-order interactions in statistical physics, machine learning and biomedicine

---

Ava Khamseh

Biomedical AI Lab

January 2023

# Collaborators

## Institute of Genetics and Cancer

Abel Jansma

Yuelin Yao

Jareth Wolf

Ava Khamseh (and School of Informatics)

Chris Ponting

## School of Physics

Luigi Del Debbio

## School of Mathematics

Sjoerd Beentjes

# A cross-disciplinary journey (2017-2023)

Study of interactions:

Part 1: Using a neural network to estimate interactions in Ising model

Part 2: model-independent estimation of interactions directly from data

Part 3: Biological interpretation of interactions and application in biomedicine

# Forward vs inverse problem

Forward problem (Statistical Physics): The goal is to provide a macroscopic description of Nature by deriving observable quantities from underlying laws.

- Ising model forward problem: Obtain observables such as magnetisation, energy and correlations, given the Hamiltonian and its parameters

Inverse problem: Starting point are observations (data), the goal is to infer microscopic properties of the system

- Estimate Ising interactions directly from data

# Interactions

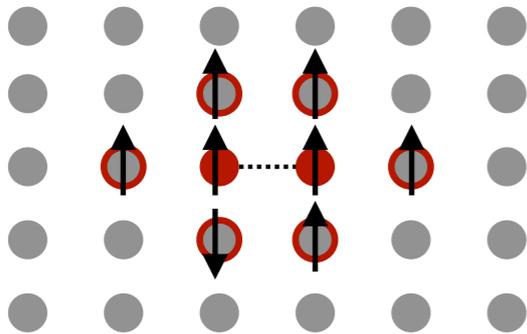
## Elementary Particles

### Standard Model of Elementary Particles

three generations of matter (fermions)			interactions / force carriers (bosons)	
I	II	III		
mass =2.2 MeV/c <sup>2</sup> charge 2/3 spin 1/2	mass =1.28 GeV/c <sup>2</sup> charge 2/3 spin 1/2	mass =173.1 GeV/c <sup>2</sup> charge 2/3 spin 1/2	0 0 0	=124.97 GeV/c <sup>2</sup> 0 0
<b>u</b> up	<b>c</b> charm	<b>t</b> top	<b>g</b> gluon	<b>H</b> higgs
mass =4.7 MeV/c <sup>2</sup> charge -1/3 spin 1/2	mass =96 MeV/c <sup>2</sup> charge -1/3 spin 1/2	mass =4.18 GeV/c <sup>2</sup> charge -1/3 spin 1/2	0 0 0	
<b>d</b> down	<b>s</b> strange	<b>b</b> bottom	<b>γ</b> photon	
mass =0.511 MeV/c <sup>2</sup> charge -1 spin 1/2	mass =105.66 MeV/c <sup>2</sup> charge -1 spin 1/2	mass =1.7768 GeV/c <sup>2</sup> charge -1 spin 1/2	=91.19 GeV/c <sup>2</sup> 0 1	
<b>e</b> electron	<b>μ</b> muon	<b>τ</b> tau	<b>Z</b> Z boson	
mass <1.0 eV/c <sup>2</sup> charge 0 spin 1/2	mass <0.17 MeV/c <sup>2</sup> charge 0 spin 1/2	mass <18.2 MeV/c <sup>2</sup> charge 0 spin 1/2	=80.39 GeV/c <sup>2</sup> ±1 1	
<b>ν<sub>e</sub></b> electron neutrino	<b>ν<sub>μ</sub></b> muon neutrino	<b>ν<sub>τ</sub></b> tau neutrino	<b>W</b> W boson	

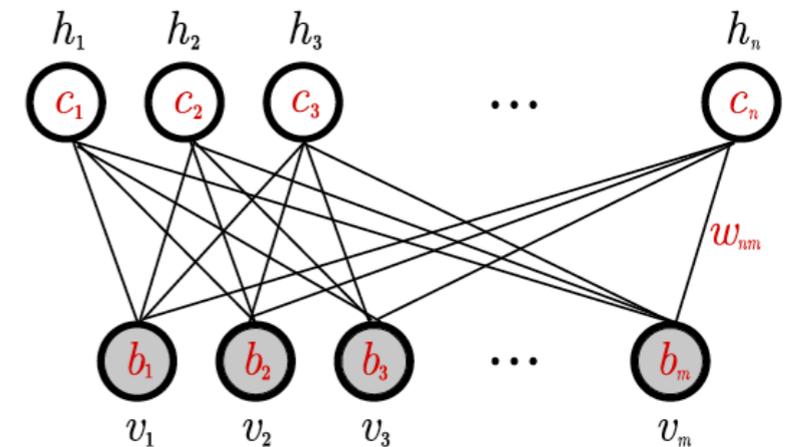
LEPTONS (left side)  
GAUGE BOSONS (middle)  
VECTOR BOSONS (right side)  
SCALAR BOSONS (far right)

## Statistical Physics

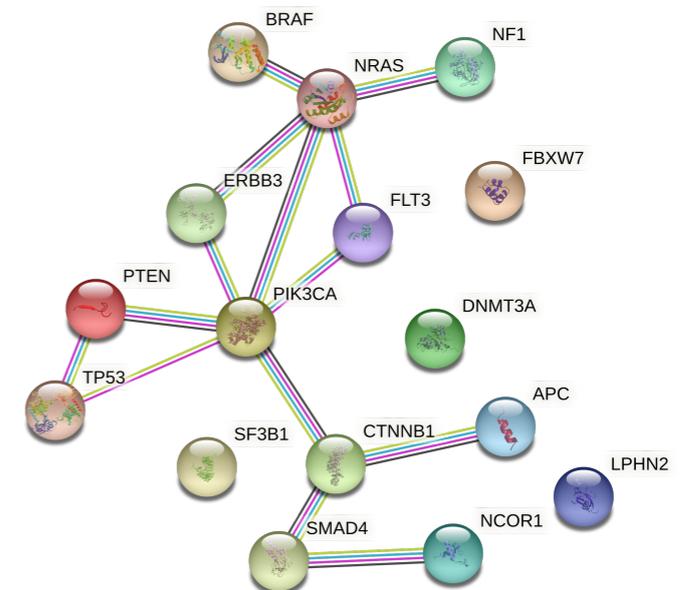


Dynamical/physical,  
vs,  
Probabilistic

## Nodes in a neural network



## Biomedicine



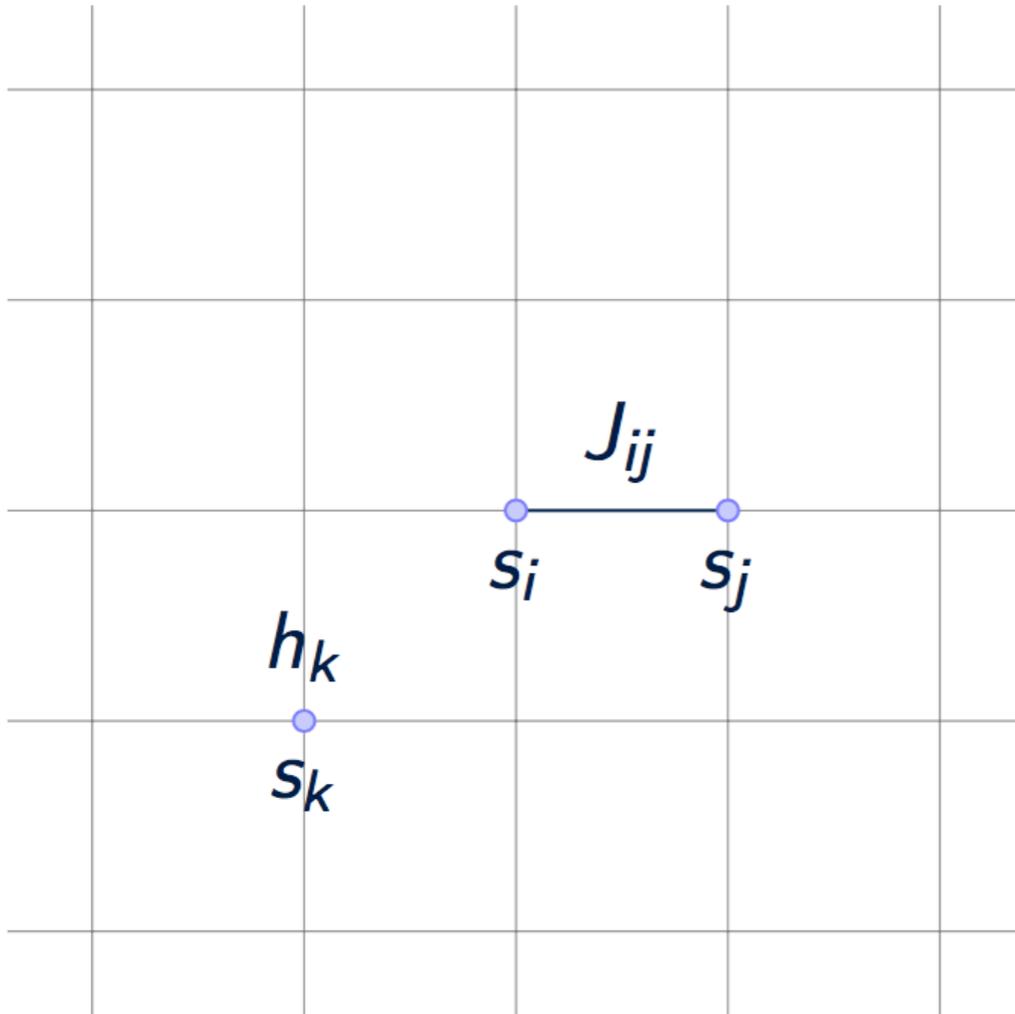
## Part 1

# Interactions: The Ising model & RBM

# Ising Model

$$p_D(s) = \frac{1}{Z(J, h)} e^{-H_{J,h}(s)}$$

$$H_{J,h} = - \sum_{i,j} J_{ij} s_i s_j - \sum_i h_i s_i \quad , \quad Z(J, h) = \sum_s e^{-H_{J,h}(s)}$$



MC simulation of the 2D Ising model at various temperatures

generate a sample

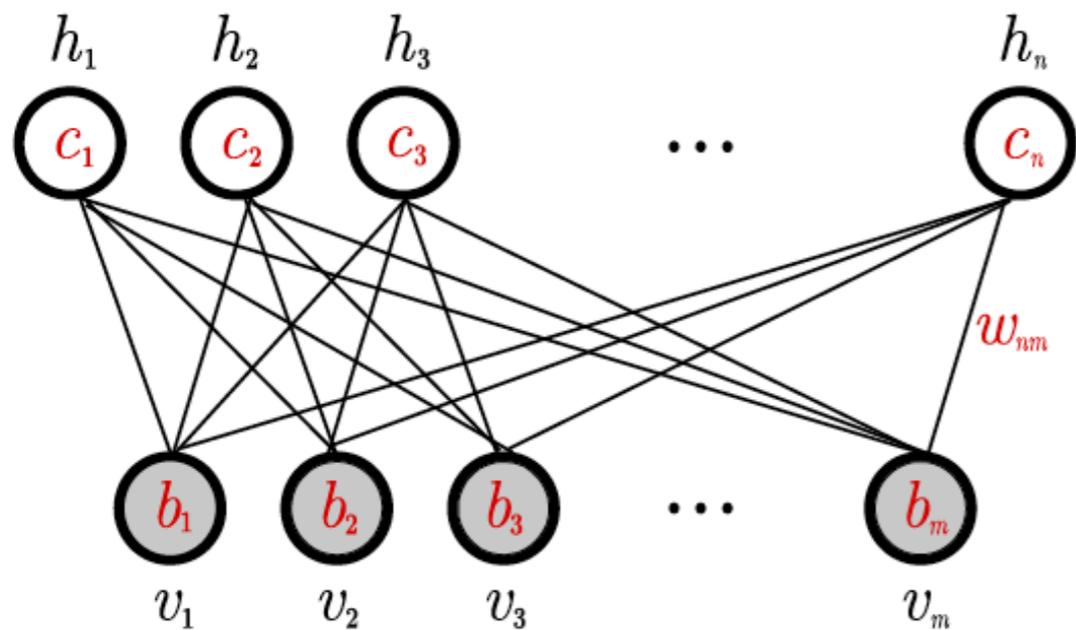
$$D = \{s^1, s^2, \dots\}, N_D \sim 10^5$$

Onsager 1944

# Restricted Boltzmann Machine (RBM)

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{i=1}^n c_i h_i - \sum_{j=1}^m b_j v_j$$

$$p_{\text{RBM}}(\mathbf{v}, \mathbf{h} | \theta) = \frac{1}{Z_{\text{RBM}}} e^{-E_{\theta}(\mathbf{v}, \mathbf{h})}$$



$$D_{\text{KL}}(q_{\text{data}}(\mathbf{v}) || p_{\text{RBM}}(\mathbf{v} | \theta)) = \sum_{\mathbf{v}} q_{\text{data}}(\mathbf{v}) \log \left( \frac{q_{\text{data}}(\mathbf{v})}{p_{\text{RBM}}(\mathbf{v} | \theta)} \right)$$
$$= \sum_{\mathbf{v}} \left( q_{\text{data}}(\mathbf{v}) \log(q_{\text{data}}) - q_{\text{data}}(\mathbf{v}) \log(p_{\text{RBM}}(\mathbf{v} | \theta)) \right)$$

**Max likelihood**  $\iff$  **Min KL divergence**

**Ising configuration training data**

# Observables

$$\langle m \rangle = \frac{1}{L^2} \left\langle \left| \sum_{i=1}^{L^2} s_i \right| \right\rangle,$$

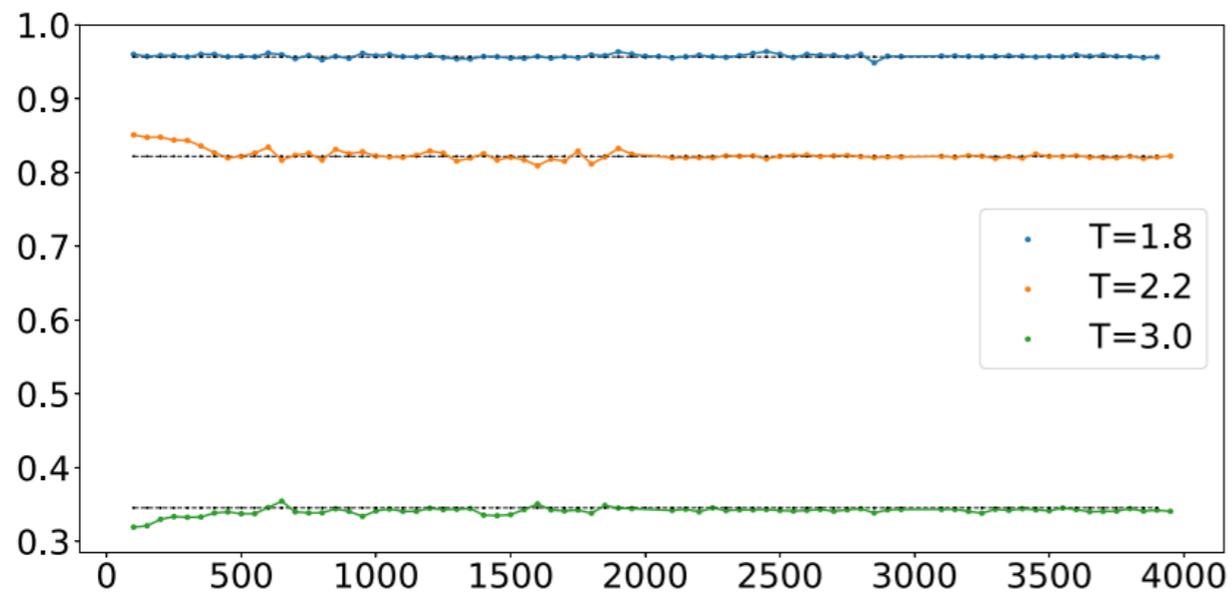
$$\langle \chi \rangle = \frac{L^2}{T} \left\langle \left\langle m^2 \right\rangle - \langle m \rangle^2 \right\rangle,$$

$$\langle E \rangle = -\frac{1}{L^2} \left\langle \sum_{\langle i,j \rangle} s_i s_j \right\rangle,$$

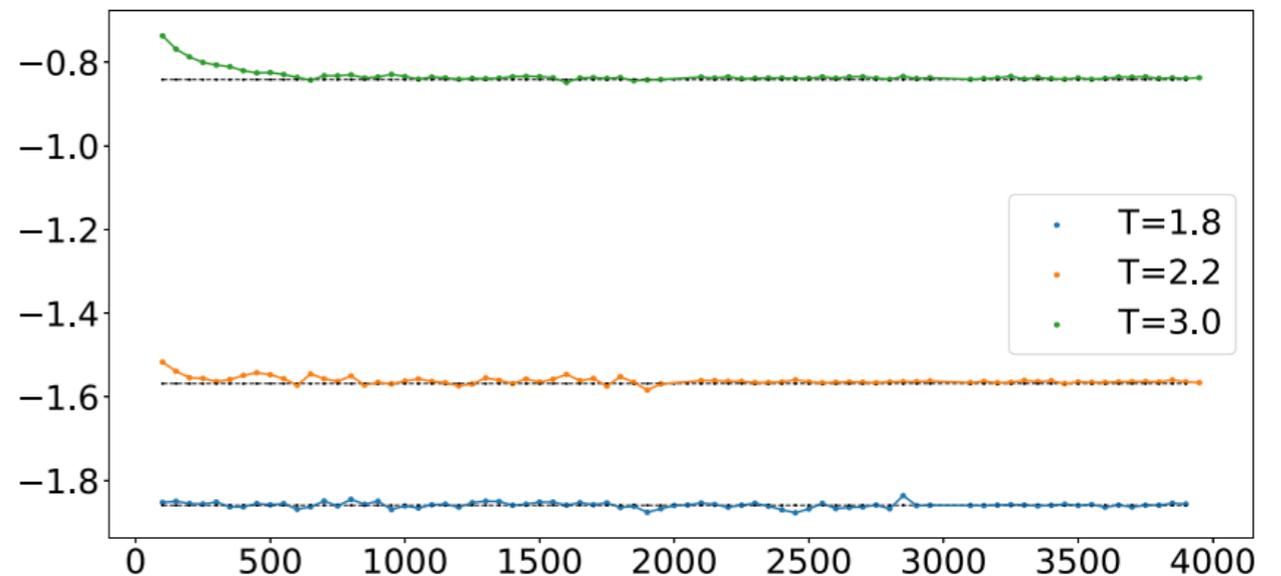
$$\langle C_V \rangle = \frac{L^2}{T^2} \left\langle \left\langle E^2 \right\rangle - \langle E \rangle^2 \right\rangle.$$

# RBM: Observables

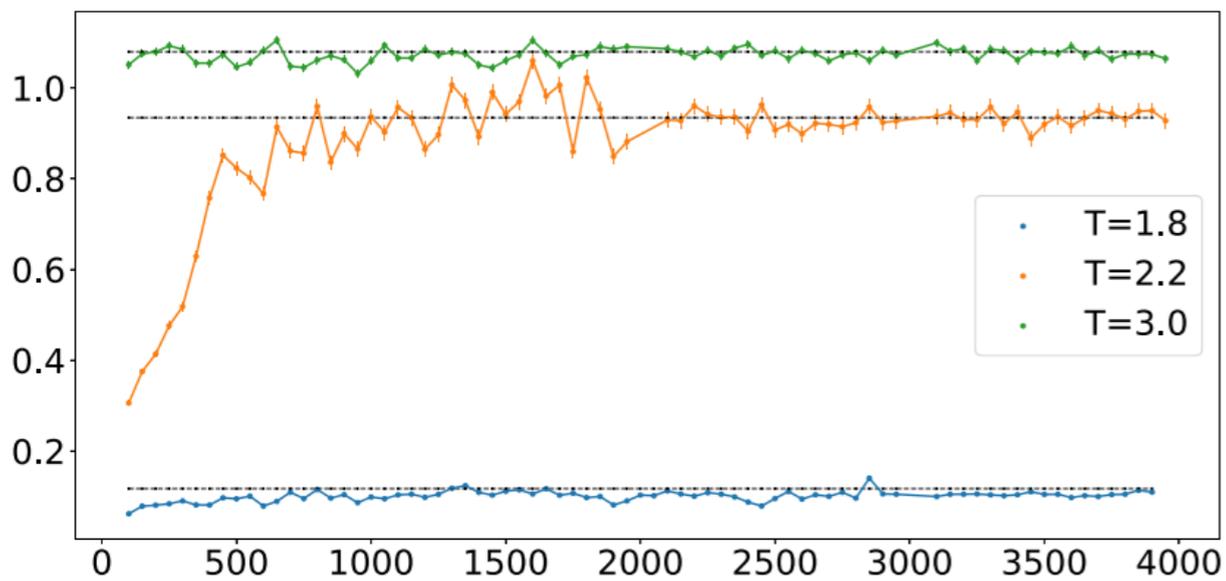
Magnetisation vs number of epochs



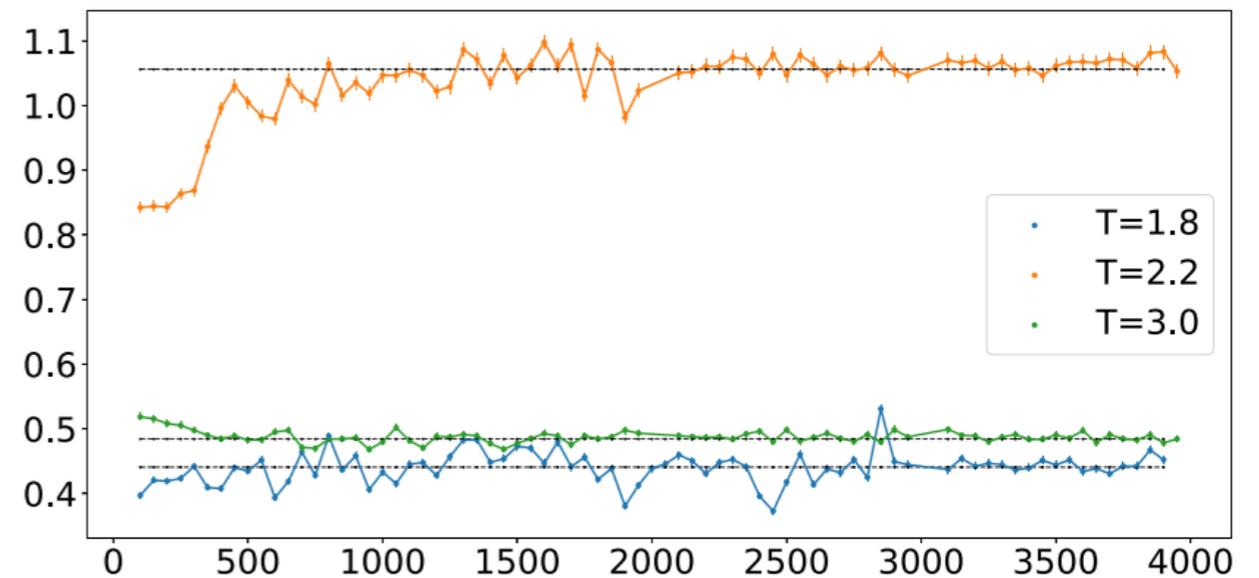
Energy vs number of epochs



Susceptibility vs number of epochs

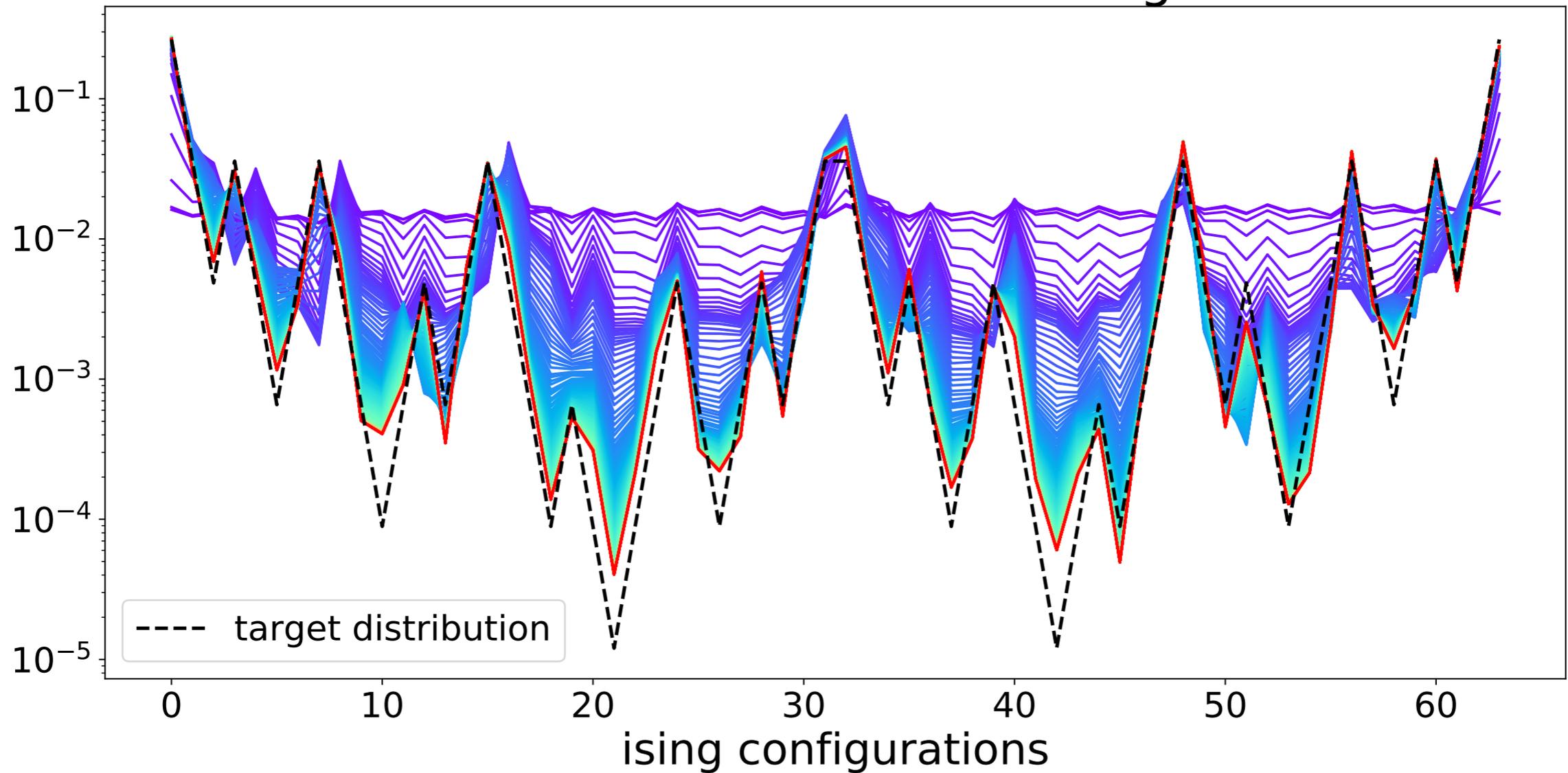


Heat capacity vs number of epochs



# Example: 1D Ising model in 6 variables

Learned distribution vs target



# RBM: Estimation of interactions

$$E(\mathbf{v}) = - \sum_j b_j v_j - \sum_j \left( \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2} \sum_{jk} \left( \sum_i \kappa_i^{(2)} W_{ik} W_{ij} \right) v_j v_k + \dots$$

Beyond pairwise, higher-order couplings

Possible to re-sum the entire series to obtain 2-point coupling!!

For binary data, using cumulant generating function ...

$$v_j^n = v_j, \quad n \in \mathbb{Z}^+$$

finite sum

# RBM: Estimation of interactions

$$E(\mathbf{v}) = - \sum_j b_j v_j - \sum_j \left( \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2} \sum_{jk} \left( \sum_i \kappa_i^{(2)} W_{ik} W_{ij} \right) v_j v_k + \dots$$

Beyond pairwise, higher-order couplings

Possible to re-sum the entire series to obtain 2-point coupling!!

For binary data, using cumulant generating function ...

$$v_j^n = v_j, \quad n \in \mathbb{Z}^+$$

finite sum

e.g., 2-point interaction:

$$H_{j_1 j_2} = \frac{1}{8} \sum_i \ln \frac{(1 + e^{c_i + W_{ij_1} + W_{ij_2}})(1 + e^{c_i})}{(1 + e^{c_i + W_{ij_1}})(1 + e^{c_i + W_{ij_2}})}$$

**Closed form expression!**

# RBM: Estimation of interactions

$$E(\mathbf{v}) = - \sum_j b_j v_j - \sum_j \left( \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2} \sum_{jk} \left( \sum_i \kappa_i^{(2)} W_{ik} W_{ij} \right) v_j v_k + \dots$$

Beyond pairwise, higher-order couplings

Possible to re-sum the entire series to obtain 2-point coupling!!

For binary data, using cumulant generating function ...

$$v_j^n = v_j, \quad n \in \mathbb{Z}^+$$

finite sum

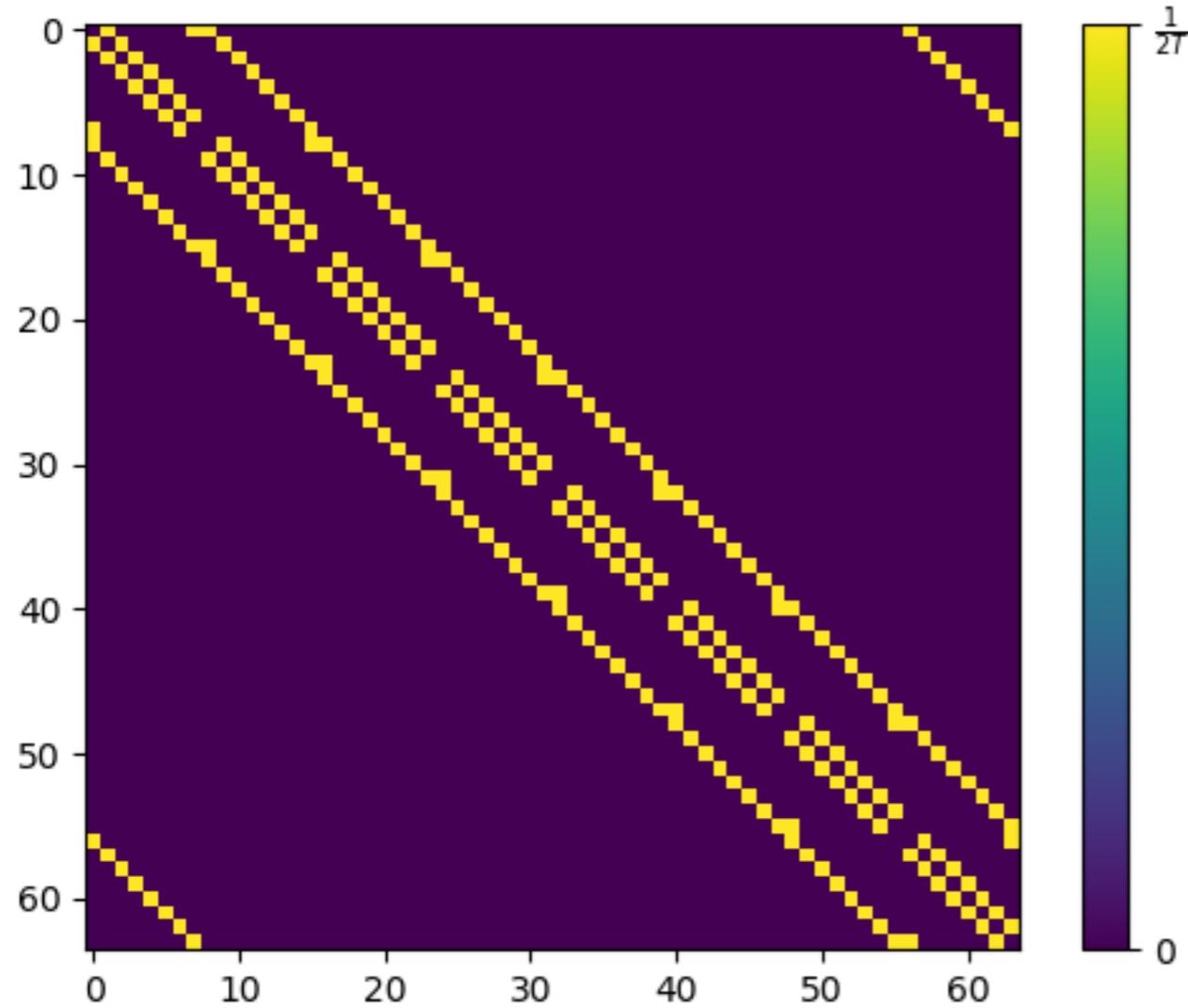
e.g., 3-point interaction:

$$\frac{1}{6} \sum_i \ln \frac{(1 + e^{c_i + W_{ij_1} + W_{ij_2} + W_{ij_3}})(1 + e^{c_i + W_{ij_1}})(1 + e^{c_i + W_{ij_2}})(1 + e^{c_i + W_{ij_3}})}{(1 + e^{c_i + W_{ij_1} + W_{ij_2}})(1 + e^{c_i + W_{ij_1} + W_{ij_3}})(1 + e^{c_i + W_{ij_2} + W_{ij_3}})(1 + e^{c_i})}$$

Closed form expression!

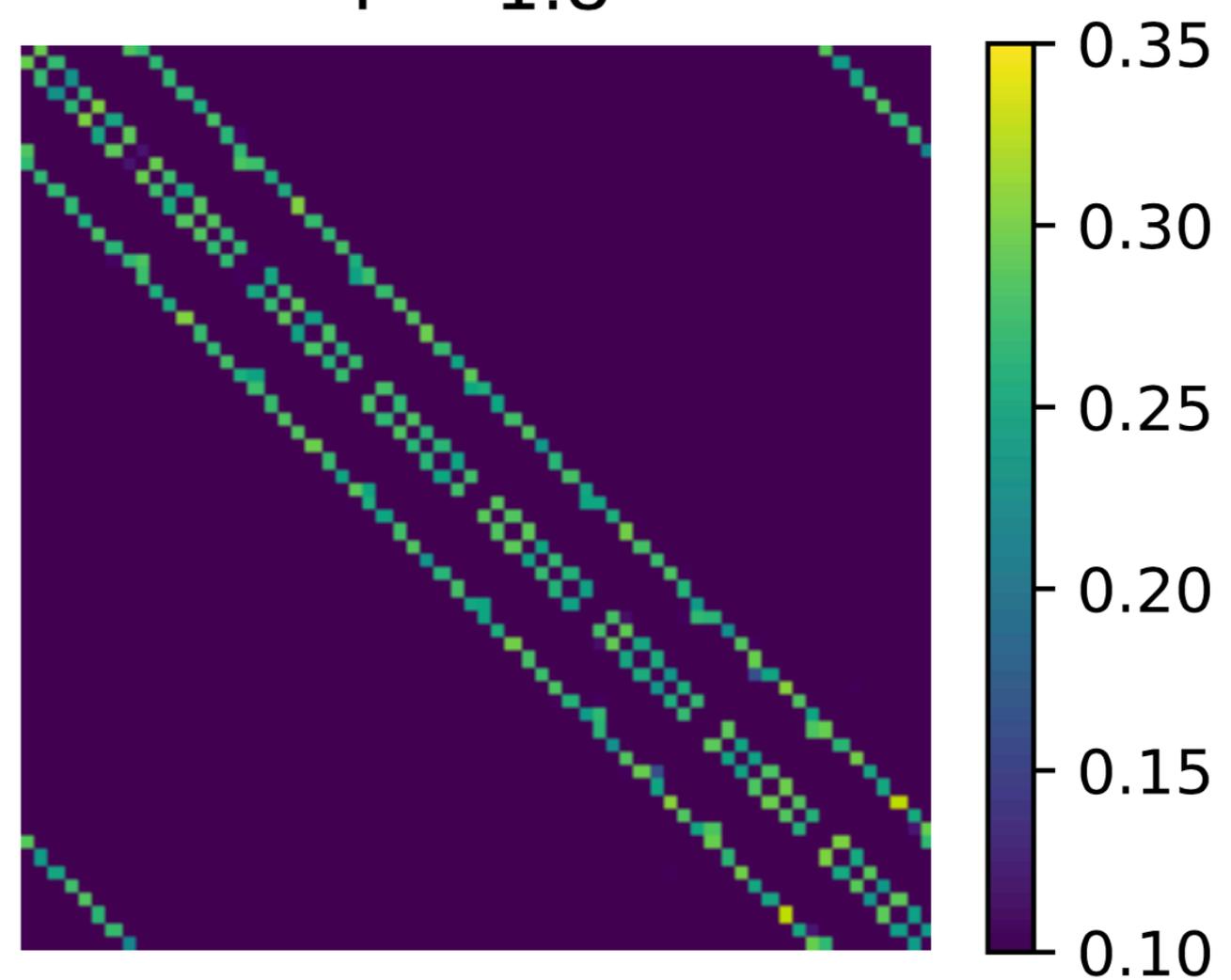
# RBM: Couplings $J_{ij}$

Ising Model

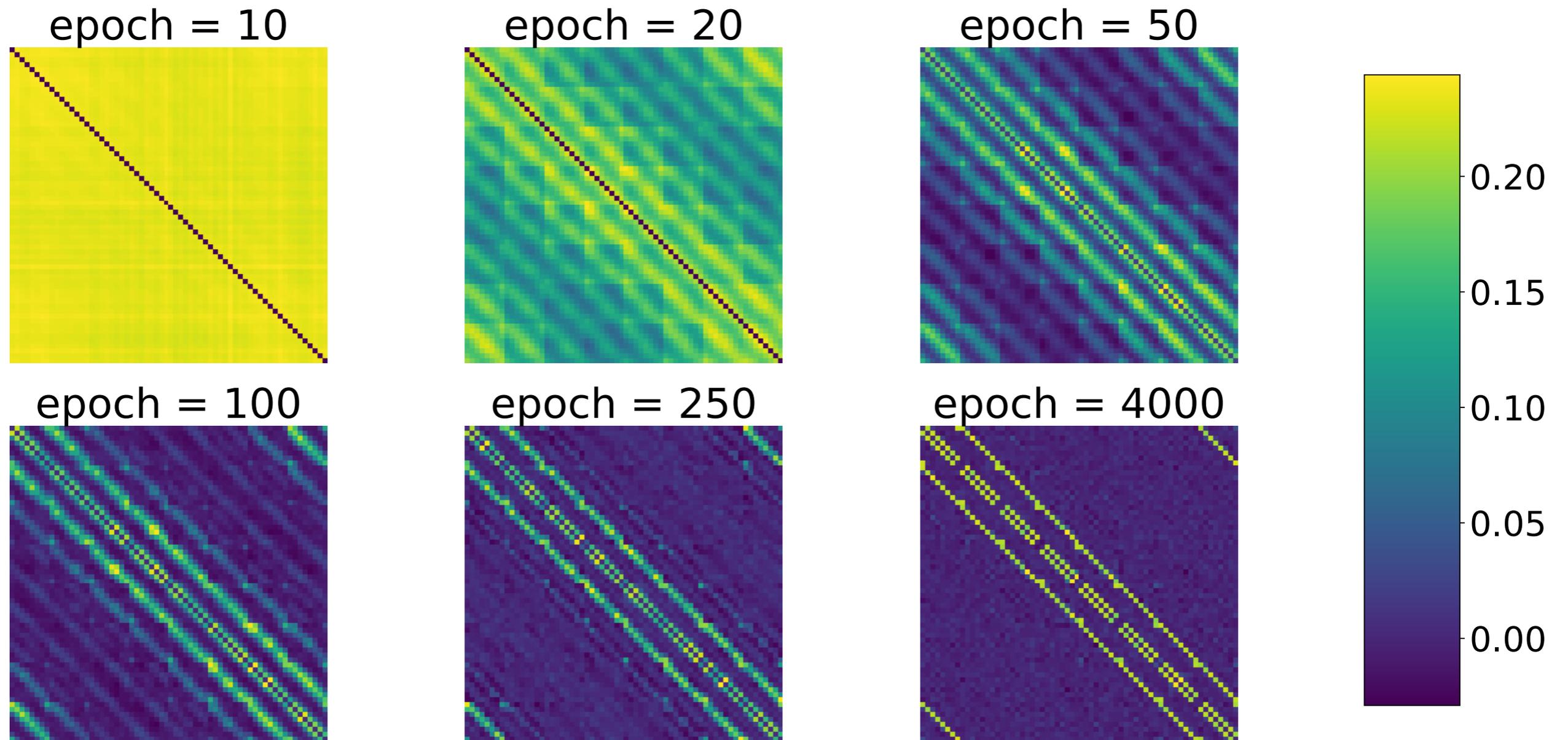


RBM prediction

$T = 1.8$

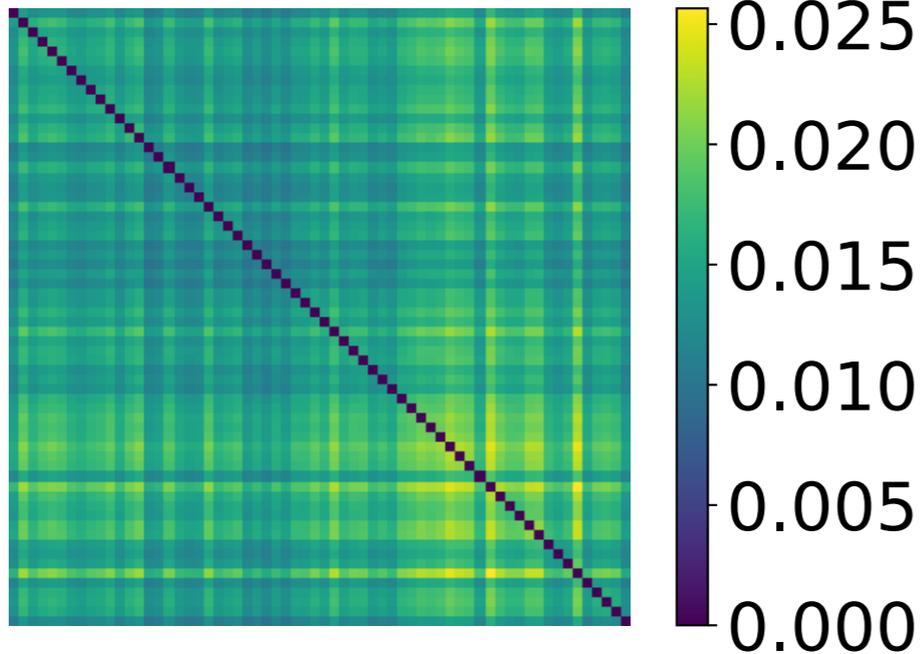


# RBM: Couplings during training



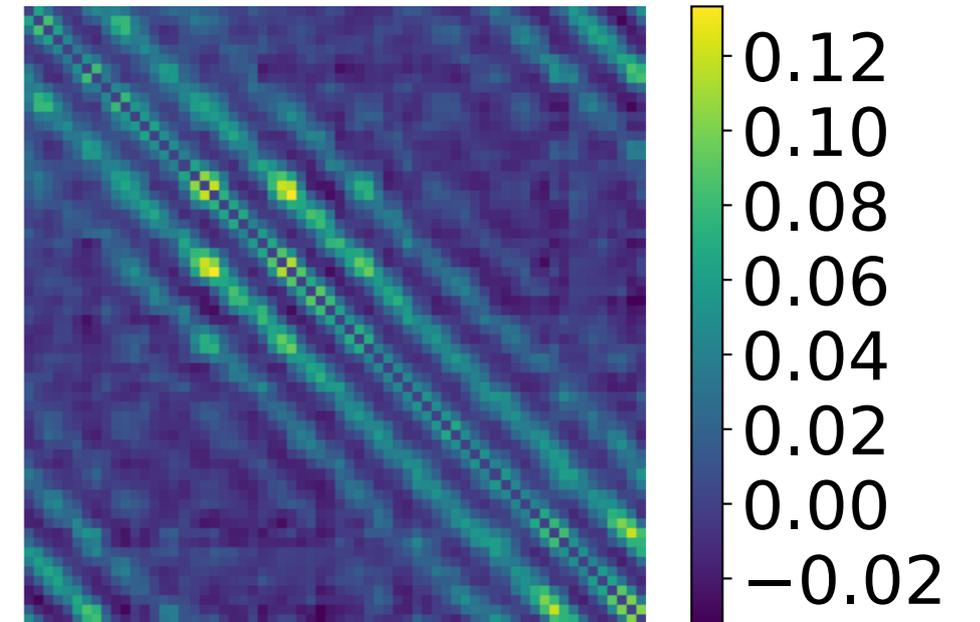
# RBM: Number of training examples

$T = 2.2$



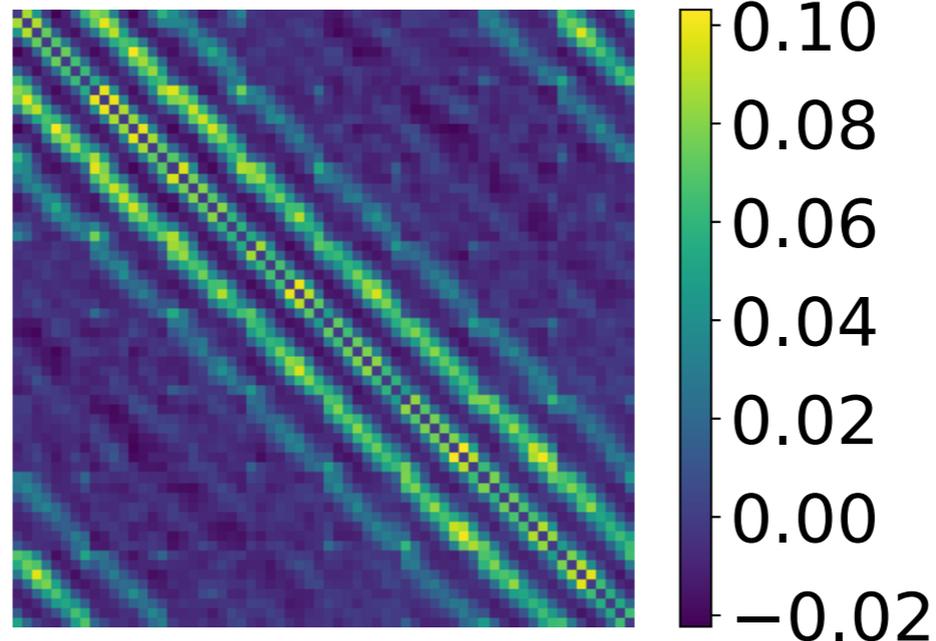
200 Examples

$T = 2.2$



2000 Examples

$T = 2.2$



10000 Examples

# RBM: Lessons Learnt

- Understand well the training criteria from RBMs: Log-likelihood, Loss, free energy, reconstruction error + moments generated by the machine
- RBMs are successful at estimating (higher-order) interactions in a given system of binary variables
- Generally, need lots of training examples

## AND

- Still need to deal with potentially very large numbers of dependent variables (e.g. Gene Networks). **RBM interactions changing depending on gene included!**
- RBMs are not particularly convenient to train ... (e.g. including time on hyper-parameter tuning)

## Part 2

# Interactions: Model-independent definition and estimation

# Defining the target

Aim: Formulate the target quantity of interest:

not as a property of a parametric statistical model

The target quantity can often be identified **without** ever specifying the functional or distributional form of the model: **model-independent**

Why is this important?

**Targeted Learning**

- 1) Be clear about what we are actually after.
- 2) Don't waste computational, analytical and data resources on irrelevant aspects of a problem (here the full joint distribution!)

# There is no “theory” in biology

Come up with a ‘sensible’ model-independent statistical definition

For 1 spin:

$$I_i^m = \ln \left( \frac{p(G_i = 1 \mid \underline{G} = 0)}{p(G_i = 0 \mid \underline{G} = 0)} \right)$$

Here **G** = spins  
Later, **G** = genes

‘odds ratio’: What is the likelihood of spin *i* being 1 vs 0

# There is no “theory” in biology

Come up with a ‘sensible’ model-independent statistical definition

For 1 spin:

$$I_i^m = \ln \left( \frac{p(G_i = 1 \mid \underline{G} = 0)}{p(G_i = 0 \mid \underline{G} = 0)} \right)$$

Here  $\mathbf{G}$  = spins  
Later,  $\mathbf{G}$  = genes

‘odds ratio’: What is the likelihood of spin  $i$  being 1 vs 0

For 2 spins:

$$I_{i,j}^m = \ln \left( \frac{p(G_{ij} = (1, 1) \mid \underline{G} = 0)}{p(G_{ij} = (0, 1) \mid \underline{G} = 0)} \right) - \ln \left( \frac{p(G_{ij} = (1, 0) \mid \underline{G} = 0)}{p(G_{ij} = (0, 0) \mid \underline{G} = 0)} \right)$$

‘odds ratio’ of spin  $i$   
with spin  $j$  being 1

‘odds ratio’ of spin  $i$   
with spin  $j$  being 0

‘generalised odds ratio’: Does the likelihood of spin  $i$  being 1 increase/decrease depending on whether spin  $j$  is 1/0. Generalisable to higher-orders.

# There is no “theory” in biology

Come up with a ‘sensible’ model-independent statistical definition

For 1 spin:

$$I_i^m = \ln \left( \frac{p(G_i = 1 \mid \underline{G} = 0)}{p(G_i = 0 \mid \underline{G} = 0)} \right)$$

Here  $\mathbf{G}$  = spins  
Later,  $\mathbf{G}$  = genes

‘odds ratio’: What is the likelihood of spin  $i$  being 1 vs 0

For 2 spins:

$$I_{i,j}^m = \ln \left( \frac{p(G_{ij} = (1, 1) \mid \underline{G} = 0)}{p(G_{ij} = (0, 1) \mid \underline{G} = 0)} \right) - \ln \left( \frac{p(G_{ij} = (1, 0) \mid \underline{G} = 0)}{p(G_{ij} = (0, 0) \mid \underline{G} = 0)} \right)$$

If two spins are independent:  $p(G_i, G_j \mid \underline{G} = 0) = p(G_i \mid \underline{G} = 0)p(G_j \mid \underline{G} = 0)$

There is no interaction:  $I_{i,j}^m = 0$

# There is no “theory” in biology

Come up with a ‘sensible’ model-independent statistical definition

For 1 spin:

$$I_i^m = \ln \left( \frac{p(G_i = 1 \mid \underline{G} = 0)}{p(G_i = 0 \mid \underline{G} = 0)} \right)$$

Here  $\mathbf{G}$  = spins  
Later,  $\mathbf{G}$  = genes

‘odds ratio’: What is the likelihood of spin  $i$  being 1 vs 0

For 2 spins:

$$I_{i,j}^m = \ln \left( \frac{p(G_{ij} = (1, 1) \mid \underline{G} = 0)}{p(G_{ij} = (0, 1) \mid \underline{G} = 0)} \right) - \ln \left( \frac{p(G_{ij} = (1, 0) \mid \underline{G} = 0)}{p(G_{ij} = (0, 0) \mid \underline{G} = 0)} \right)$$

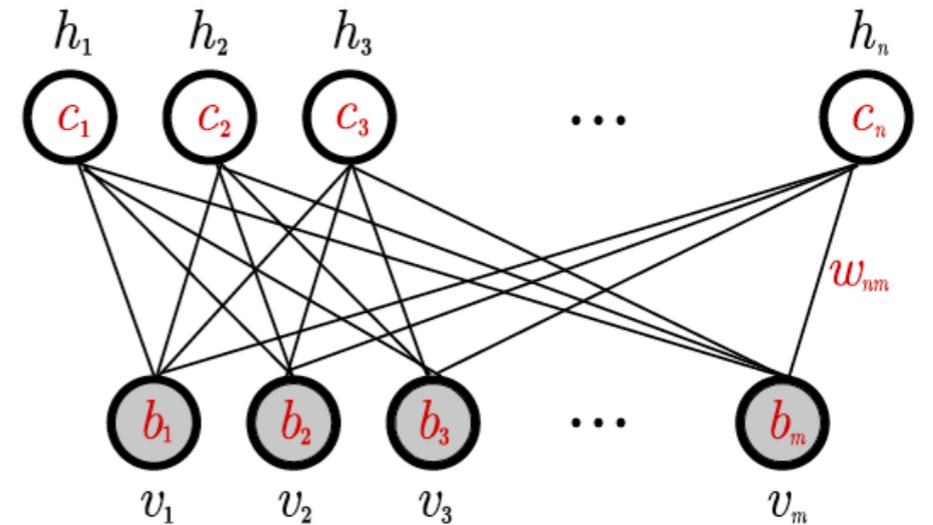
If two spins are independent:  $p(G_i, G_j \mid \underline{G} = 0) = p(G_i \mid \underline{G} = 0)p(G_j \mid \underline{G} = 0)$

There is no interaction:  $I_{i,j}^m = 0$

**Spoiler:  $J_{ij}$  in Ising!**

# Recall analytical formula for RBM interactions

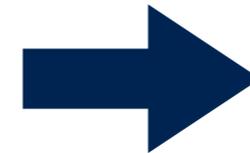
$$E_{\theta}(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{i=1}^n c_i h_i - \sum_{j=1}^m b_j v_j$$



$$p_{\text{RBM}}(\mathbf{v}, \mathbf{h} | \theta) = \frac{1}{Z_{\text{RBM}}} e^{-E_{\theta}(\mathbf{v}, \mathbf{h})}$$

Marginal: 
$$p(\mathbf{v} | \theta) = \frac{1}{Z(\theta)} \prod_{j=1}^m (e^{b_j v_j}) \prod_{i=1}^n \left( 1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j} \right)$$

Asymptotic expansion, resummation, ...



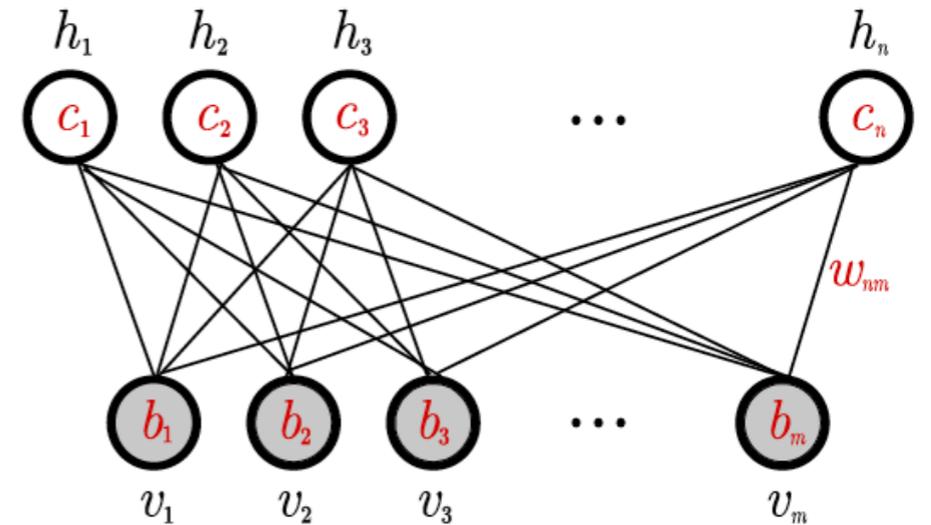
Analytical closed-form expression for n-point interactions, e.g. 2-point:

$$J_{j_1, j_2} \propto \ln \prod_{i=1}^n \frac{(1 + e^{c_i + w_{ij_1} + w_{ij_2}})(1 + e^{c_i})}{(1 + e^{c_i + w_{ij_1}})(1 + e^{c_i + w_{ij_2}})}$$

# Recall analytical formula for RBM interactions

$$E_{\theta}(\mathbf{v}, \mathbf{h}) = - \sum_{i=1}^n \sum_{j=1}^m w_{ij} h_i v_j - \sum_{i=1}^n c_i h_i - \sum_{j=1}^m b_j v_j$$

$$p_{\text{RBM}}(\mathbf{v}, \mathbf{h} | \theta) = \frac{1}{Z_{\text{RBM}}} e^{-E_{\theta}(\mathbf{v}, \mathbf{h})}$$



Marginal: 
$$p(\mathbf{v} | \theta) = \frac{1}{Z(\theta)} \prod_{j=1}^m (e^{b_j v_j}) \prod_{i=1}^n \left( 1 + e^{c_i + \sum_{j=1}^m w_{ij} v_j} \right)$$

Instead, use the TL formulation to directly read-off the coupling!!

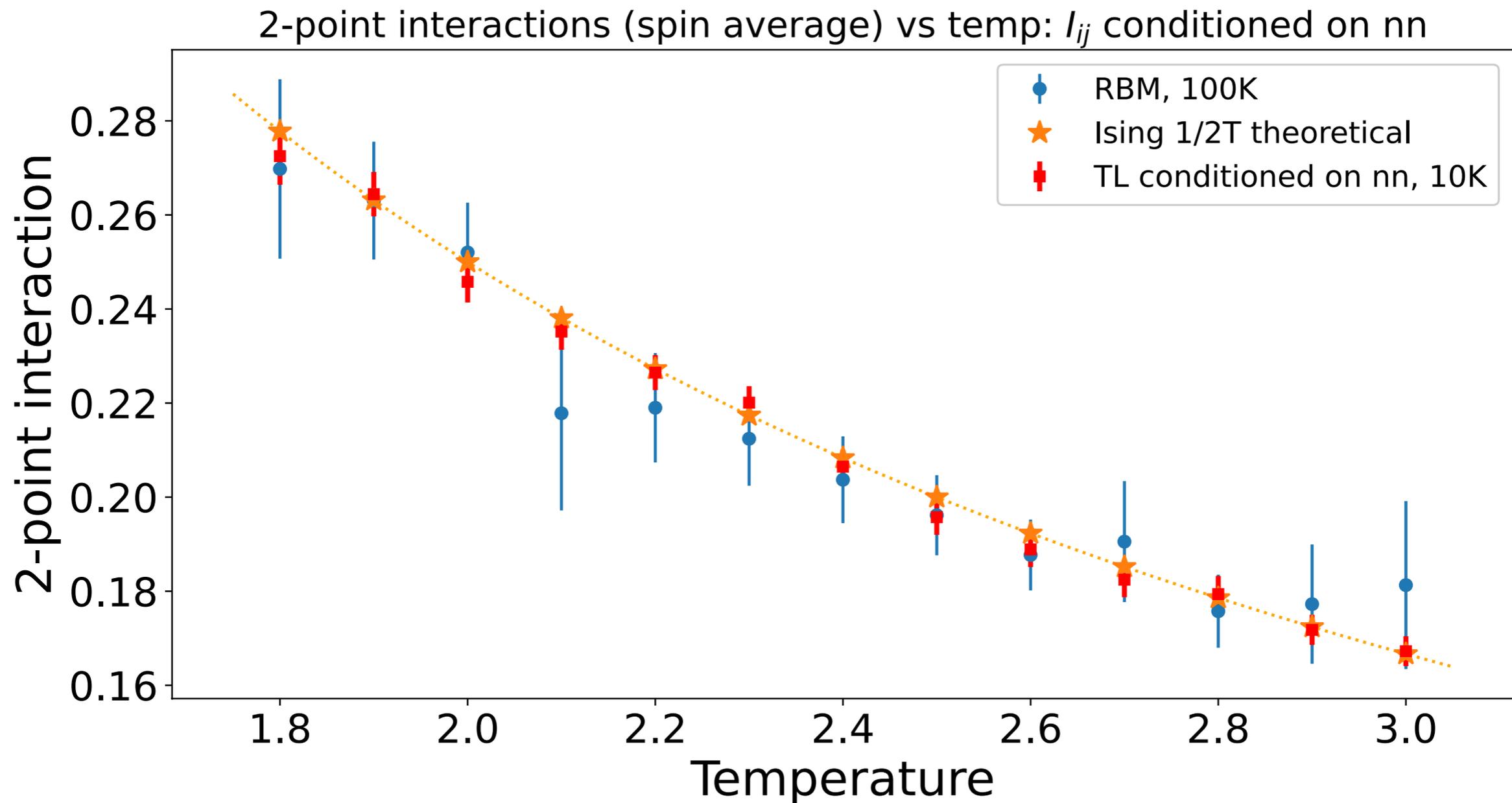
$$I_{j_1, j_2}^m = \frac{p(v_{j_1 j_2} = (1, 1), \underline{v} = 0) p(v_{j_1 j_2} = (0, 0), \underline{v} = 0)}{p(v_{j_1 j_2} = (1, 0), \underline{v} = 0) p(v_{j_1 j_2} = (0, 1), \underline{v} = 0)} = \prod_{i=1}^n \frac{(1 + e^{c_i + w_{ij_1} + w_{ij_2}})(1 + e^{c_i})}{(1 + e^{c_i + w_{ij_1}})(1 + e^{c_i + w_{ij_2}})}$$

No asymptotic expansion and re-summation required ...

Applies to other energy based models

# Model-independent estimation results

Conditioning on parent spins to isolate pairs from the rest of the system (Markovian). Run time: Few seconds per temperature.



10K samples

# Biology: Large number of dependent variables

Estimating intricate interaction structure amongst many genes

Certain approximation no longer possible:  $p(G_i, G_j) \neq p(G_i)p(G_j)$

Number of variables  $\gg$  data, (and high temperatures) **G binarised!**

$$I_{i,j}^m = \ln \left( \frac{p(G_{ij} = (1, 1) \mid \underline{G} = 0) p(G_{ij} = (0, 0) \mid \underline{G} = 0)}{p(G_{ij} = (0, 1) \mid \underline{G} = 0) p(G_{ij} = (1, 0) \mid \underline{G} = 0)} \right)$$

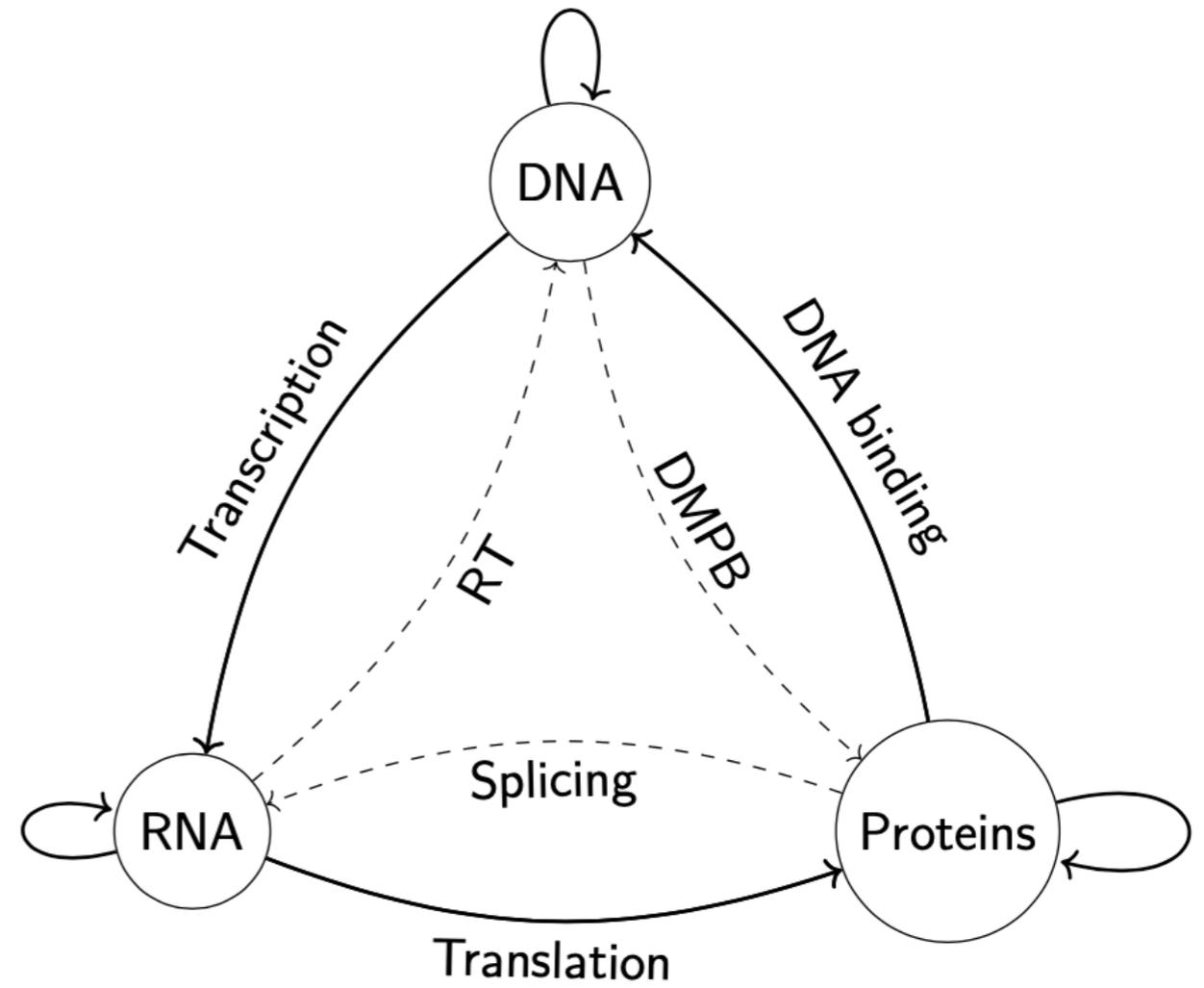
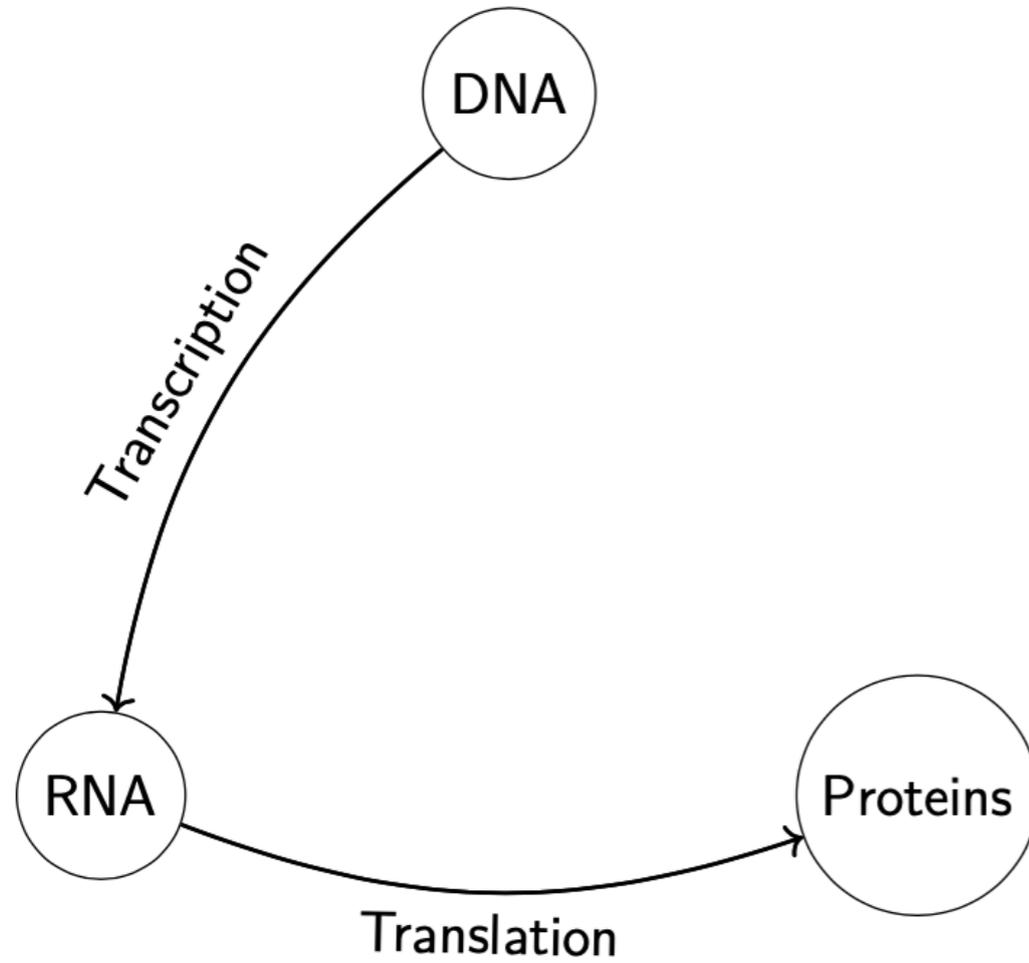
**Estimate** conditional dependencies directly from data, using efficient causal discovery algorithms (e.g. PC, Score-based MCMC)

Nothing comes for free! These come with their own assumptions/bias  
Keep in mind to be conservative.

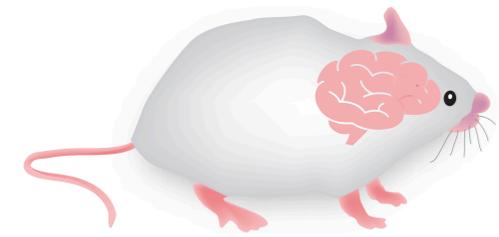
## Part 3

# Biological data: Gene expression

# The central dogma



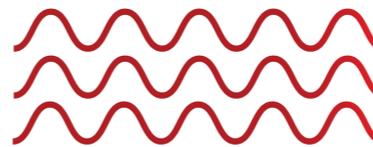
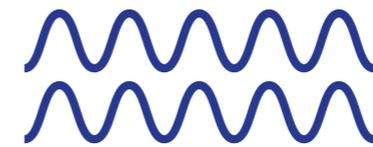
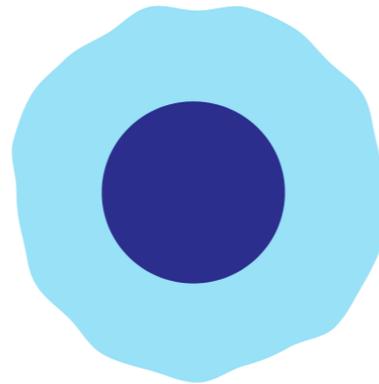
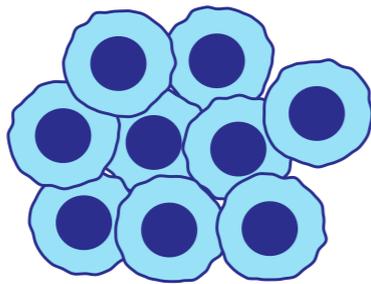
# Biology: Large number of dependent variables



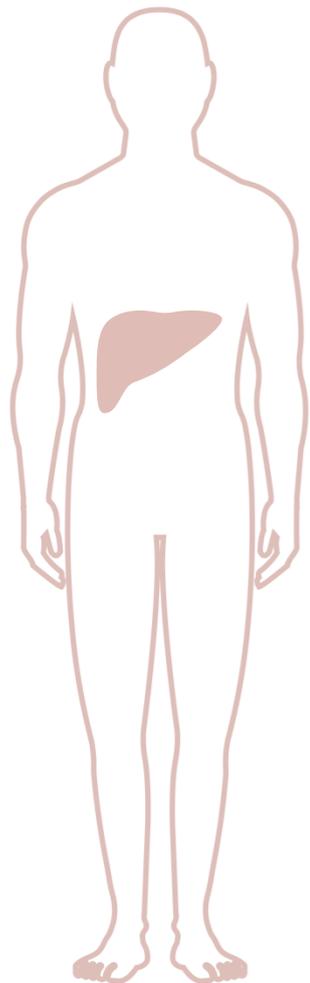
Single Cell

RNA

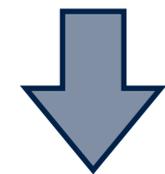
Sequencing



Illumina



Count data  
(Later binarised)



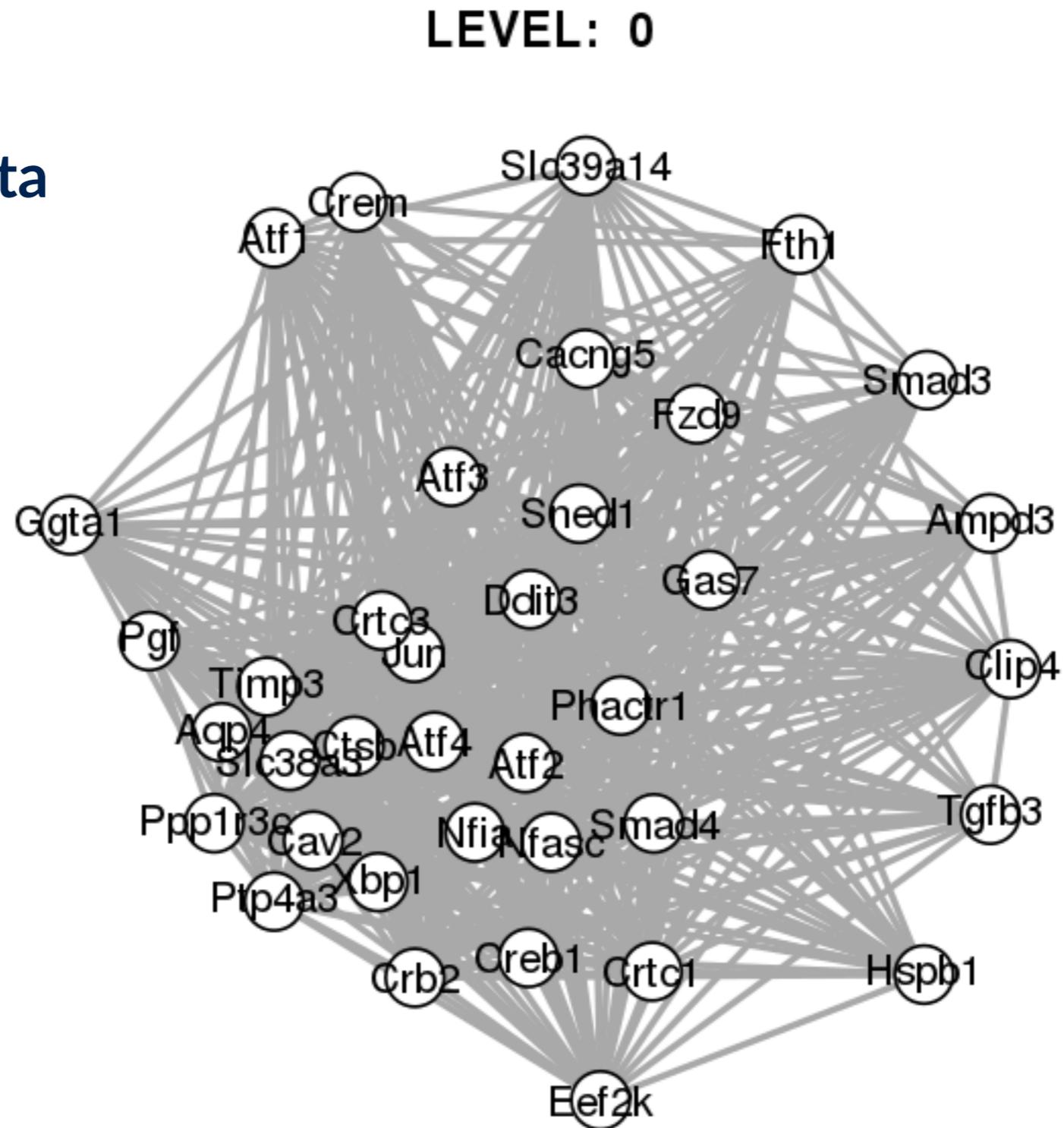
Gene

Cell

NxG

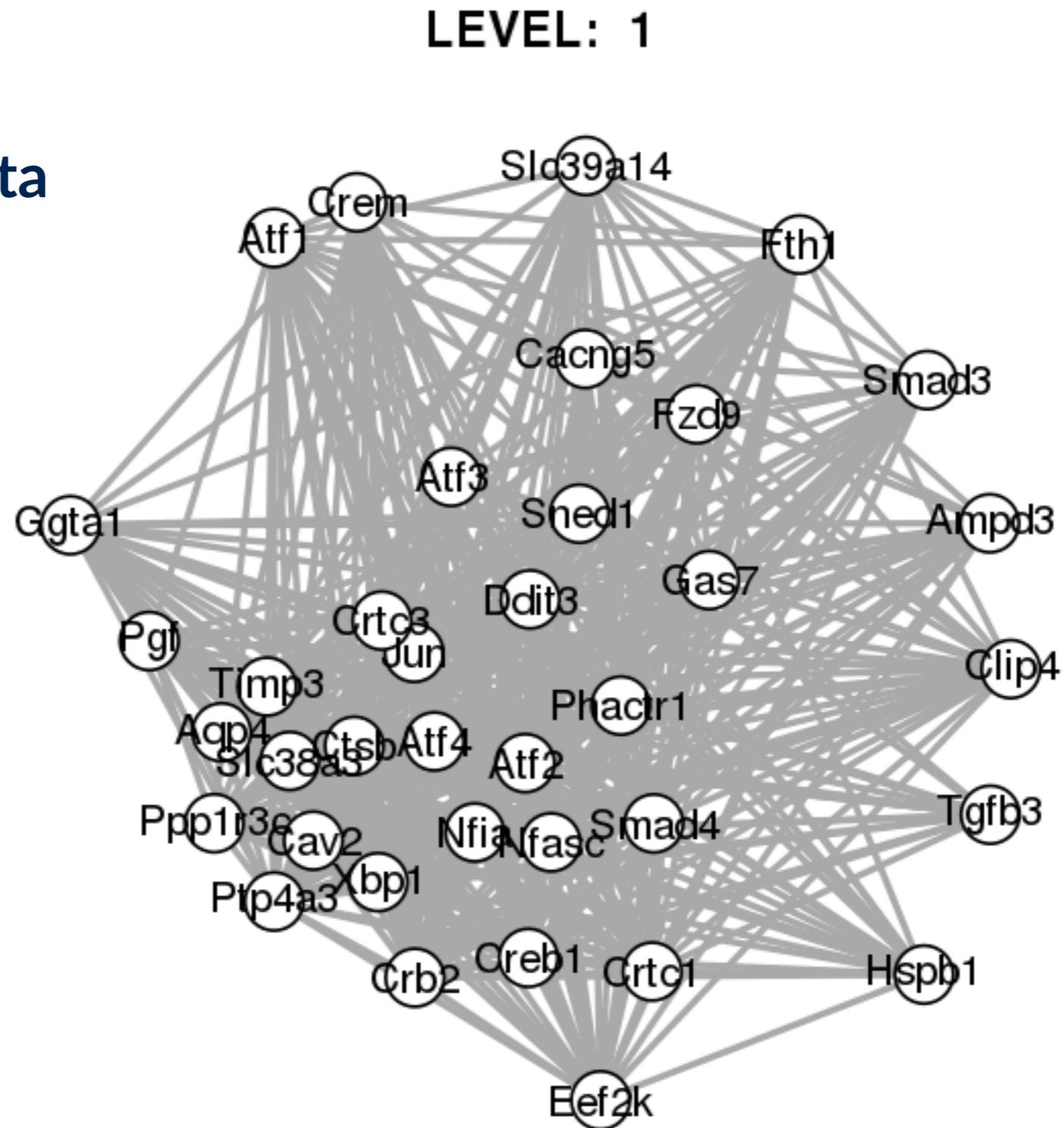
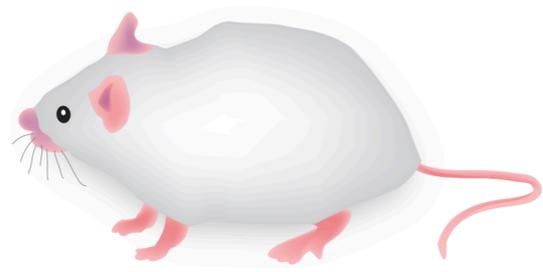
# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data



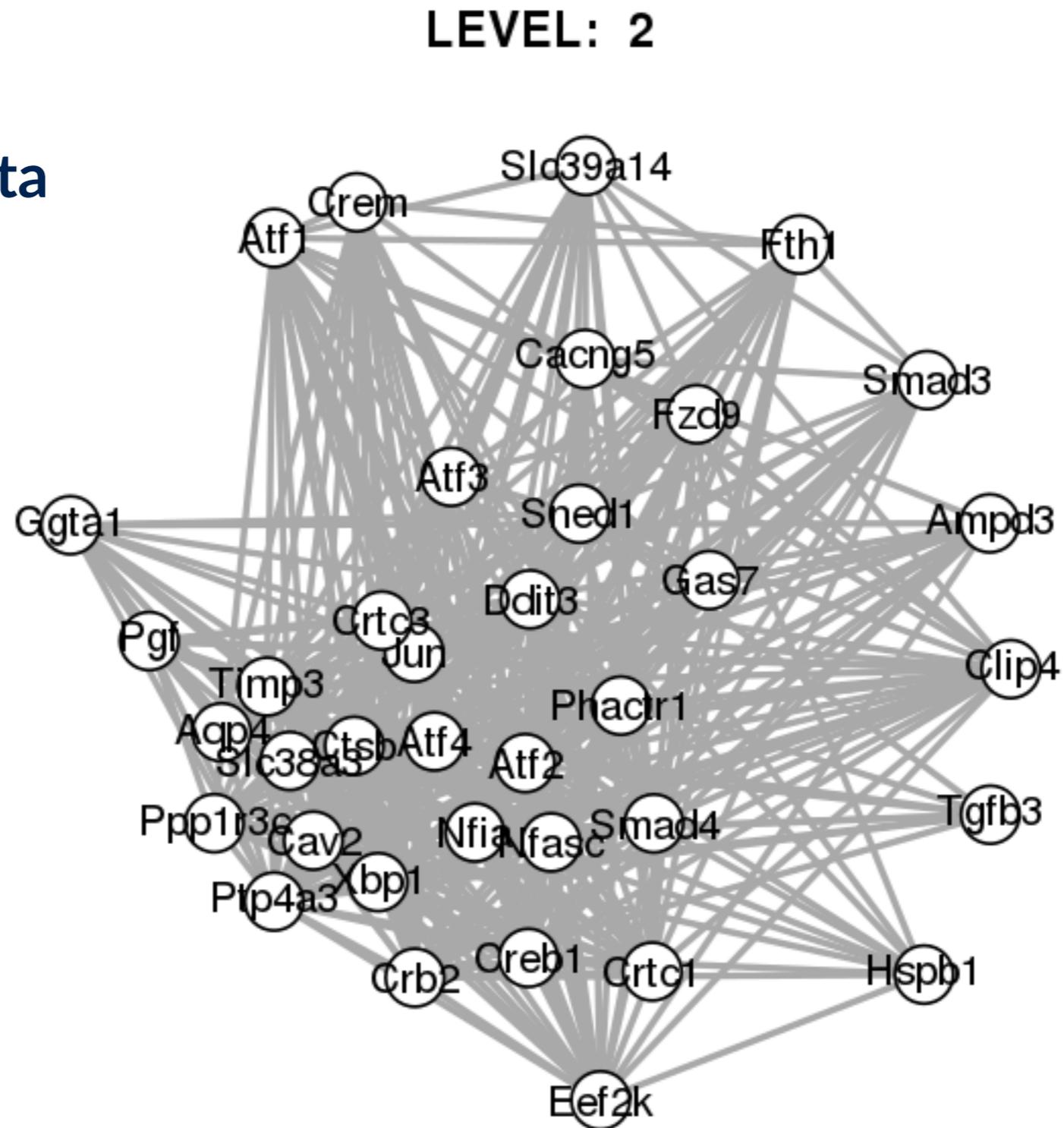
# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data



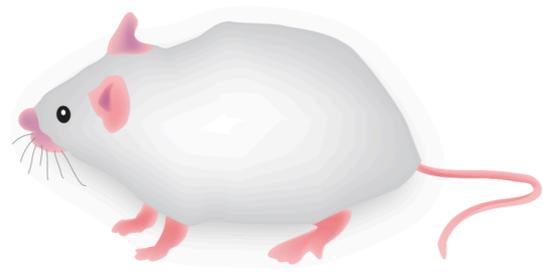
# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data

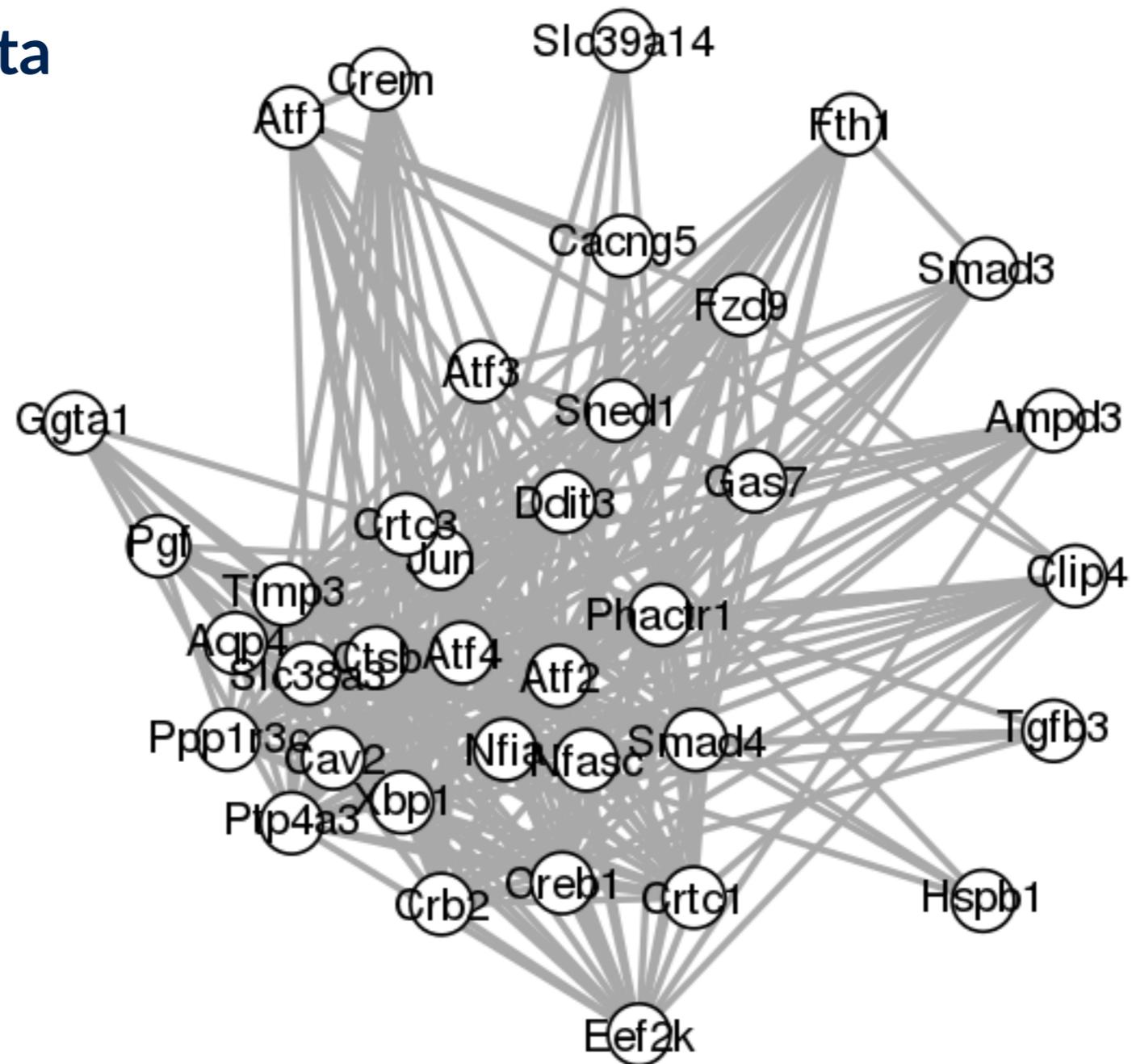


# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data

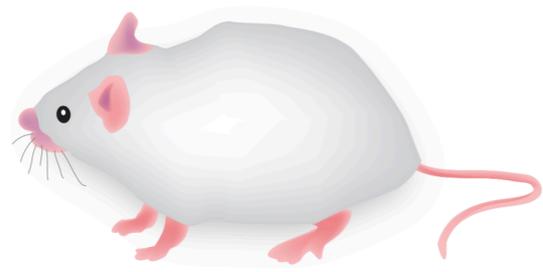


LEVEL: 3

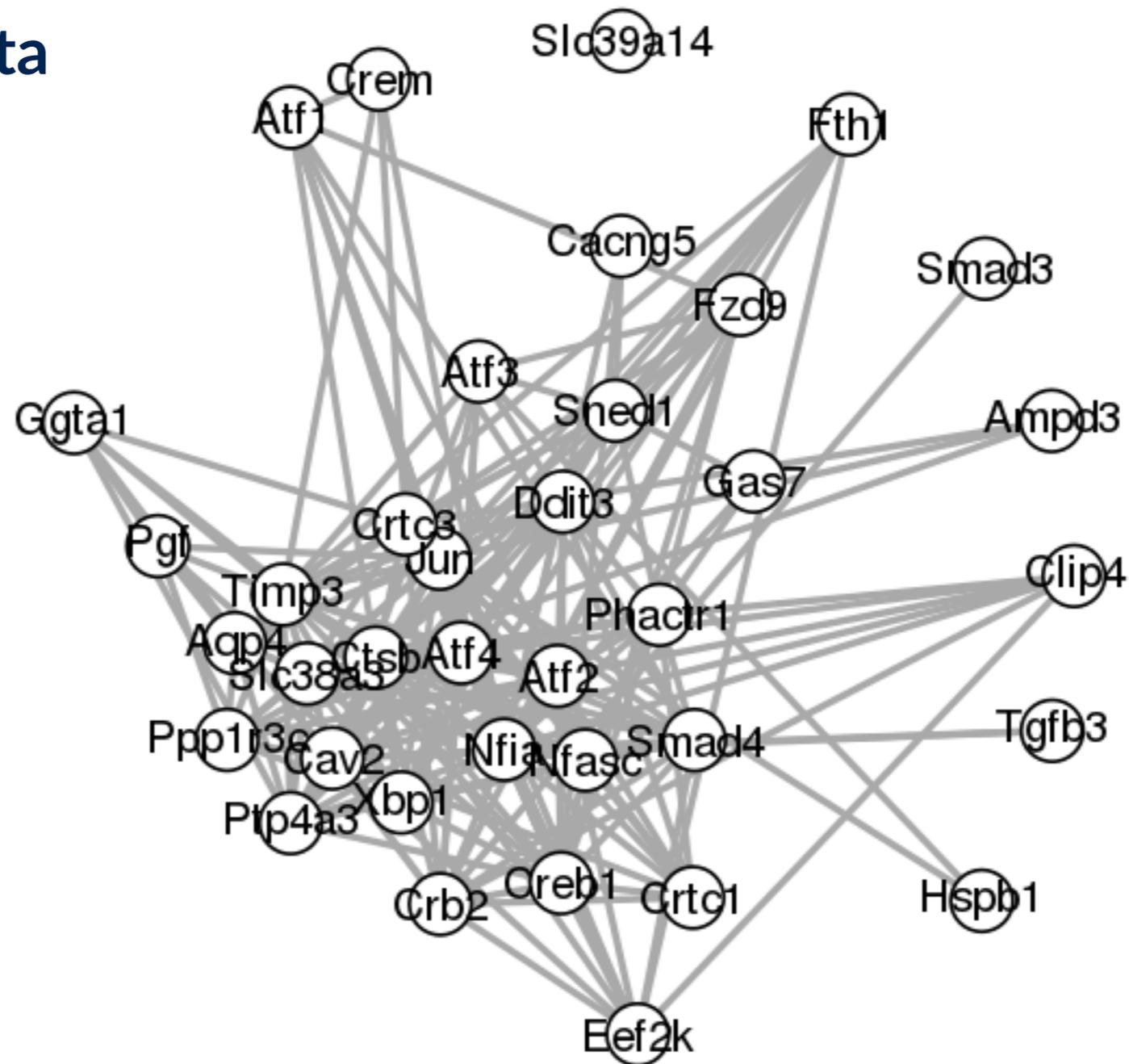


# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data

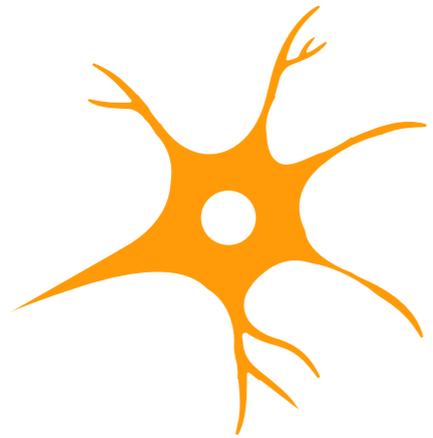
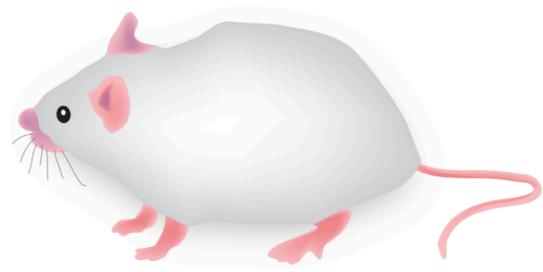


LEVEL: 4

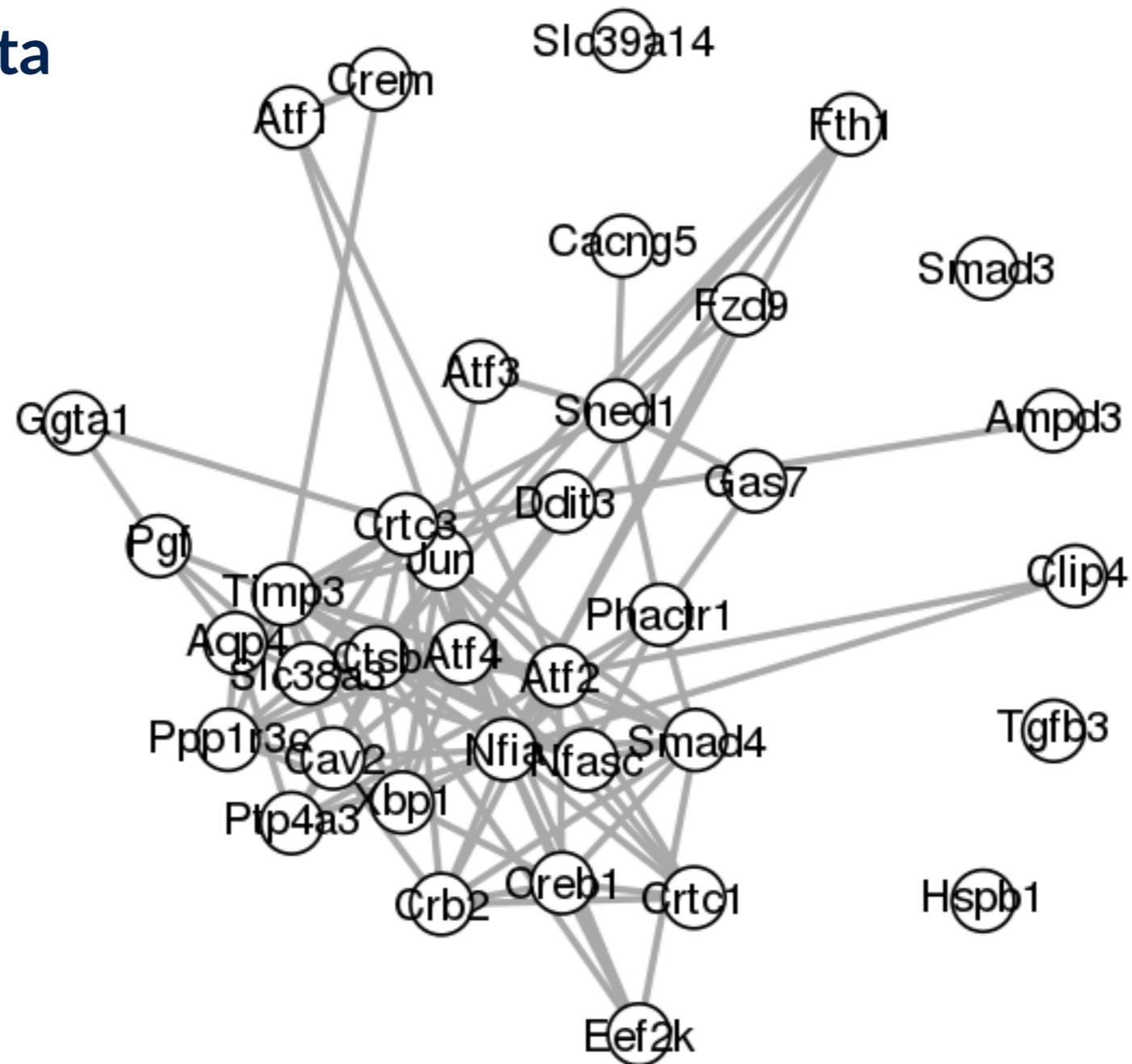


# (In)dependence of gene expressions

10X single-cell  
1M mouse brain  
developmental data



LEVEL: 5

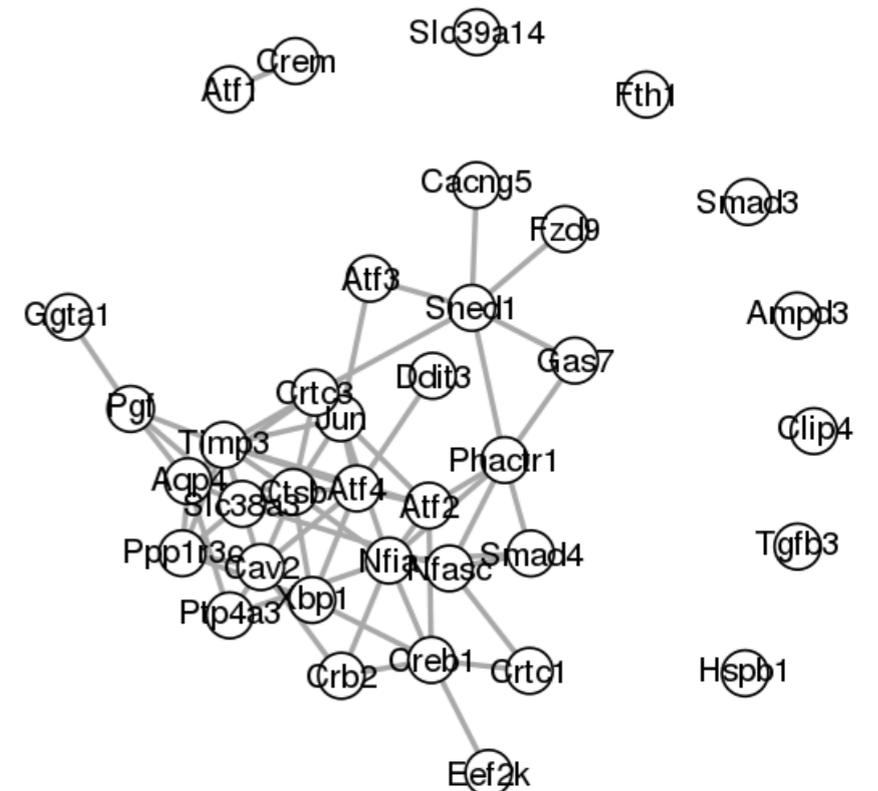




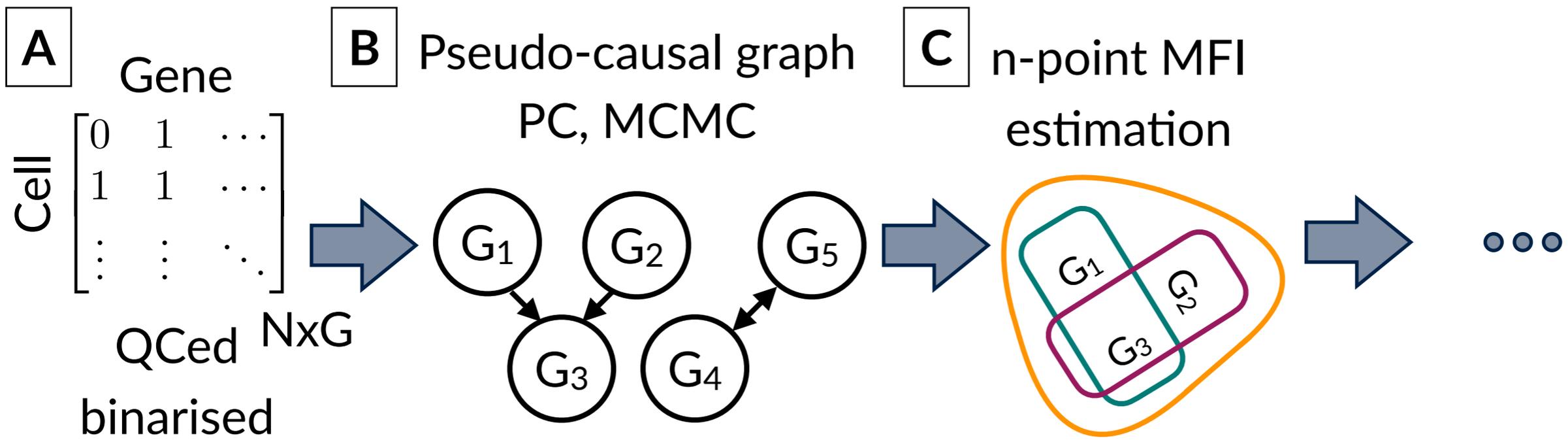
# Estimate model-free n-point interactions

$$I_{i,j}^m = \ln \left( \frac{p(G_{ij} = (1, 1) \mid \underline{G} = 0) p(G_{ij} = (0, 0) \mid \underline{G} = 0)}{p(G_{ij} = (0, 1) \mid \underline{G} = 0) p(G_{ij} = (1, 0) \mid \underline{G} = 0)} \right)$$

2-point up to 7-point interactions



# Stator workflow



**What is the biological interpretation of these n-point interactions?**

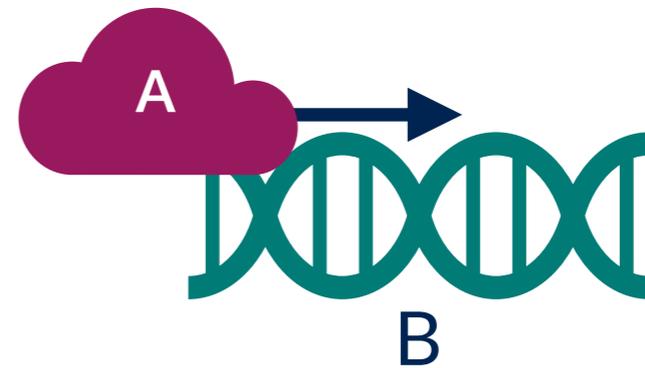
**Does it help answer questions of molecular & cell biologists?**

# Regulation vs Cell State

1. MFIs represent physical interactions amongst molecules **within** a cell



2. MFIs represent a biochemical network:  
Transcription factor A -> Target gene B

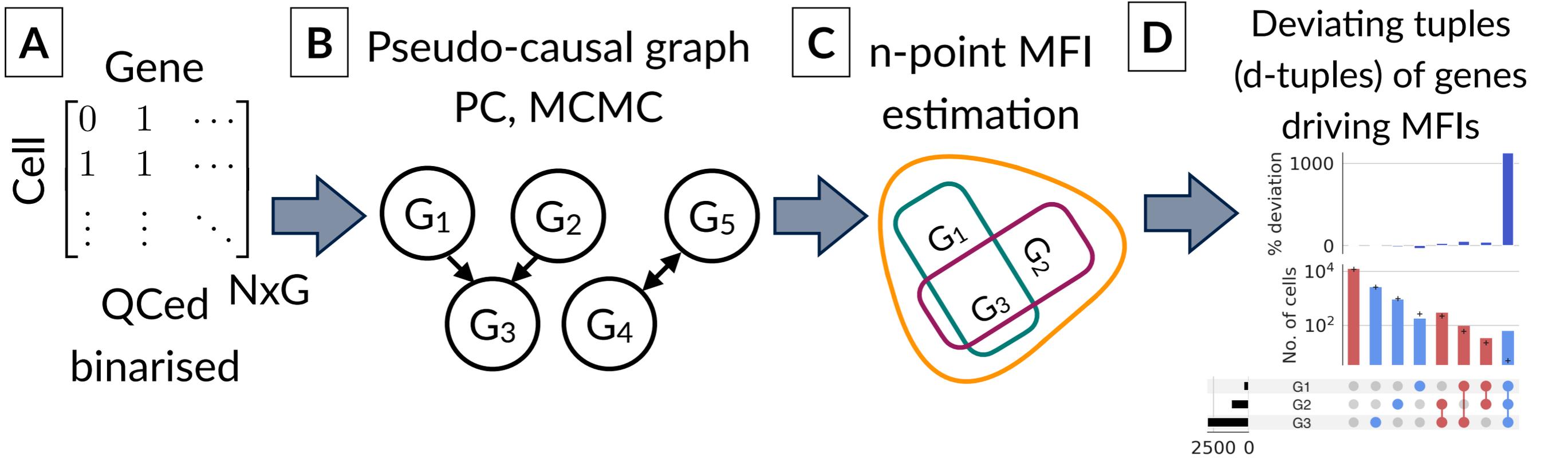


3. MFIs represent dependence structure amongst genes that imply **cell types, subtypes or states**

(MFIs estimated as an average across **diverse** cell populations)

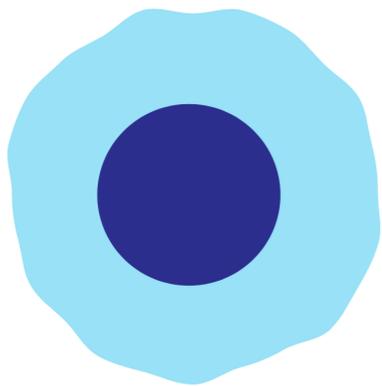
i.e. the statistical interpretation of interactions, rather than dynamical/physical

# Stator workflow



d-tuple :  $(G_1, G_2, G_3) = (1, 1, 1)$

For each cell



if  $(G_1, G_2, G_3) = (1, 1, 1) \Rightarrow 1$

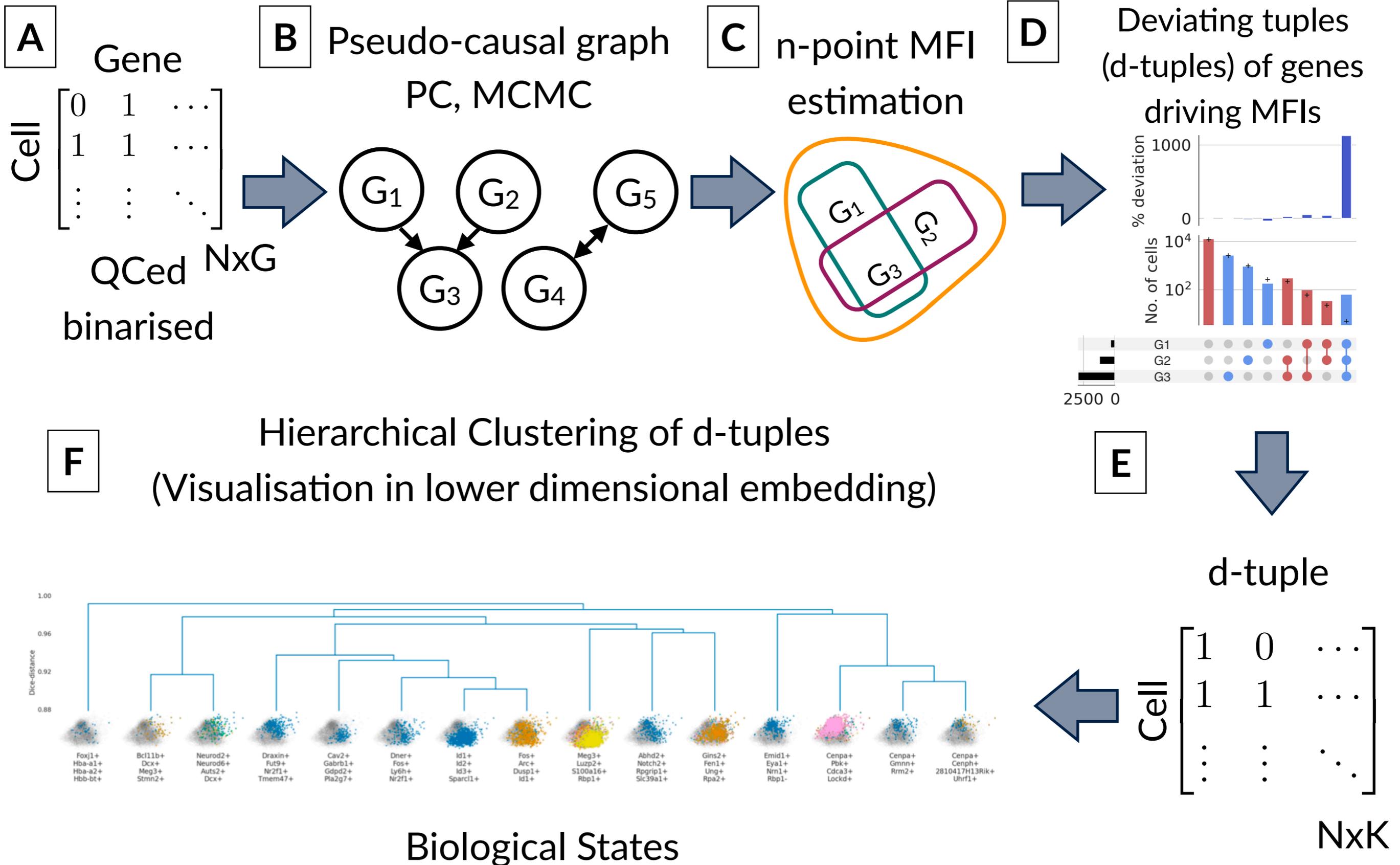
if  $(G_1, G_2, G_3) \neq (1, 1, 1) \Rightarrow 0$

**E**

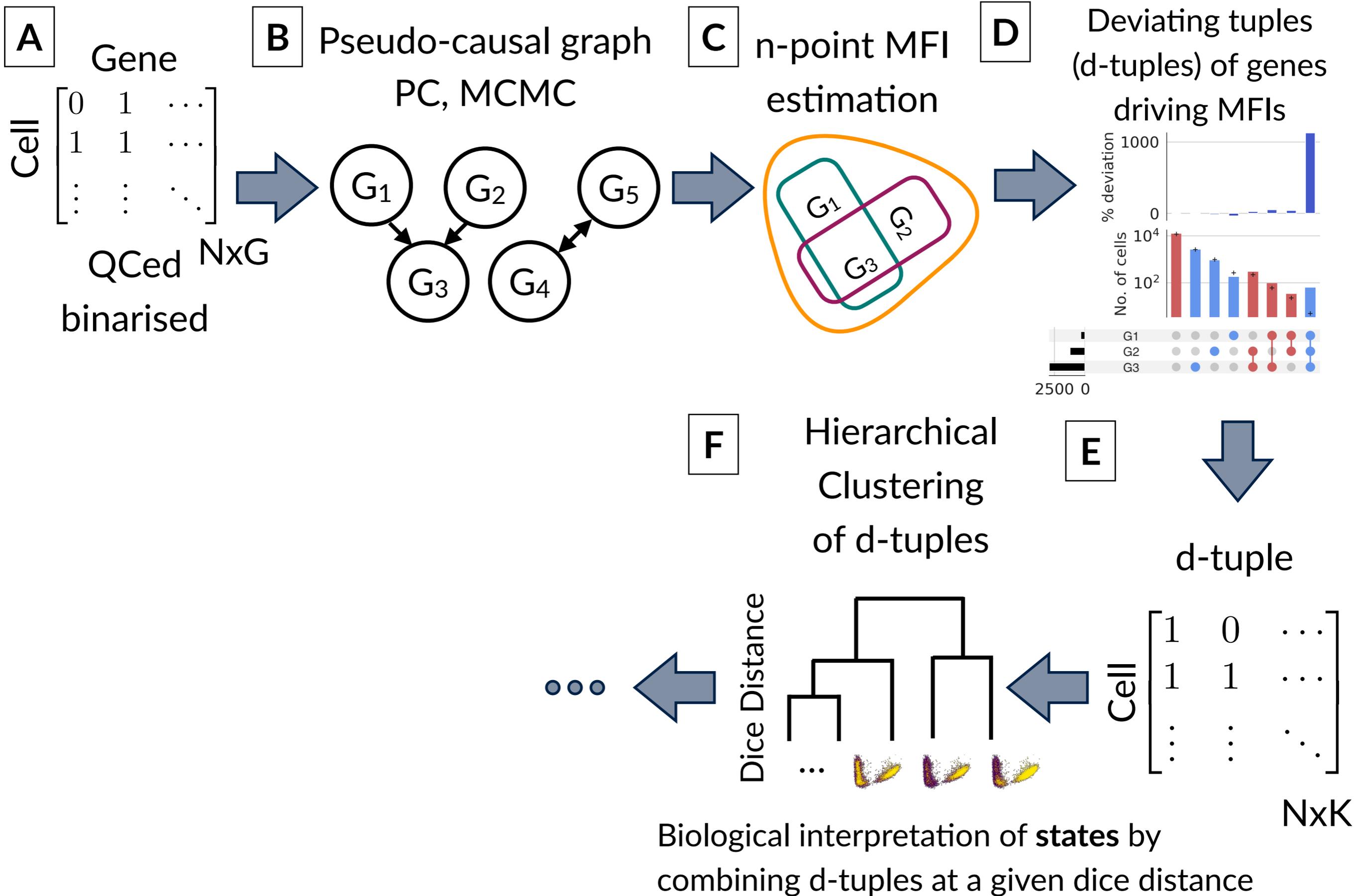
d-tuple

Cell  $\begin{bmatrix} 1 & 0 & \dots \\ 1 & 1 & \dots \\ \vdots & \vdots & \ddots \end{bmatrix}$   
 NxK

# Stator workflow



# Stator workflow



**Will biologists bother using this methodology?**

# Criteria:

1. Ease of use  
(easy-to-follow documentation, easy copy/paste code, press of a button)
2. Speed
3. Output with good visualisation
4. Biology that “makes sense”

...

...

...

# Criteria:

1. Ease of use

(easy-to-follow documentation, easy copy/paste code, press of a button)

2. Speed

3. Output with good visualisation

4. Biology that “makes sense”

...

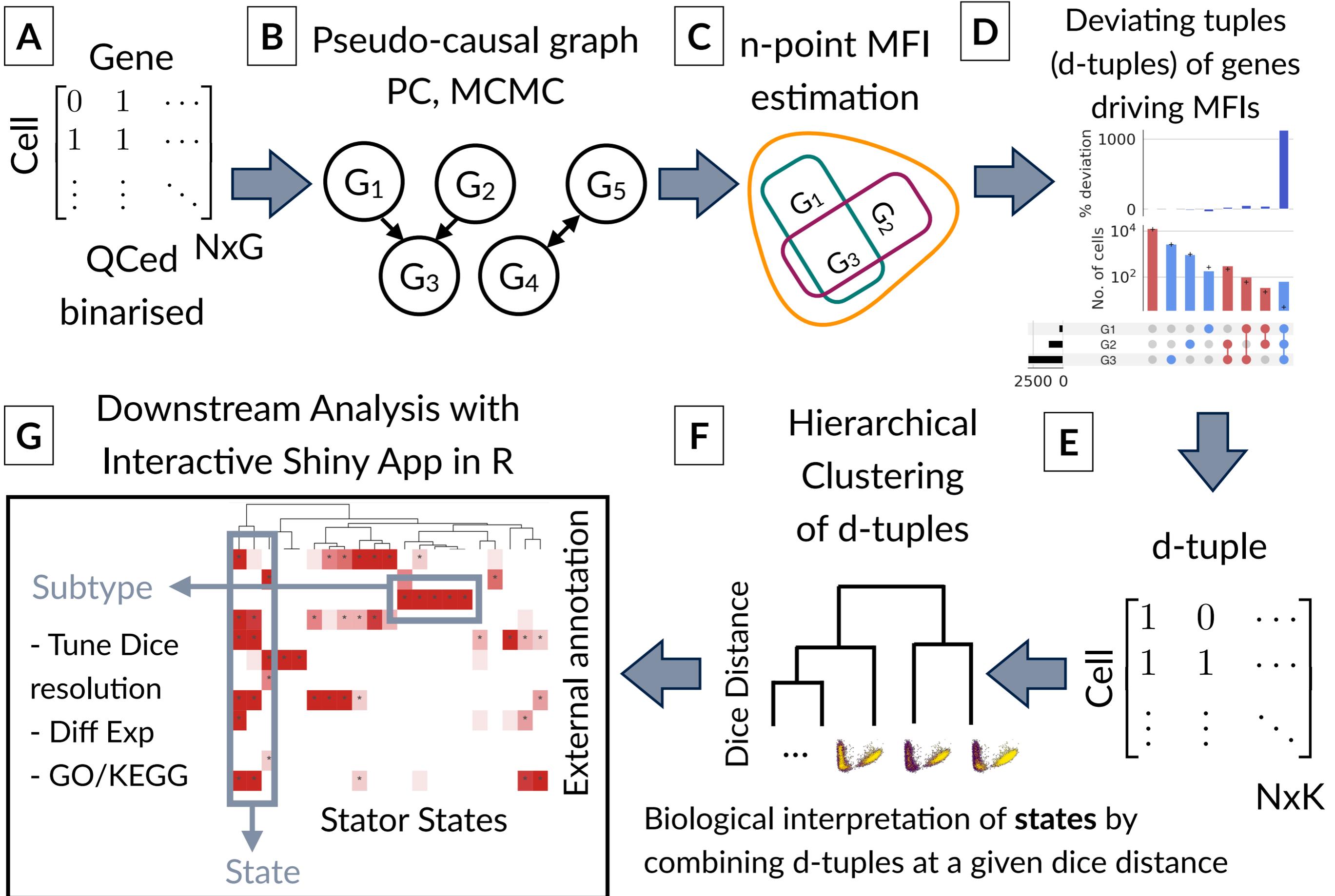
...

...

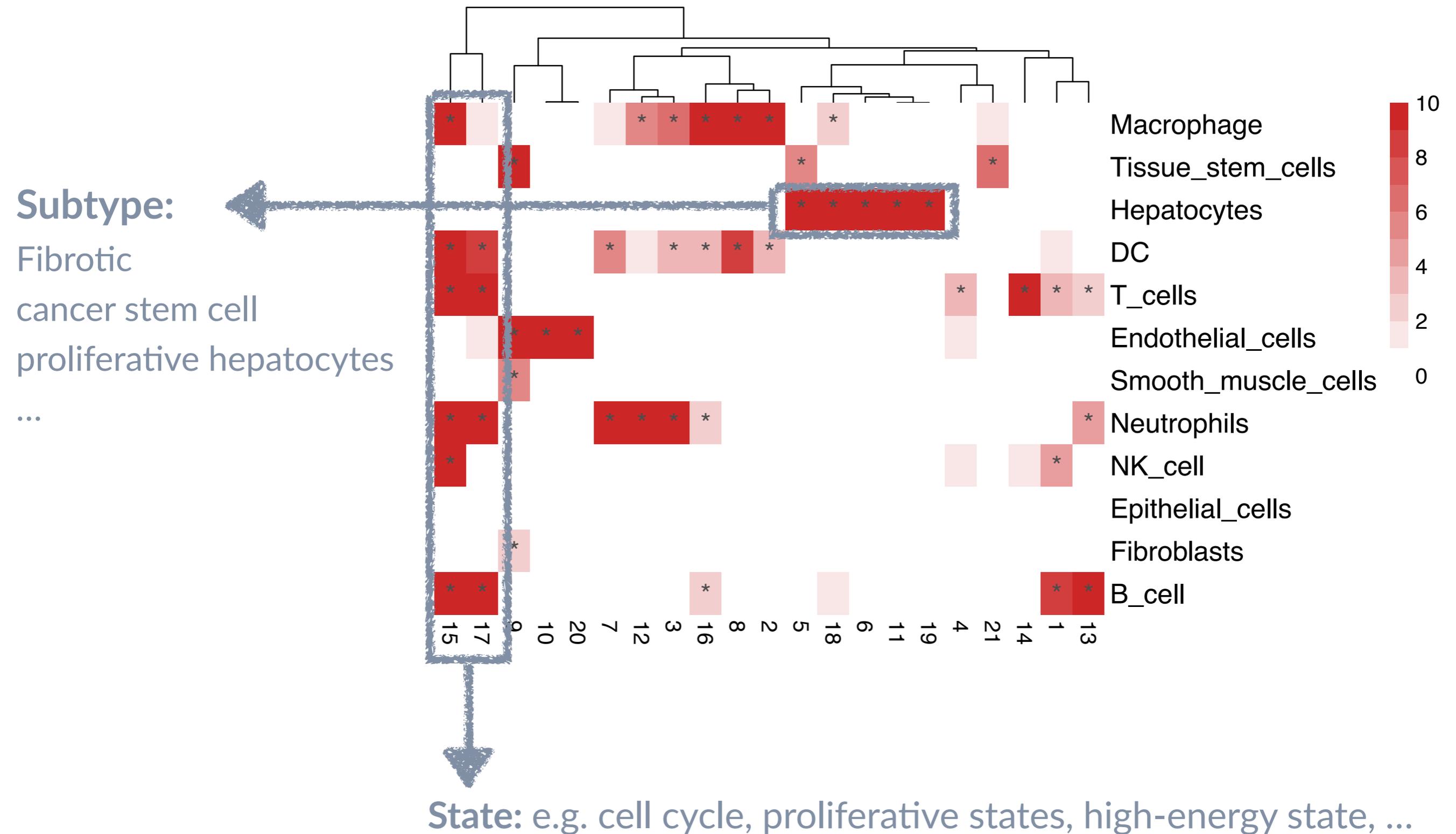
**n. The methodology behind the software**

(where n can be very large!)

# Stator workflow



# Liver Cancer: Cell types and states



# Stator app interface



About

Table

Heatmap

GO & KEGG

Using rrvgo

Upset Plot

DE analysis

## Explore cell states by MFIs

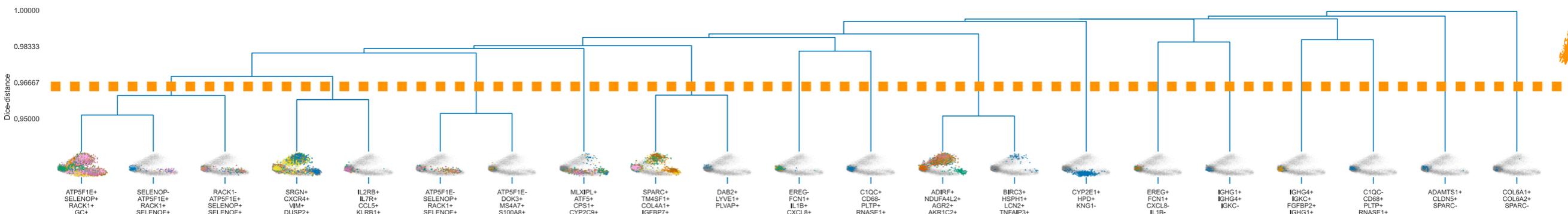
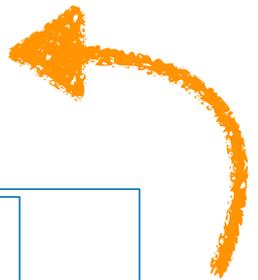
**MFIs** takes in scRNA-seq count matrix and estimate gene interactions. Here we show how to use these MFIs to explore cell states.

### Data Visualization & Analysis

- Table - A Summary statistics for deviating state
- Heatmaps - Over-representation test for MFIs and other cell annotations
- GO & KEGG for genes in each state
- rrvgo - Simplifying the redundance of GO sets
- Upset Plot
- DE analysis for mutually exclusive states

### Tutorial

Resolution (dice distance) cut-off



# Stator app interface

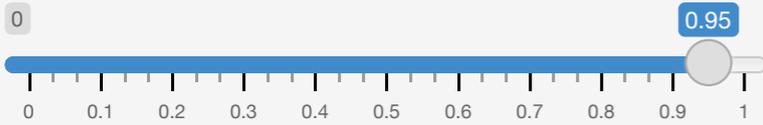


THE UNIVERSITY of EDINBURGH



Here we use scRNA-seq HCC dataset. To upload your data, click the box:

Dice distance:



Submit



Resolution cut-off

- About
- Table**
- Heatmap
- GO & KEGG
- Using rrvgo
- Upset Plot
- DE analysis

## Table for deviating MFIs

391 deviating MFIs in total, 21 clusters.

Show  entries

Search:

	genes	state	dev	pval	cluster
47	IGHG4_IGHG1_IGKC	111	216.09886998247	2.29736631083133e-28	Cluster:1
150	IGHG4_FGFBP2_IGKC	111	21.9122324209873	1.04827895292059e-19	Cluster:1
341	PLTP_CD68_C1QC_RNASE1	1101	6.41225201584546	1.39266594426104e-9	Cluster:2
213	EREG_CXCL8_IL1B_FCN1	1011	14.7648630077289	1.05933261179815e-8	Cluster:3
220	EREG_CXCL8_IL1B_FCN1	1101	13.9465435194403	3.12038785771356e-13	Cluster:3
282	EREG_IL1B_FCN1	101	9.95679034074909	2.46844896187663e-13	Cluster:3
292	EREG_CXCL8_FCN1	101	9.09789155019342	9.41630198774148e-13	Cluster:3
177	IGHG4_IGHG1_IGKC	110	18.5991179035581	8.521545688936e-16	Cluster:4
63	CXCL1_IER3_CXCL3	111	99.0782359764262	2.01516836396423e-47	Cluster:5
74	MB_NDUFA4L2_TFF2	111	79.5037037037037	6.82448451177721e-18	Cluster:5

Showing 1 to 10 of 391 entries

Previous  2 3 4 5 ... 40 Next

Download as .csv

# Stator app interface



Medical Research Council



THE UNIVERSITY of EDINBURGH



CANCER RESEARCH UK

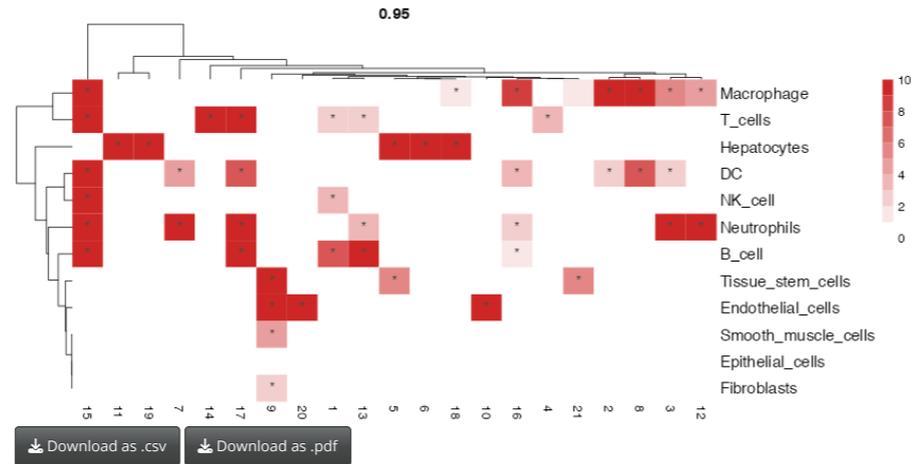
Here we use scRNA-seq HCC dataset. To upload your data, click the box:

Dice distance:

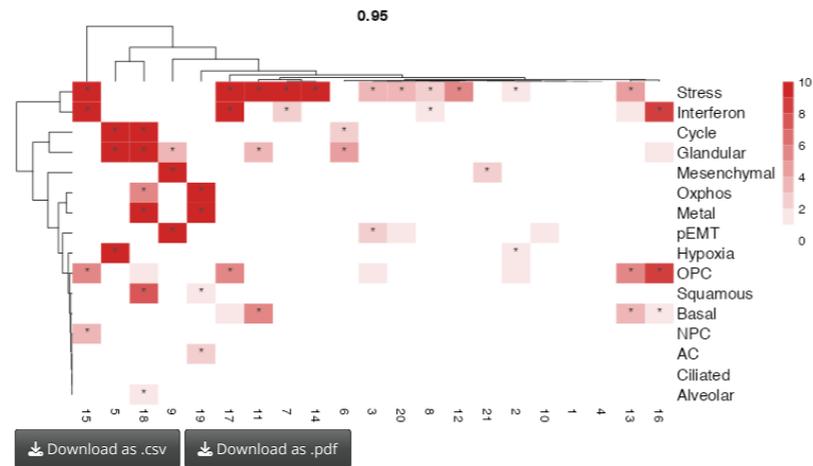
0 0.1 0.2 0.3 0.4 0.5 0.6 0.7 0.8 0.9 1

[About](#) [Table](#) [Heatmap](#) [GO & KEGG](#) [Using rvgo](#) [Upset Plot](#) [DE analysis](#)

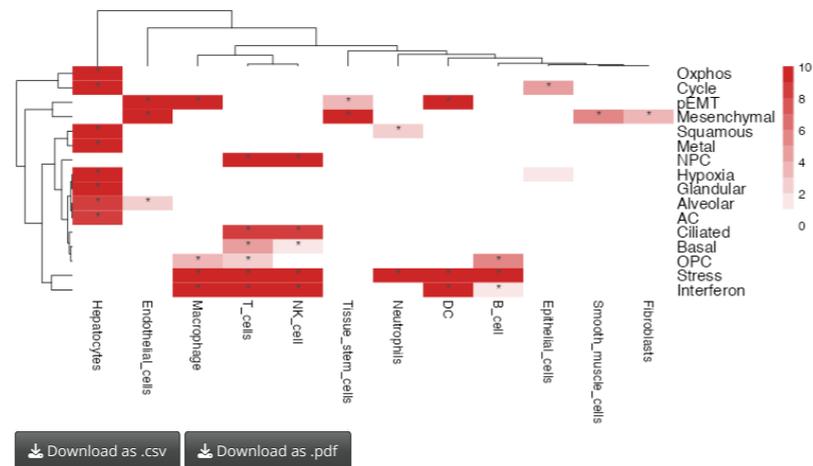
Heatmap: Cell types (clustering + singleR)



Heatmap: Cell states (NMF)



Heatmap: Cell states vs. Cell types

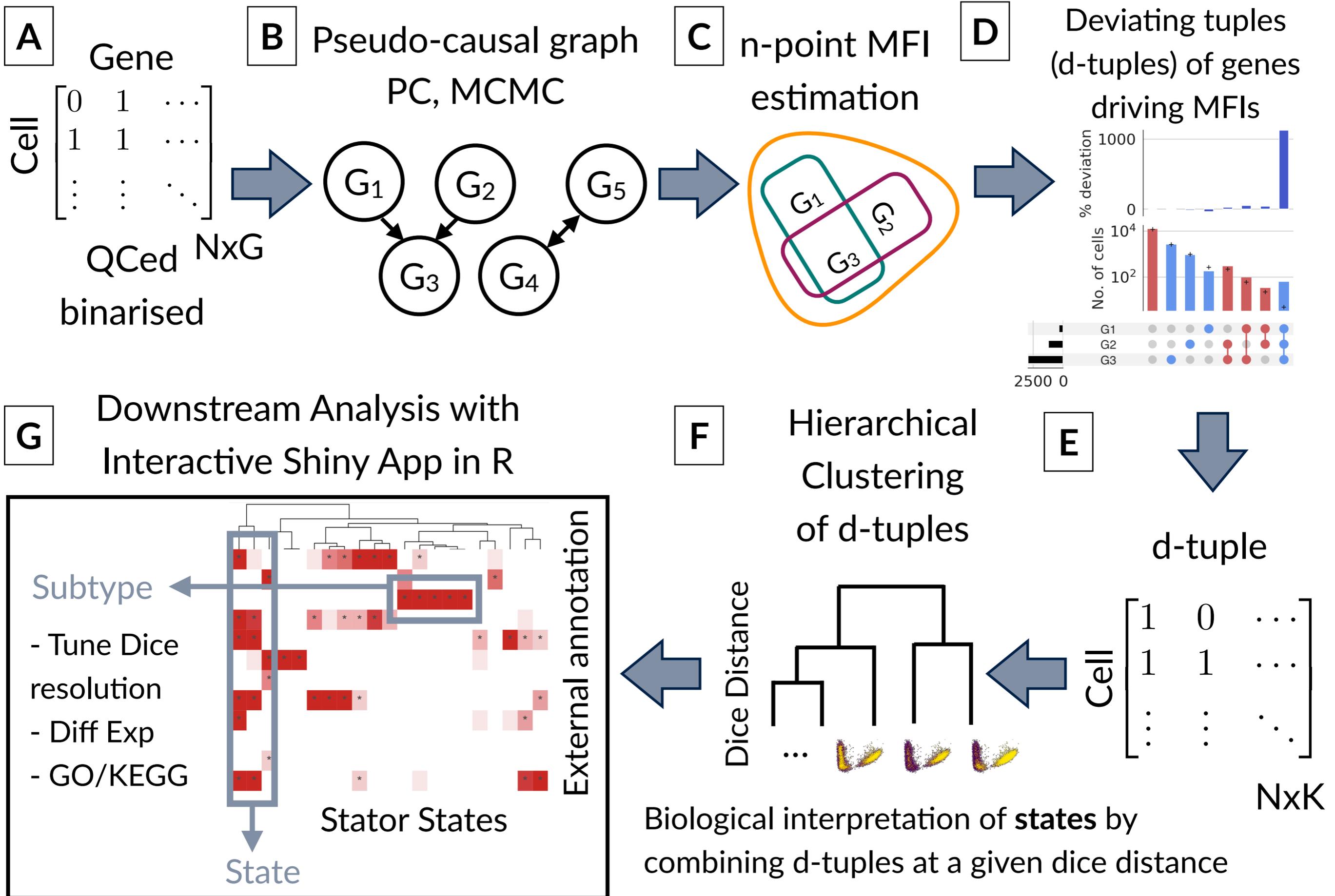


Compare to known cell types from the literature

Compare to other methods used in the literature, e.g. NMF

Compare to expert annotated cell types

# Stator workflow



**Backup slides**

# RBM Prediction: n-point interactions

- A non pair-wise treatment
- Higher order couplings
- Not accessible via standard statistical techniques

$$E(\mathbf{v}) = - \sum_j b_j v_j - \sum_j \left( \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2} \sum_{jk} \underbrace{\left( \sum_i \kappa_i^{(2)} W_{ik} W_{ij} \right)} v_j v_k + \dots$$

Re-sum the entire series to obtain 2-point coupling!!

# Derivation of n-point interactions in closed form

$$\begin{aligned} E(\mathbf{v}) &= \ln \sum_{\mathbf{h}} e^{E(\mathbf{v}, \mathbf{h})} \\ &= \ln \sum_{\mathbf{h}} e^{-\sum_j b_j v_j - \sum_i c_i h_i - \sum_{i,j} h_i W_{ij} v_j} \end{aligned}$$

$$\begin{aligned} E(\mathbf{v}) &= -\sum_j b_j v_j - \sum_i \ln \sum_{h_i} e^{c_i h_i} e^{\sum_j h_i W_{ij} v_j} \\ &= -\sum_j b_j v_j - \sum_i \ln \sum_{h_i} q(h_i) e^{t h_i}, \end{aligned}$$

$$t \equiv \sum_j W_{ij} v_j \text{ and } q(h_i) \equiv e^{c_i h_i}$$

**Cumulant generating function:**

$$K_i(t) \equiv \ln \sum_{h_i} q(h_i) e^{t h_i} = \sum_n \frac{\kappa_i^{(n)} t^n}{n!}$$

$$\kappa_i^{(n)} = \partial_t^n K_i(t) |_{t=0}$$

# Derivation of n-point interactions in closed form

$$\begin{aligned}
 E(\mathbf{v}) &= - \sum_j b_j v_j - \sum_i \kappa_i^{(0)} - \sum_i \kappa_i^{(1)} t - \sum_i \frac{\kappa_i^{(2)} t^2}{2!} - \dots \\
 &= - \sum_i \kappa_i^{(0)} - \sum_j \left( b_j + \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2!} \sum_{j_1, j_2} \left( \sum_i \kappa_i^{(2)} W_{ij_1} W_{ij_2} \right) v_{j_1} v_{j_2} - \dots
 \end{aligned}$$

$$v_j^n = v_j \quad , \quad n \in \mathbb{Z}^+$$

e.g. 2-point interaction:

$$\sum_{n>1} \frac{1}{2(n!)} \sum_{0<k<n} \sum_{j_1 \neq j_2} \left( \sum_i \kappa_i^{(n)} \binom{n}{k} W_{ij_1}^k W_{ij_2}^{n-k} \right) v_{j_1} v_{j_2}$$

...

$$H_{j_1 j_2} = \frac{1}{8} \sum_i \ln \frac{(1 + e^{c_i + W_{ij_1} + W_{ij_2}})(1 + e^{c_i})}{(1 + e^{c_i + W_{ij_1}})(1 + e^{c_i + W_{ij_2}})}$$

**Closed form expression!**

# Derivation of n-point interactions in closed form

$$\begin{aligned}
 E(\mathbf{v}) &= - \sum_j b_j v_j - \sum_i \kappa_i^{(0)} - \sum_i \kappa_i^{(1)} t - \sum_i \frac{\kappa_i^{(2)} t^2}{2!} - \dots \\
 &= - \sum_i \kappa_i^{(0)} - \sum_j \left( b_j + \sum_i \kappa_i^{(1)} W_{ij} \right) v_j - \frac{1}{2!} \sum_{j_1, j_2} \left( \sum_i \kappa_i^{(2)} W_{ij_1} W_{ij_2} \right) v_{j_1} v_{j_2} - \dots
 \end{aligned}$$

$$v_j^n = v_j \quad , \quad n \in \mathbb{Z}^+$$

e.g. 2-point interaction:

$$\sum_{n>1} \frac{1}{2(n!)} \sum_{0<k<n} \sum_{j_1 \neq j_2} \left( \sum_i \kappa_i^{(n)} \binom{n}{k} W_{ij_1}^k W_{ij_2}^{n-k} \right) v_{j_1} v_{j_2}$$

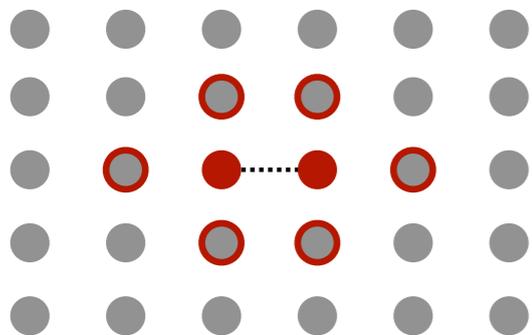
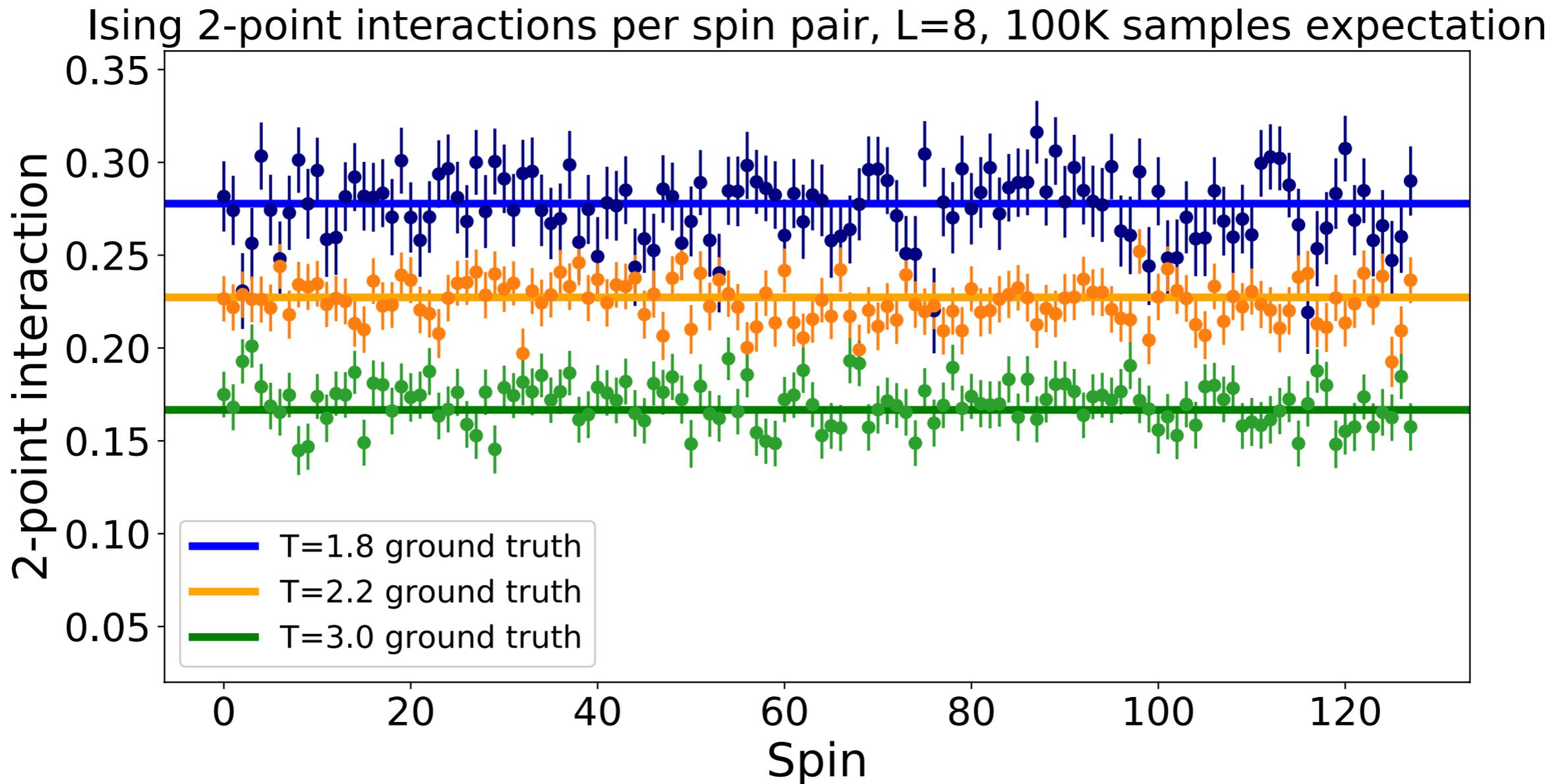
e.g. 3-point interaction:

...

$$\frac{1}{6} \sum_i \ln \frac{(1 + e^{c_i + W_{ij_1} + W_{ij_2} + W_{ij_3}})(1 + e^{c_i + W_{ij_1}})(1 + e^{c_i + W_{ij_2}})(1 + e^{c_i + W_{ij_3}})}{(1 + e^{c_i + W_{ij_1} + W_{ij_2}})(1 + e^{c_i + W_{ij_1} + W_{ij_3}})(1 + e^{c_i + W_{ij_2} + W_{ij_3}})(1 + e^{c_i})}$$

**Closed form expression!**

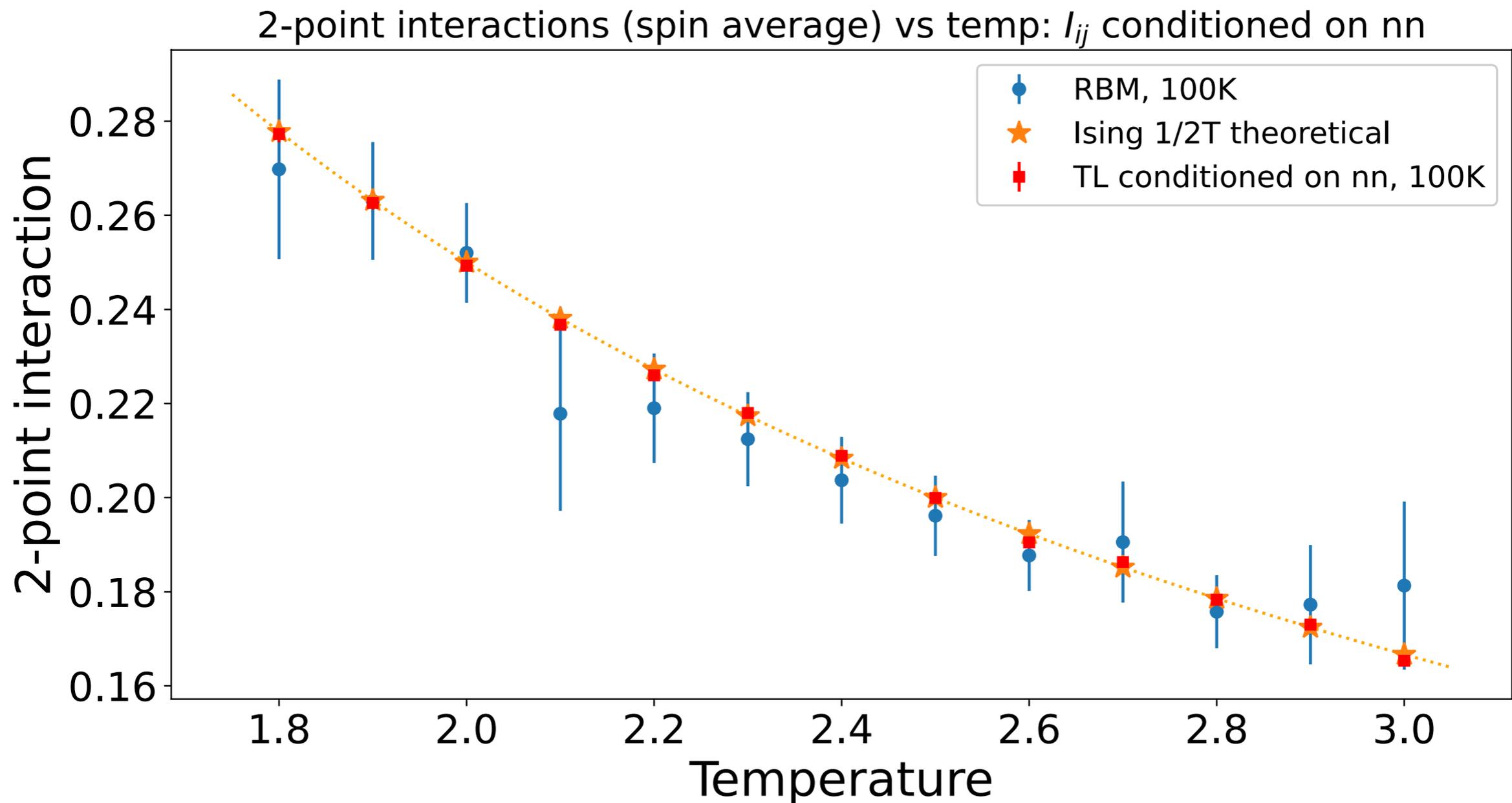
# Back to Ising ...



$$H(s) = - \sum_{\langle i,j \rangle} J_{ij} s_i s_j$$

# Model-independent estimation results

Conditioning on parent spins to isolate pairs from the rest of the system (Markovian). Run time: Few seconds per temperature.



100K samples