# A solvable model of Generative Diffusion and the memorization phenomenon

Leonardo Bandera

# Agenda

- Diffusion Sampling and problem statement

- Understanding training dynamics: the original setting

- First line of research: towards a more realistic architecture

- Second line of research: structured data

# Generative Diffusion

The Forward Process

$$\boldsymbol{x}_0 \longrightarrow \boldsymbol{x}_1 \longrightarrow \cdots \longrightarrow \boldsymbol{x}_T$$



Original Data

Complete Noise

$$\boldsymbol{x}_0 \longleftarrow \boldsymbol{x}_1 \longleftarrow \cdots \longleftarrow \boldsymbol{x}_T$$
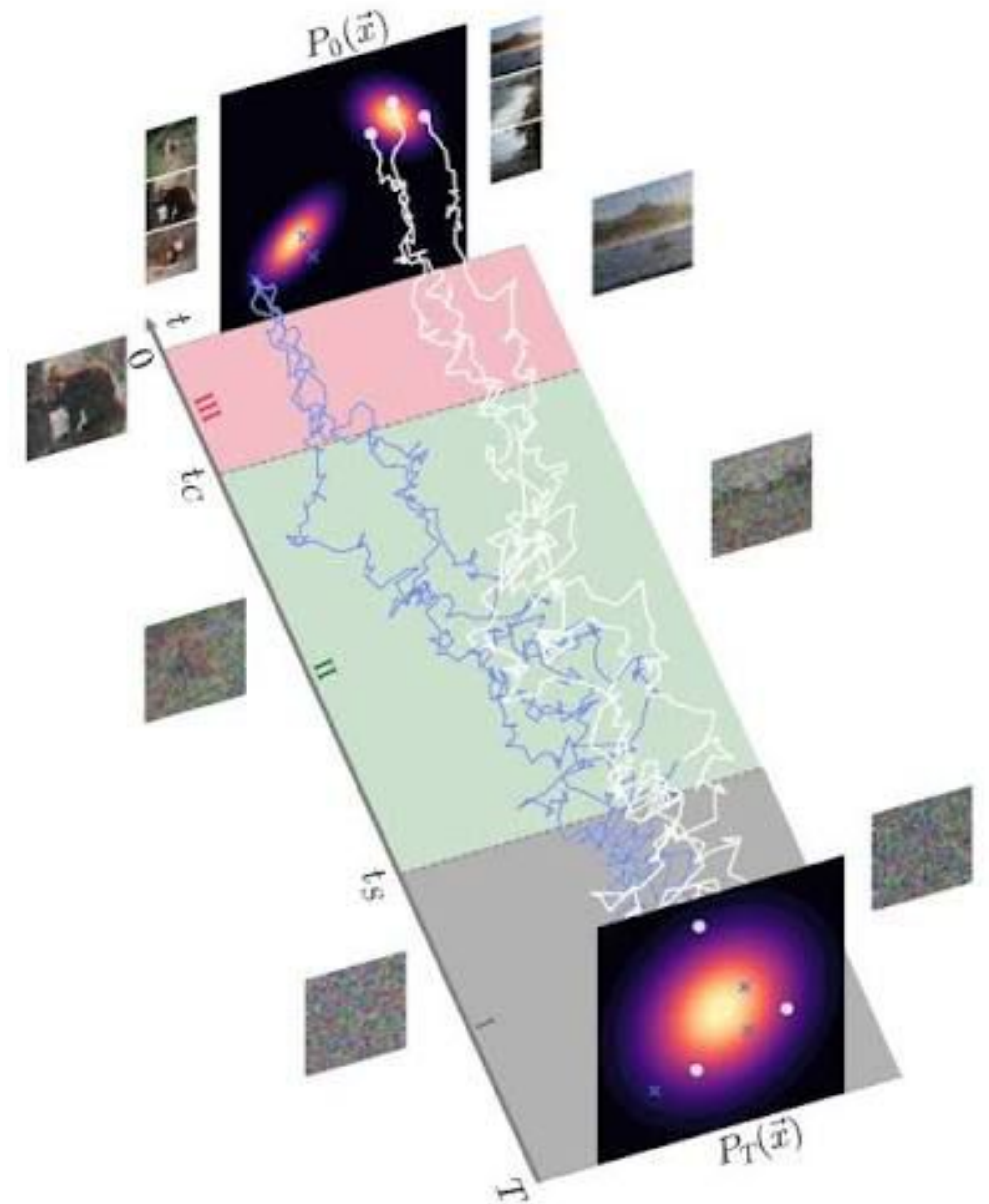
The Generative Backward Process

# Score-based Diffusion

- Forward process: $\dfrac{d\boldsymbol{x}}{dt} = -\boldsymbol{x}(t) + \boldsymbol{\eta}(t)$

- Reversed process: $\dfrac{d\boldsymbol{x}}{d\tau} = \boldsymbol{x}(\tau) + 2\nabla\log\widetilde{P}_\tau(\boldsymbol{x}) + \tilde{\boldsymbol{\eta}}(\tau)$

- The score $S(\boldsymbol{x}, t) = \nabla\log P_t(\boldsymbol{x})$ is unknown… we learn it!

$$\mathscr{L}_\lambda(\boldsymbol{\theta}) = \frac{1}{n}\sum_{\mu=1}^{n}\mathbb{E}_{t\sim Q(t)}\mathbb{E}_{\boldsymbol{\xi}\sim\mathscr{N}(0,\mathbb{I}_d)}\left\| \hat{S}^{\boldsymbol{\theta}}(\boldsymbol{x}_t^\mu(\boldsymbol{\xi}), t) + \frac{\boldsymbol{\xi}}{\sqrt{\Delta_t}} \right\|^2$$

# Problem: sampling with the empirical loss minimizer leads to memorization!



Giulio Biroli, Tony Bonnaire, Valentin de Bortoli, and Marc Mézard. Dynamical regimes of diffusion models. Nature Communications, 15(1):9957, nov 2024
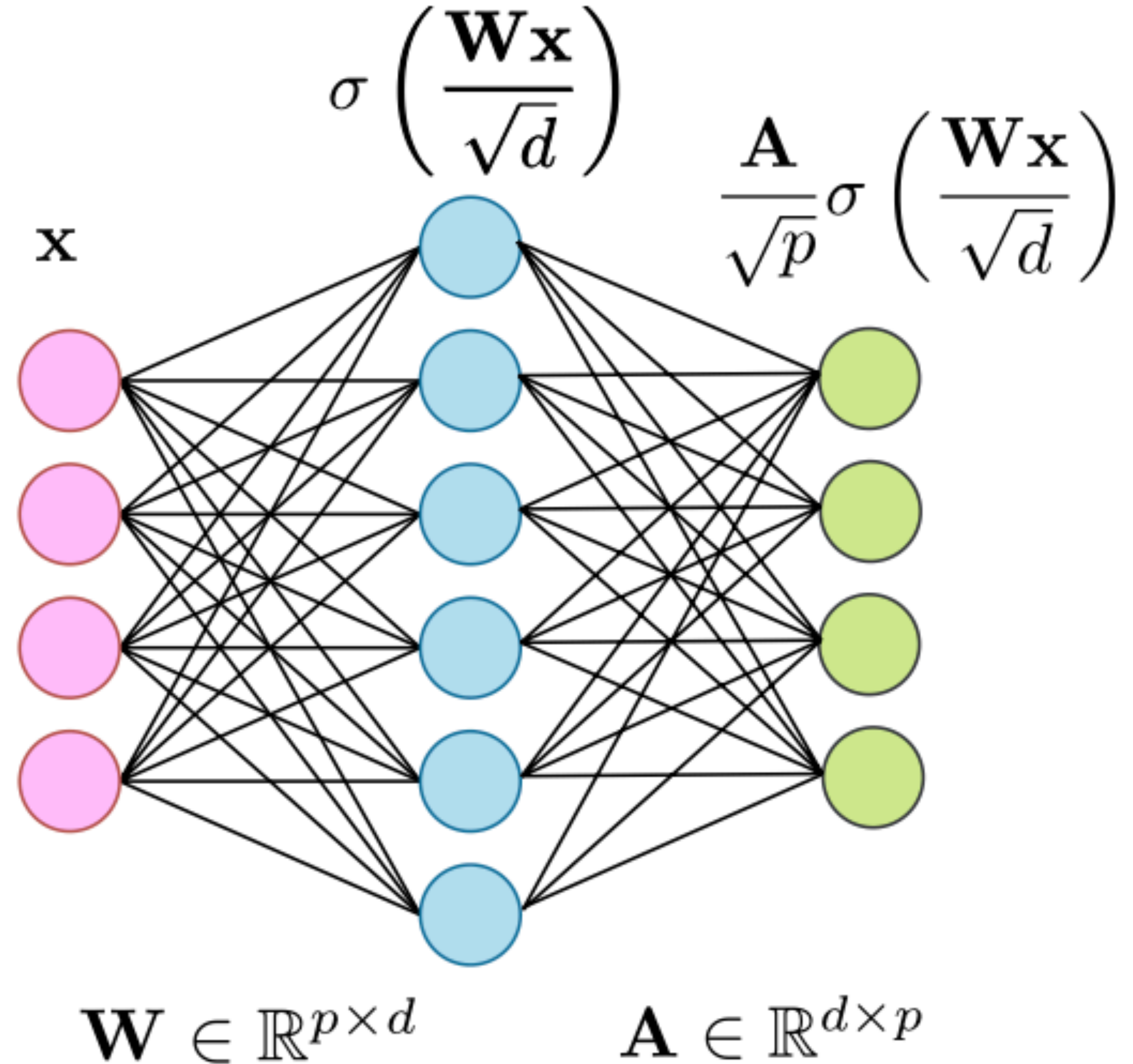
# Why real diffusion systems avoid this degeneracy?

- Architectural constraints

- **Training dynamics**

# The original setting

Gaussian data

One RFNN per time

Tony Bonnaire, Raphaël Urfin, Giulio Biroli, and Marc Mézard. Why diffusion models don't memorize: The role of implicit dynamical regularization in training. arXivpreprint arXiv:2505.17638, 2025

$$\sigma \left( \frac{\mathbf{Wx}}{\sqrt{d}} \right)$$

$$\frac{\mathbf{A}}{\sqrt{p}} \sigma \left( \frac{\mathbf{Wx}}{\sqrt{d}} \right)$$

$\mathbf{x}$

$$\mathbf{W} \in \mathbb{R}^{p \times d} \qquad \mathbf{A} \in \mathbb{R}^{d \times p}$$

Gradient Flow

$$\dot{A}(\tau) = -2\Delta_t \frac{d}{p} A U - \frac{2d\sqrt{\Delta_t}}{\sqrt{p}} V^\top$$

where

$$U = \frac{1}{n} \sum_{\nu=1}^{n} \mathbb{E}_{\boldsymbol{\xi}}[\sigma(\frac{Wx_t^\nu(\boldsymbol{\xi})}{\sqrt{d}})\sigma(\frac{Wx_t^\nu(\boldsymbol{\xi})}{\sqrt{d}})^\top]$$

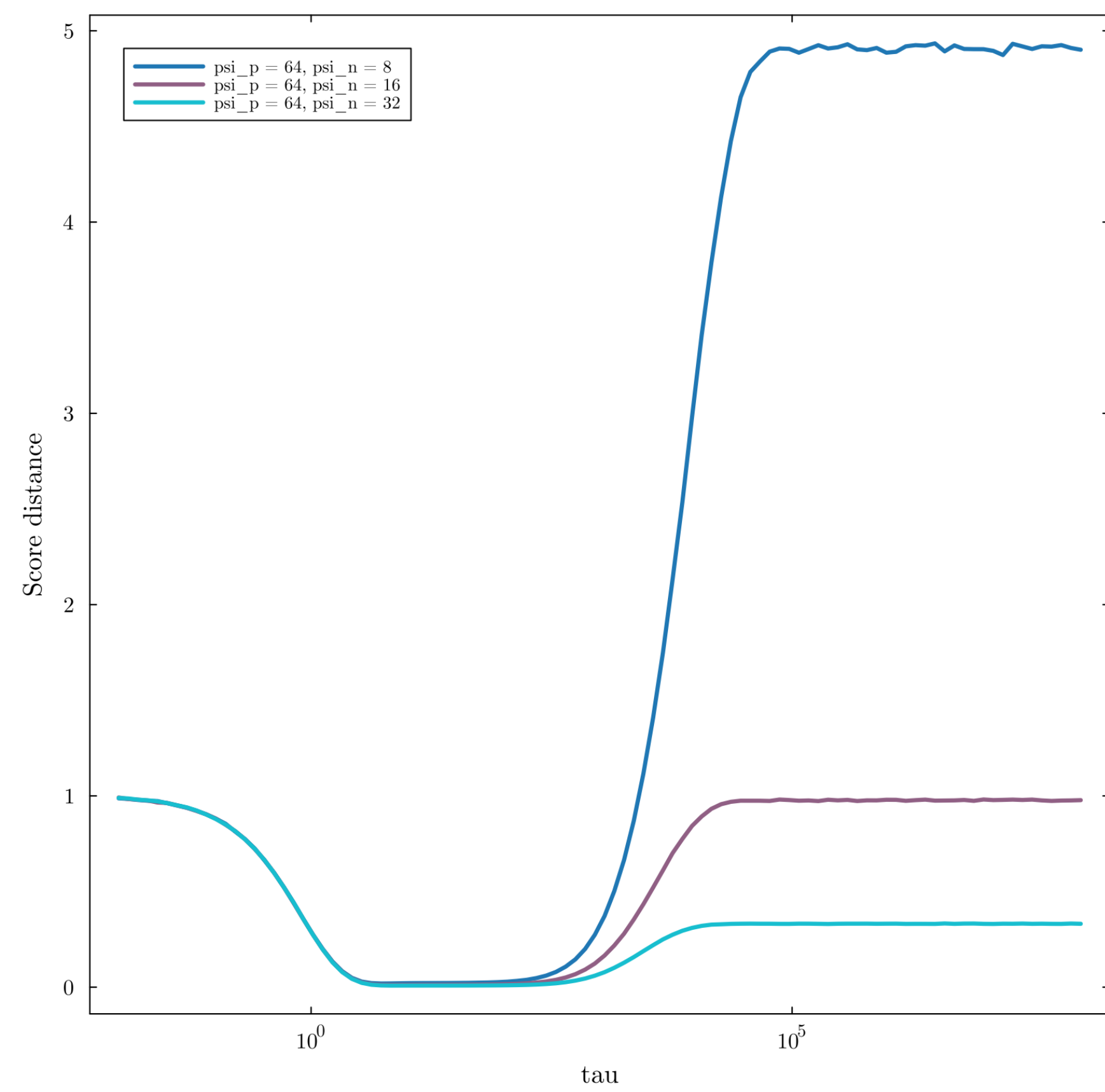The timescales of the training dynamics are given by the inverse eigenvalues of the matrix $\Delta_t \frac{d}{p} U$ !
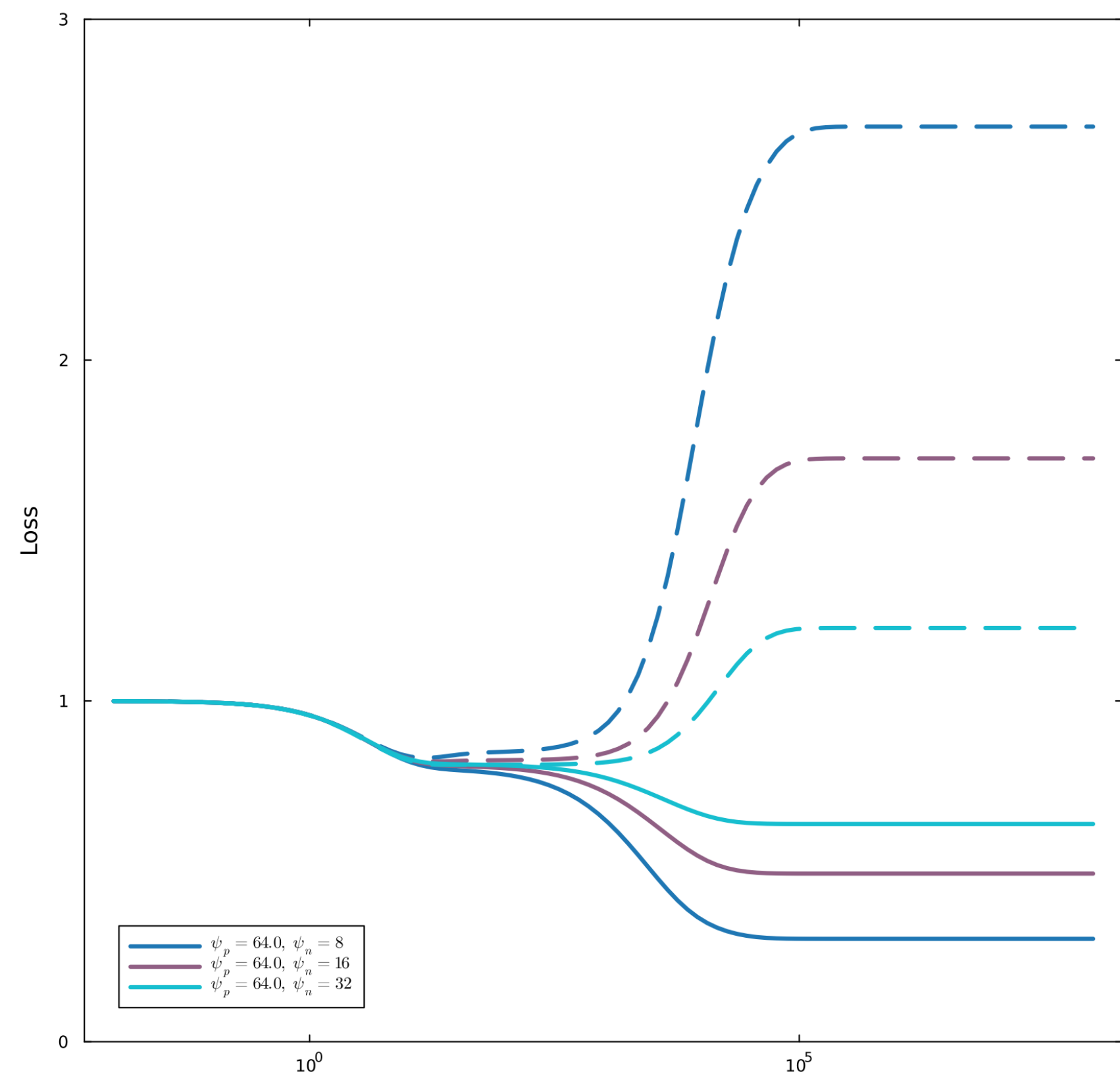
# $U$ spectrum

- Stieltjes transform

- GEP

- Replica

*Two bulks!*

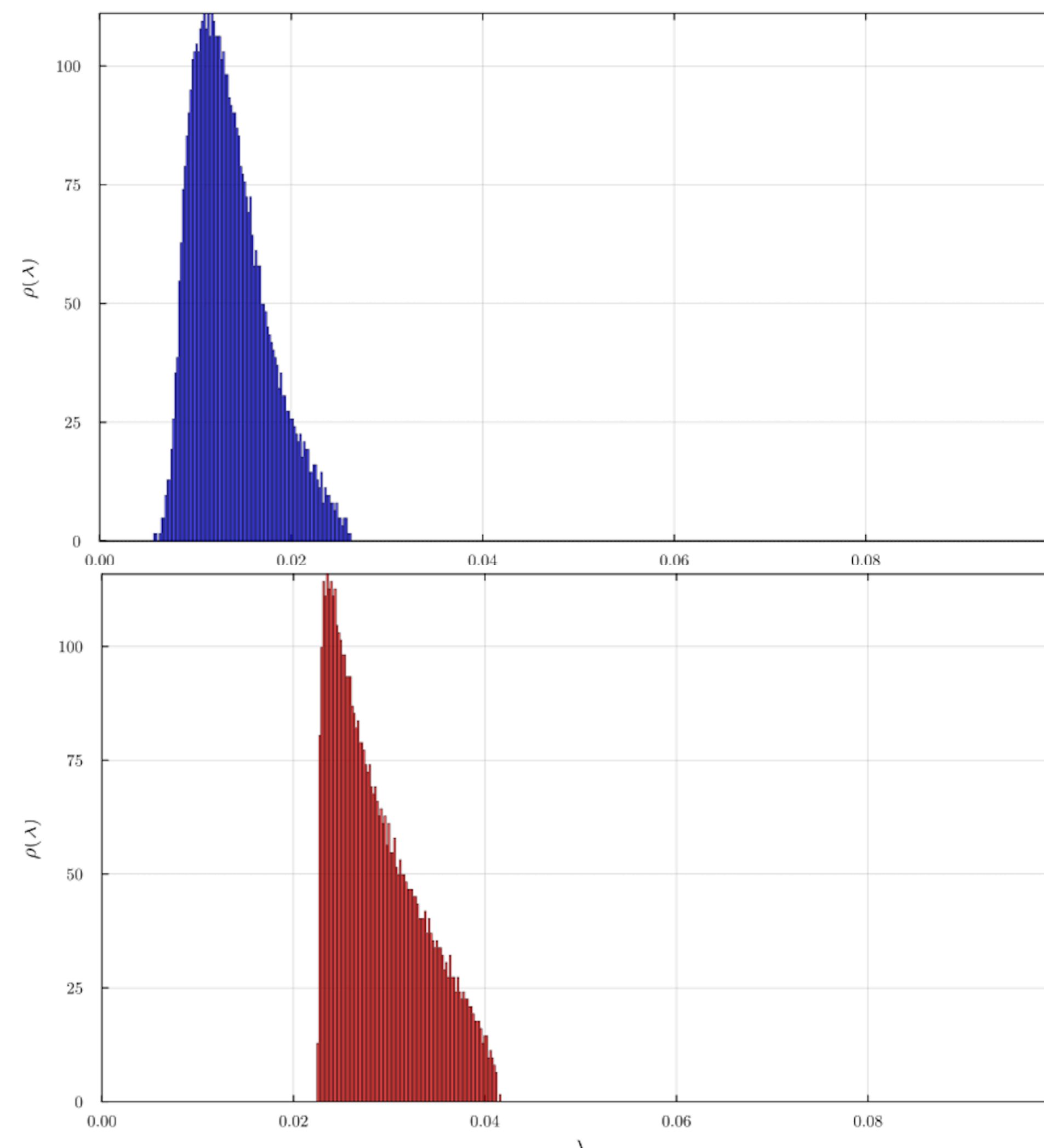What's happening at the metrics at the timescales associated to the bulks?

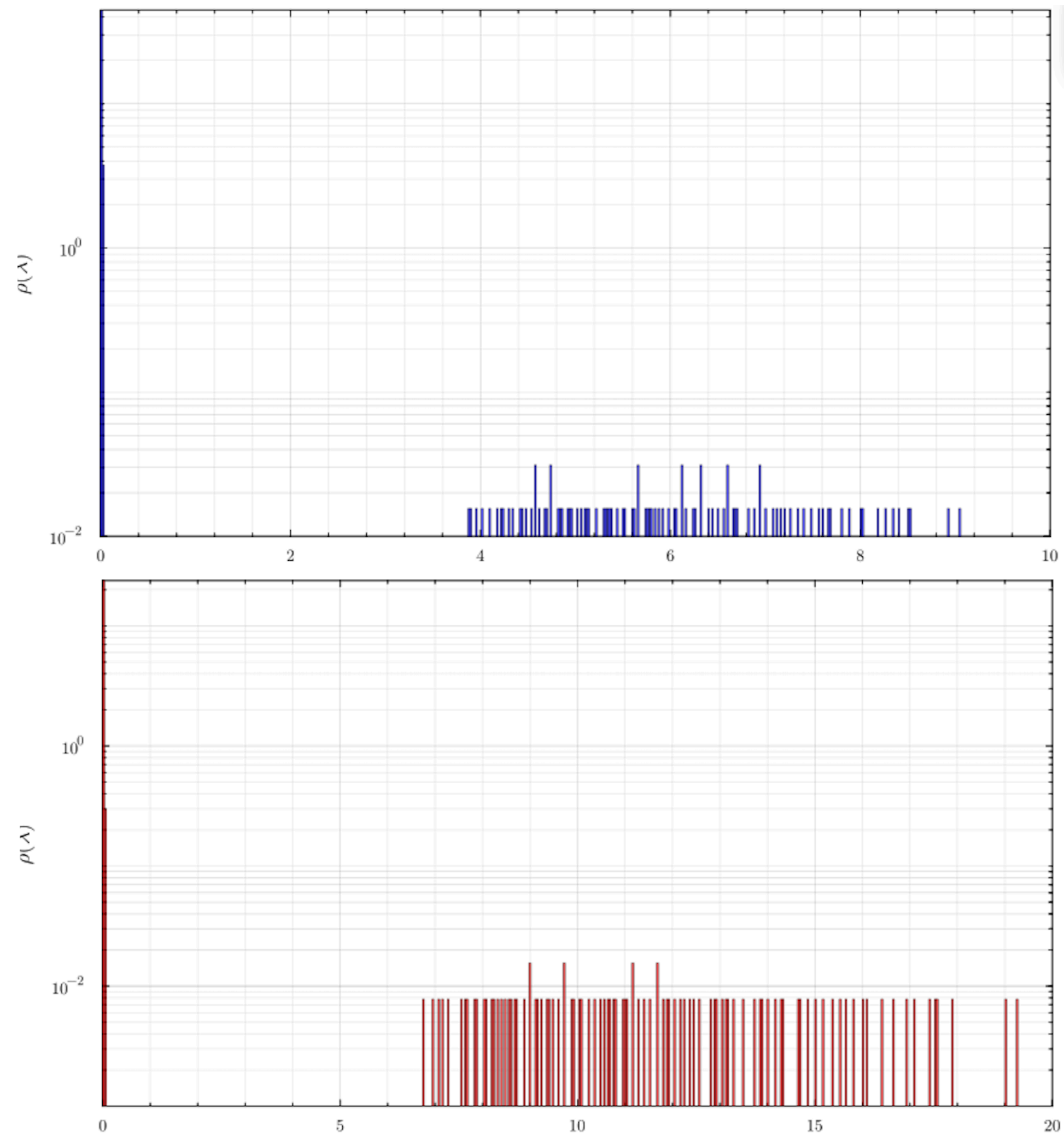# Time Integrated RFNN

$$S_A(x, t) = \alpha_t \frac{A}{\sqrt{p}} \sigma\left( \frac{Wx}{\sqrt{d}} + tb \right) + \beta_t x$$

$$\mathcal{L}(A, \{x_i^\nu\}) = \frac{1}{n} \sum_{\nu=1}^{n} \frac{1}{d} \mathbb{E}_{\xi, t} \parallel \sqrt{\Delta_t} S_A(x_t^\nu(\xi), t) + \xi \parallel^2$$

- Same phenomenology

- Time integration brakes GEP: we cannot compute $U$ spectrum!

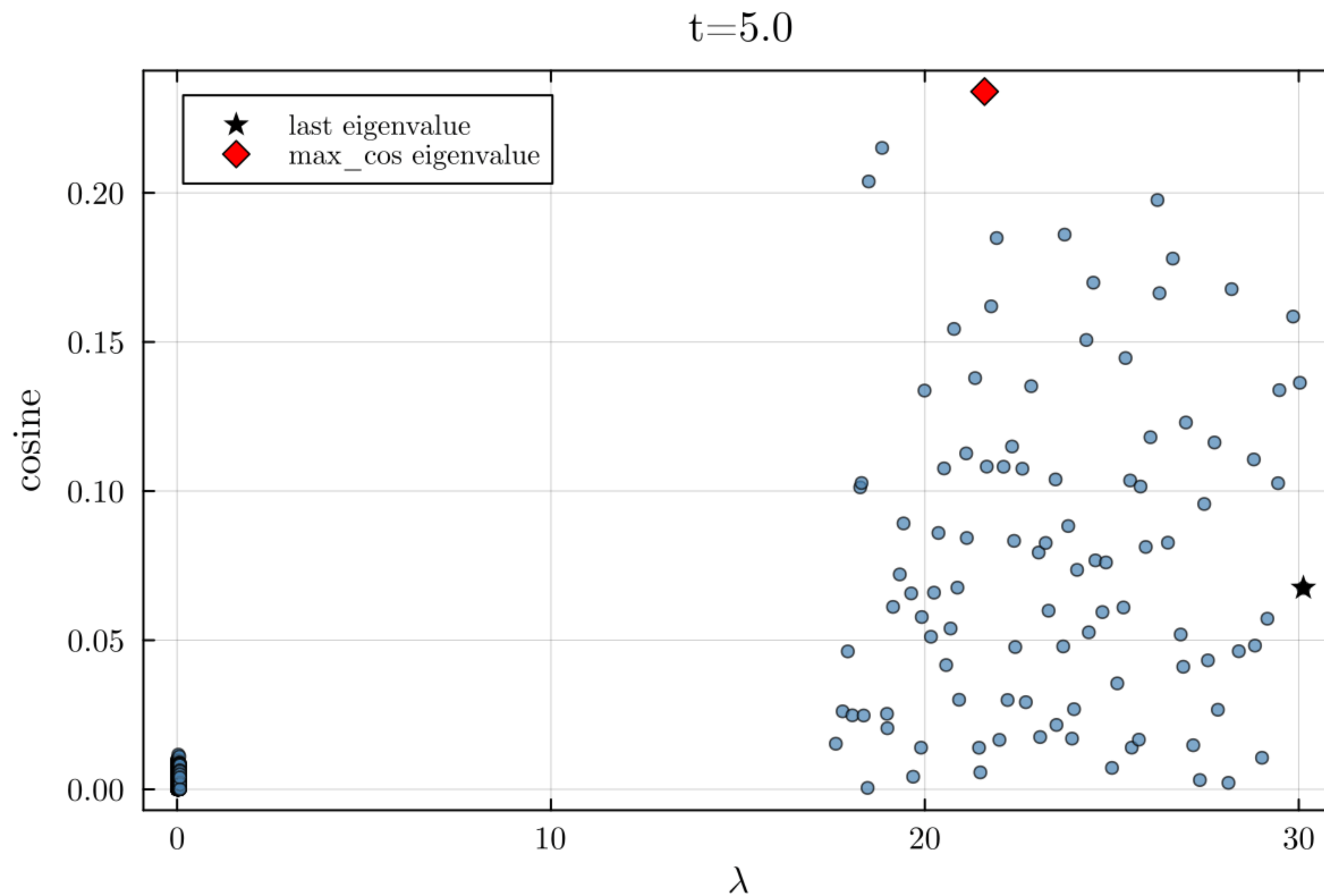# MoG Data

$$P_0 = \frac{1}{2}\mathcal{N}(\boldsymbol{m}, \sigma_{\boldsymbol{x}}^2 \mathbb{I}_d) + \frac{1}{2}\mathcal{N}(-\boldsymbol{m}, \sigma_{\boldsymbol{x}}^2 \mathbb{I}_d)$$

or equivalently

$$\boldsymbol{x}^{\mu} = c^{\mu}\boldsymbol{m} + \sigma_x \boldsymbol{z}^{\mu}$$

# Phenomenology:
# $U$ spectrum BBP transition

# GEP holds!

- $U_{GEP}$ for MoG data is $U_{GEP}$ for gaussian data plus rank-1 terms: the bulks are the same!

- We must focus on the outlier…

We need to compute this:

$$\overline{G}_{ij} = \lim_{n \to 0} \mathbb{E}_x \int \prod_{a=1}^{n} d\psi^a \exp\left[-\frac{1}{2} \sum_a (\psi^a)^T (z\mathbb{I}_N - \boldsymbol{U}) \psi^a\right] \psi_i^1 \psi_j^1$$
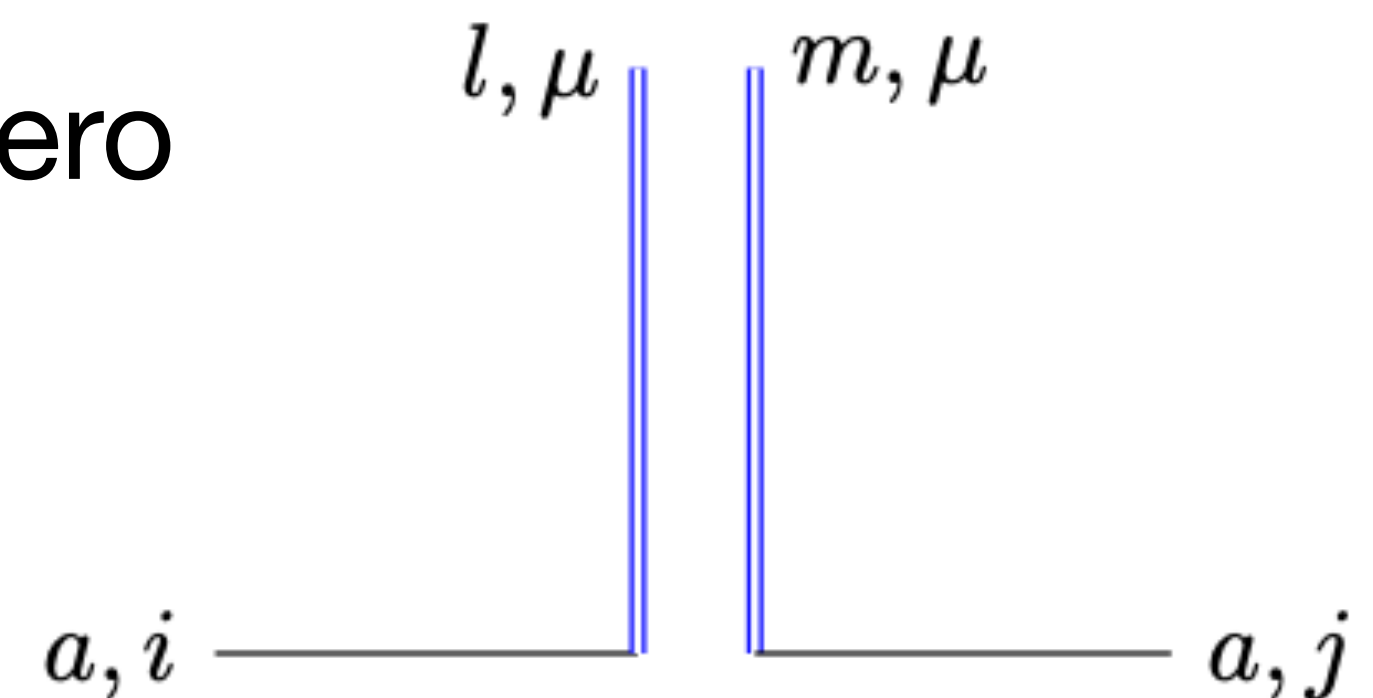
# It's a matter of gaussian averages!

- Wick theorem to reduce averages over many fields to averages of pairs

- Feynman diagrams to keep track of all possible pairings and interactions between different gaussian fields
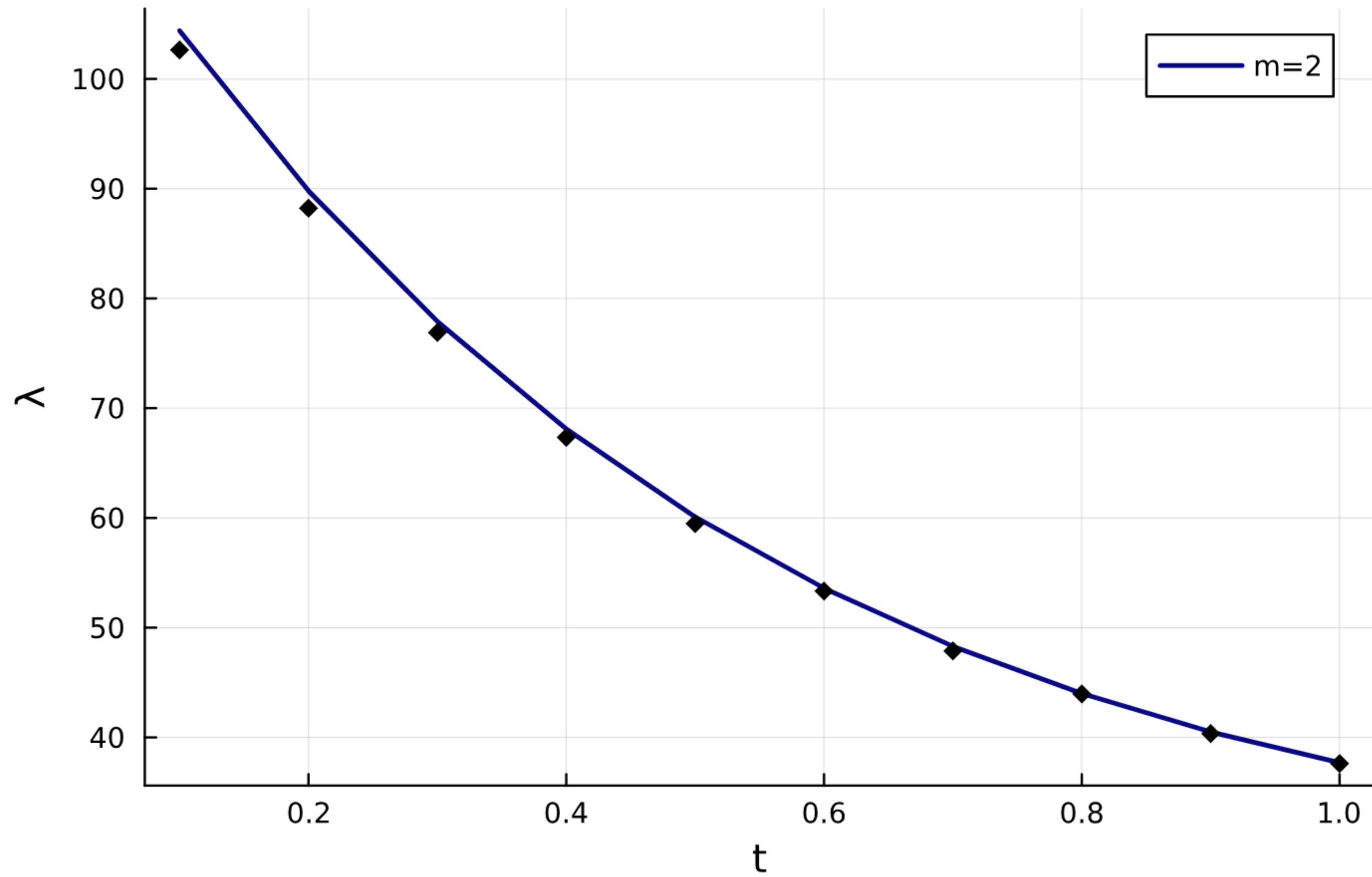
$$\Sigma_{ij}^b = \sum_{\mu l} \; i \overset{l}{\underset{\mu}{\frown}} j \; + \; \sum_{\mu,l,m,n,k} i - \mu \overset{l}{\underset{\substack{k \; n \\ \bar{G}_{kn}^b}}{\overset{m}{\bigcirc}}} \mu - j \; + \; \sum_{\substack{\mu,l,m,n,k \\ p,q,r}} i - \mu \overset{l}{\underset{\substack{k \; n \\ \bar{G}_{kn}^b}}{\overset{m}{\bigcirc}}} \mu \; \mu \overset{q}{\underset{\substack{r \; p \\ \bar{G}_{pr}^b}}{\bigcirc}} \mu - j \; + \cdots$$

We get to a self-consistent equation for $\overline{G}_{ij}$ :

- solving it gives the bulks

- the condition for the outlier is setting a denominator to zero

Solutions for the outlier equation

# *This is interesting at two levels:*

- *Diffusion dynamics*: what's happening at trajectories at transition time?

- *Training dynamics*: what's happening at the training metrics at the timescale associated to the outlier?
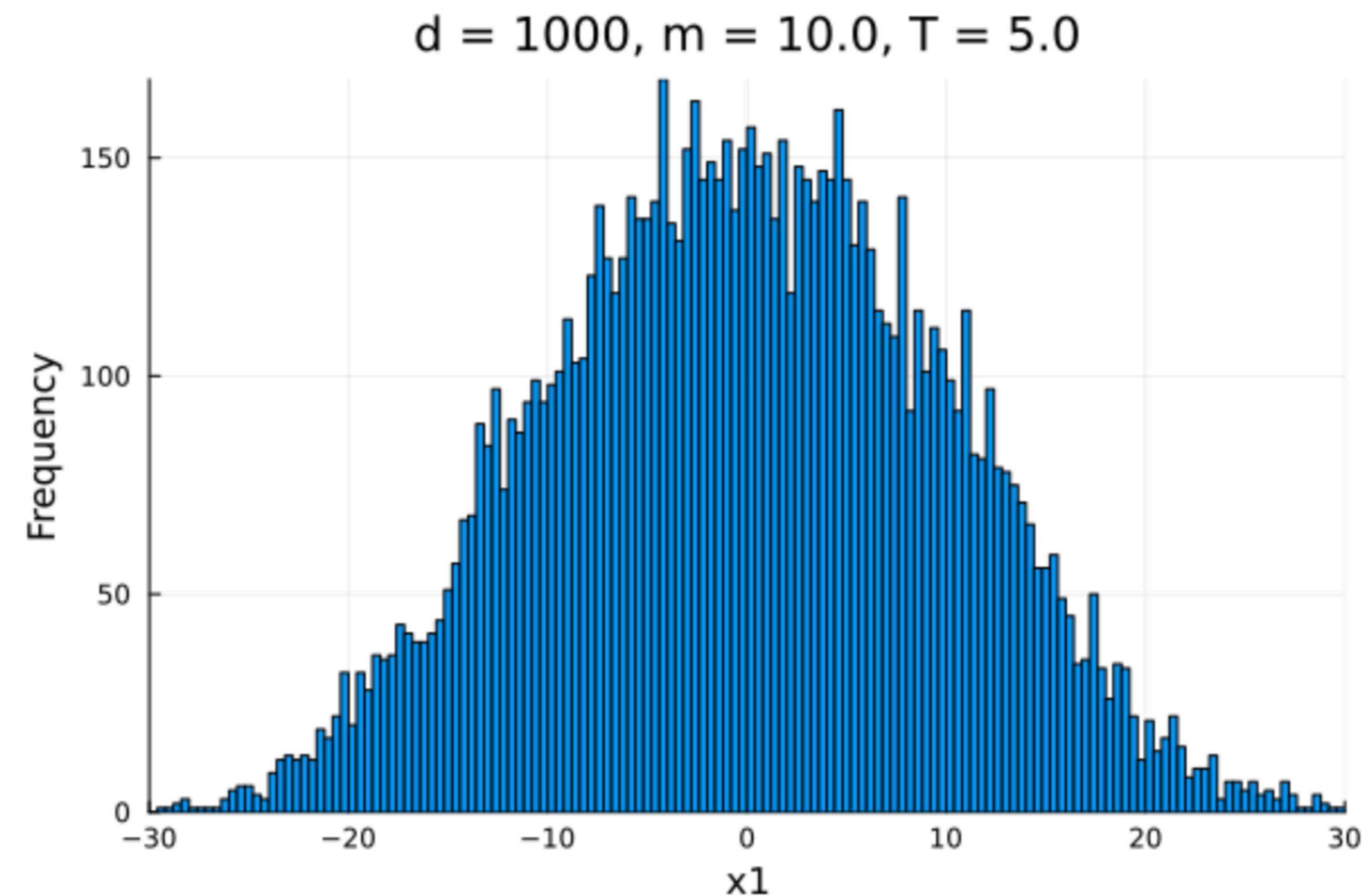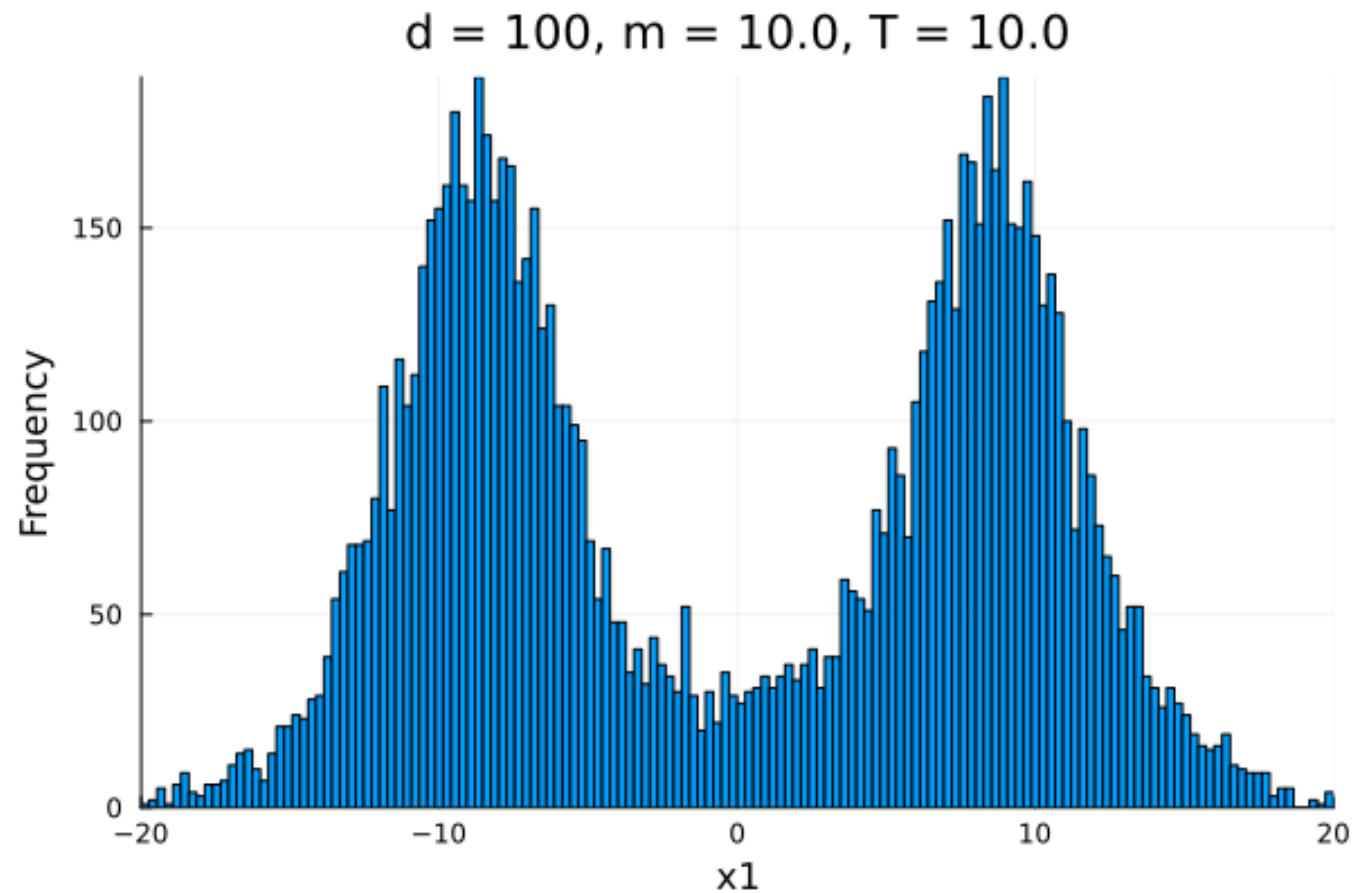
# But…

## … a legit doubt:

# Can RFNN learn a mixture?

# It depends on the Gaussians separation!

If $\| m \| \sim \mathcal{O}(1)$ the random features score is approximately linear, so the backward process can at most produce a gaussian distribution!



d = 1000, m = 10.0, T = 5.0

If $\|\boldsymbol{m}\| \sim \mathcal{O}(\sqrt{d})$ the random features score is non linear, and we can sample effectively from a mixture.



d = 100, m = 10.0, T = 10.0

# Conclusions

The behavior of the dominant eigenvalue reveals the point (in terms of $t$ and $\tau$) at which the diffusion process becomes sensitive to the informative direction.

- **Non-separated regime**: the process cannot recover the bimodal structure and diffusion collapses to a single component whose variance aligns with the signal direction

- **Well-separated regime**: the process recovers the multimodal structure and separates the modes