**Title:** Architectural bias in a transport-based generative model : an asymptotic perspective

**Speakers:** Hugo Cui

**Collection/Series:** Theory + AI Workshop: Theoretical Physics for AI

**Date:** April 10, 2025 - 9:45 AM

**URL:** https://pirsa.org/25040092

**Abstract:**

We consider the problem of learning a generative model parametrized by a two-layer auto-encoder, and trained with online stochastic gradient descent, to sample from a high-dimensional data distribution with an underlying low-dimensional structure. We provide a tight asymptotic characterization of low-dimensional projections of the resulting generated density, and evidence how mode(l) collapse can arise. On the other hand, we discuss how in a case where the architectural bias is suited to the target density, these simple models can efficiently learn to sample from a binary Gaussian mixture target distribution.

# Architectural bias in generative models
## —an asymptotic viewpoint
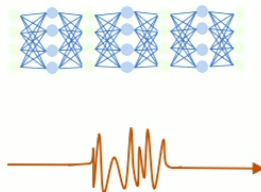
*Hugo Cui*

PI Theory + AI workshop

Based on

**HC**, Pehlevan, Lu, *Precise asymptotics of learning diffusion models: theory & insights,* **ArXiv 2025**
**HC**, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity,* **ICLR 2024**

1

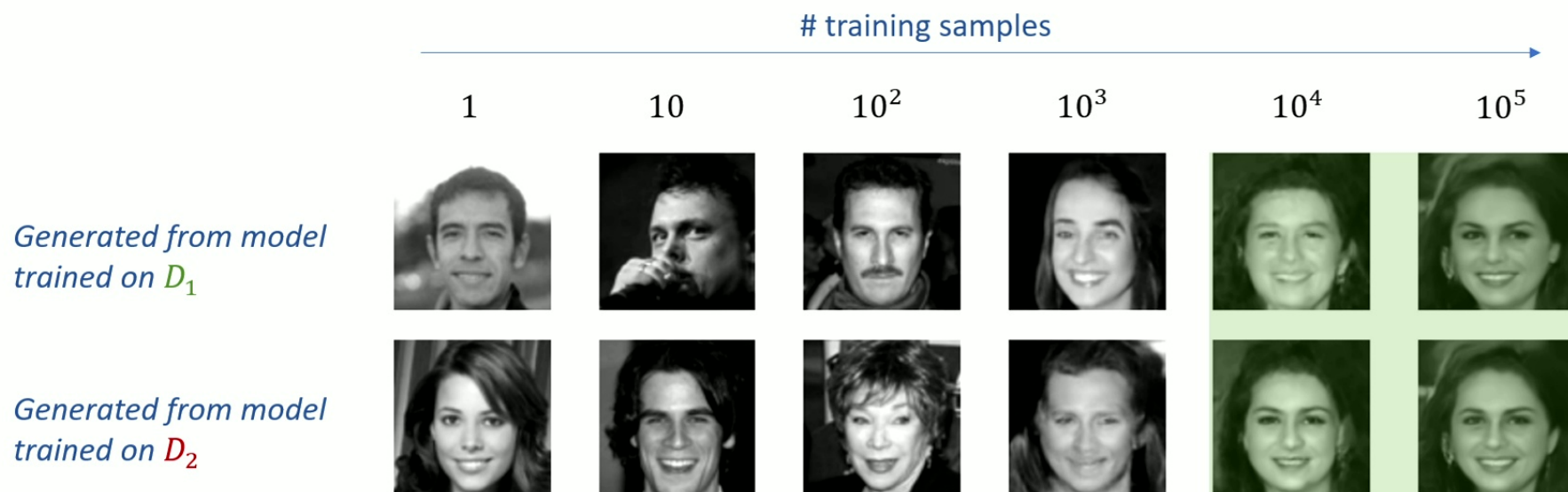Training set $\sim \rho$          new generated samples

***Transport-based generative models*** learn to sample (generate) complex distributions in high-dimensions from moderate training sets.

→Existence of strong *inductive biases* in the architecture.

Ho, Jain, Abbeel, *Denoising Diffusion Probabilistic Models*, NeurIPS 2020
Sohl-Dickstein et al., *unsupervised learning using nonequi-librium thermodynamics,* ICML 2015
Song and Ermon, *Generative modeling by estimating gradients of the data distribution*. NeurIPS 2019

2

# training samples

$$1 \qquad 10 \qquad 10^2 \qquad 10^3 \qquad 10^4 \qquad 10^5$$

*Generated from model trained on $D_1$*



*Generated from model trained on $D_2$*

Kadkhodaie et al., *Generalization in diffusion models arises from geometry-adaptive harmonic representation*, ICLR 2024

Two models trained on disjoint training sets $D_1$ and $D_2$ generate the **same image** from a given prompt when trained with sufficiently many samples ($\sim 2 - 16$x dimension).

3

How is the distribution of generated samples shaped by the network architecture?

→Try to understand in simple models.

4

## Analysis of the transport only

Chen et al,. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. arXiv:2209.11215, 2022.
Biroli et al, Dynamical regimes of diffusion models. Nature Communications, 15(1):9957, 2024
...

5

## Analysis of the transport only

Chen et al,. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. arXiv:2209.11215, 2022.
Biroli et al, Dynamical regimes of diffusion models. Nature Communications, 15(1):9957, 2024
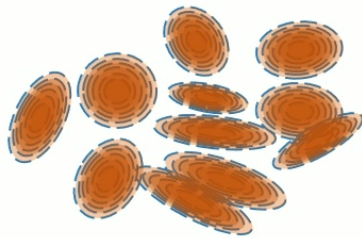...

## Sample bounds when density can be perfectly learnt by the model class with enough samples:

Boffi et al., Shallow diffusion networks provably learn hidden low-dimensional structure., arXiv:2410.11275, 2024.
Chen et al., *Score approx-imation, estimation and distribution recovery of diffusion models on low-dimensional data*, ICML 2023
Oko, Akiyama and Suzuki, *Diffusion models are minimax optimal distribution estimators,* ICML 2023

6

Target density $\rho$

---- Generated density $\hat{\rho}$

## Analysis of the transport only

Chen et al,. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. arXiv:2209.11215, 2022.
Biroli et al, Dynamical regimes of diffusion models. Nature Communications, 15(1):9957, 2024
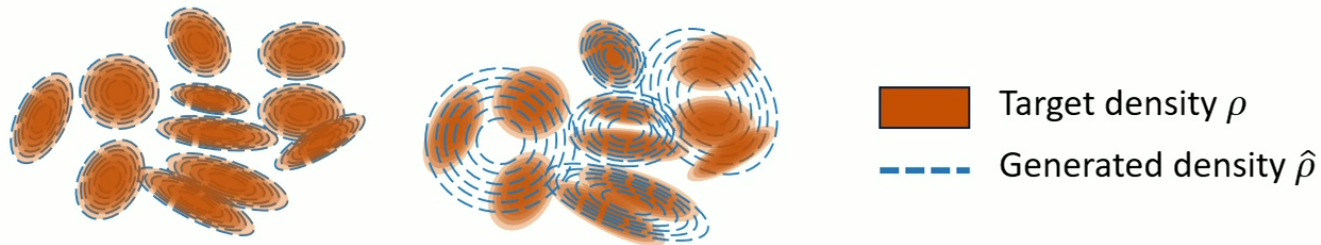…

## Sample bounds when density can be perfectly learnt by the model class with enough samples:

Boffi et al., Shallow diffusion networks provably learn hidden low-dimensional structure., arXiv:2410.11275, 2024.
Chen et al., *Score approx-imation, estimation and distribution recovery of diffusion models on low-dimensional data*, ICML 2023
Oko, Akiyama and Suzuki, *Diffusion models are minimax optimal distribution estimators,* ICML 2023

7

Target density $\rho$

Generated density $\hat{\rho}$

<u>Analysis of the transport only</u>

Chen et al,. Sampling is as easy as learning the score: theory for diffusion models with minimal data assumptions. arXiv:2209.11215, 2022.
Biroli et al, Dynamical regimes of diffusion models. Nature Communications, 15(1):9957, 2024
…

<u>Sample bounds when density can be perfectly learnt by the model class with enough samples:</u>

Boffi et al., Shallow diffusion networks provably learn hidden low-dimensional structure., arXiv:2410.11275, 2024.
Chen et al., *Score approx-imation, estimation and distribution recovery of diffusion models on low-dimensional data*, ICML 2023
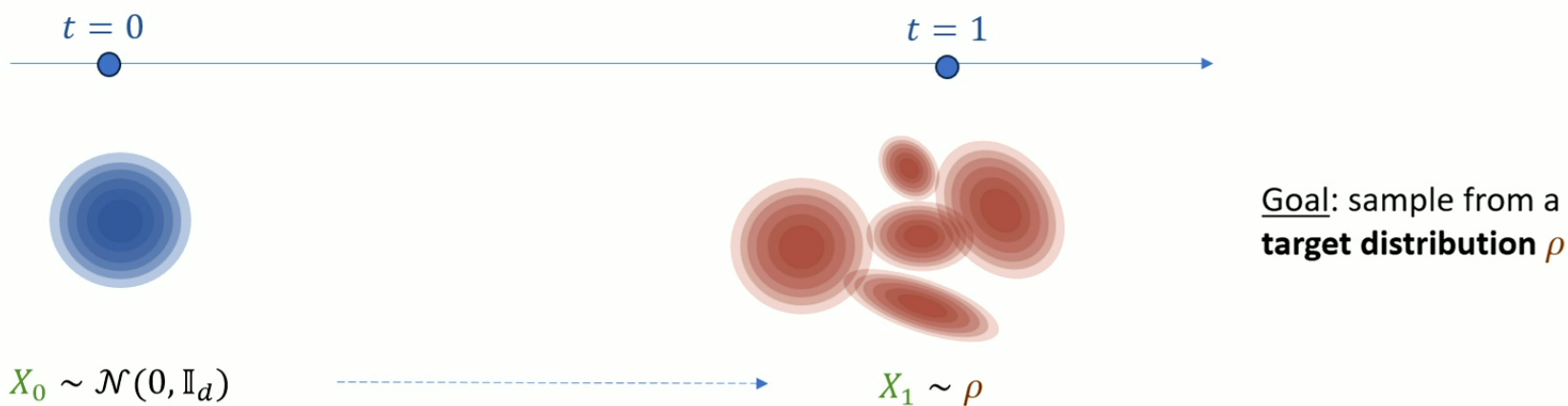Oko, Akiyama and Suzuki, *Diffusion models are minimax optimal distribution estimators,* ICML 2023

→*To complement these results* : a **tight** characterization of the generated density in the case **where architecture and target distribution are not perfectly matched**.

8

1. Generated density for an auto-encoder parametrized model

2. Failure modes : mode(l) collapse

3. Aligned case: binary isotropic Gaussian mixture distribution.

**HC**, Pehlevan, Lu, *Precise asymptotics of learning diffusion models: theory & insights,* **ArXiv 2025**
**HC**, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity,* **ICLR 2024**

9

$t = 0$

$t = 1$

Goal: sample from a
**target distribution** $\rho$

$X_0 \sim \mathcal{N}(0, \mathbb{I}_d)$

$X_1 \sim \rho$

Albergo, Boffi, and Vanden-Eijnden, *Stochastic interpolants: A unifying framework for flows and diffusions*. arXiv:2303.08797, 2023.
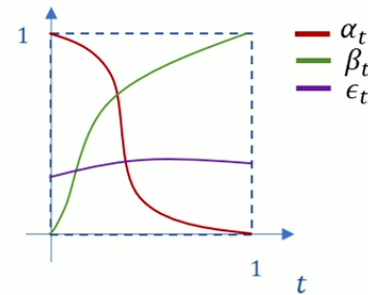
10

$t = 0$

$X_0 \sim \mathcal{N}(0, \mathbb{I}_d)$

$t = 1$

$X_1 \sim \rho$

The sampling can be done by transporting $X_0$ through the SDE for $t \in [0,1]$

$$\frac{d}{dt}X_t = \left(\dot{\beta}_t - \frac{\dot{\alpha}_t}{\alpha_t}\beta_t + \epsilon_t\frac{\beta_t}{\alpha_t^2}\right)f(t, X_t) + \left(\frac{\dot{\alpha}_t}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2}\right)X_t + \sqrt{2\epsilon_t}dW_t$$

*For any choice of interpolation schedules st:*

$\alpha, \beta \in \mathcal{C}^2([0,1])$
$\alpha_0 = \beta_1 = 1, \alpha_1 = \beta_0 = 0$
$\epsilon_t \geq 0$



$\alpha_t$
$\beta_t$
$\epsilon_t$

1

1
$t$

Albergo, Boffi, and Vanden-Eijnden, *Stochastic interpolants: A unifying framework for flows and diffusions.* arXiv:2303.08797, 2023.

11

$t = 0$                                                    $t = 1$

$X_0 \sim \mathcal{N}(0, \mathbb{I}_d)$                              $X_1 \sim \rho$

The sampling can be done by transporting $X_0$ through the SDE for $t \in [0,1]$

$$\frac{d}{dt} X_t = \left( \dot{\beta}_t - \frac{\dot{\alpha}_t}{\alpha_t} \beta_t + \epsilon_t \frac{\beta_t}{\alpha_t^2} \right) \boldsymbol{f}(t, X_t) + \left( \frac{\dot{\alpha}_t}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2} \right) X_t + \sqrt{2\epsilon_t} dW_t$$

*Denoising function* is the minimizer of a **denoising objective**

$$\boldsymbol{f} = \min_h \int_0^1 \mathbb{E}_{x_1 \sim \rho, x_0 \sim \mathcal{N}(0, \mathbb{I}_d)} \| h(t, \alpha_t x_0 + \beta_t x_1) - x_1 \|^2 dt$$

Learnable from data
Empirical average
Network param.

Albergo, Boffi, and Vanden-Eijnden, *Stochastic interpolants: A unifying framework for flows and diffusions.* arXiv:2303.08797, 2023.

12

$$\forall x \in \mathbb{R}^d, \quad f_{b,w}(x) = bx + \frac{w}{\sqrt{d}}\sigma\left(\frac{w^\top x}{\sqrt{d}}\right)$$
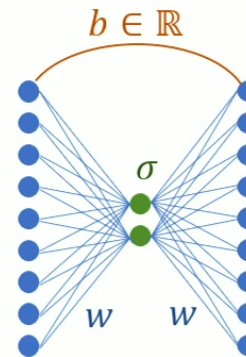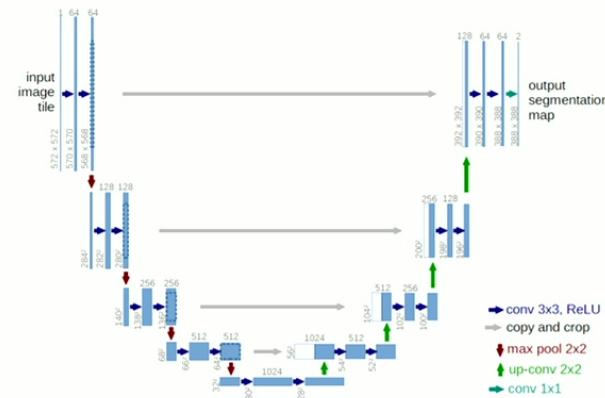
Trainable skip connection $b \in \mathbb{R}$

Weight matrix $w \in \mathbb{R}^{d \times r}$



Vincent et al., *Stacked denoising AEs:Learning useful representations in a deep net. with a local denoising criterion,* JMLR 2010     13

$$\forall x \in \mathbb{R}^d, \quad \boxed{f_{b,w}(x) = bx + \frac{w}{\sqrt{d}} \sigma\left(\frac{w^\top x}{\sqrt{d}}\right)}$$

Trainable skip connection $b \in \mathbb{R}$

Weight matrix $w \in \mathbb{R}^{d \times r}$

<u>Remark</u>: U-Nets are used in practice.

- skip connections
- bottlenecks
- convolutional layers



Ronneberger, Fischer, and Brox *U-net: Convolutional networks for biomedical image segmentation.* MICCAI 2015
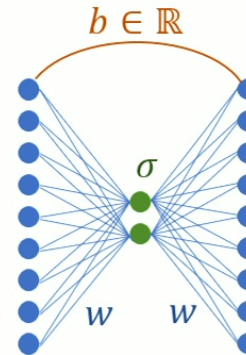
Vincent et al., *Stacked denoising AEs:Learning useful representations in a deep net. with a local denoising criterion,* JMLR 2010          14

$$\forall x \in \mathbb{R}^d, \quad f_{b,w}(x) = bx + \frac{w}{\sqrt{d}} \sigma\left(\frac{w^\top x}{\sqrt{d}}\right)$$



Trainable skip connection $b \in \mathbb{R}$

Weight matrix $w \in \mathbb{R}^{d \times r}$

Given a training set of $n$ i.i.d samples $\left\{x_1^\mu \sim \rho, x_0^\mu \sim \mathcal{N}(0, \mathbb{I}_d)\right\}_{\mu=1}^n$ one can train the network $f_{b,w}(x)$ with <u>online SGD</u>

$$b_{\mu+1} - b_\mu = -\frac{\eta}{d^2}\left(\partial_b \mathbb{E}_t \left\| x_1^\mu - f_{b_\mu, w_\mu}(\alpha_t x_0^\mu + \beta_t x_1^\mu)\right\|^2\right)$$

$$w_{\mu+1} - w_\mu = -\eta\left(\nabla_w \mathbb{E}_t \left\| x_1^\mu - f_{b_\mu, w_\mu}(\alpha_t x_0^\mu + \beta_t x_1^\mu)\right\|^2 + \lambda/d w_\mu\right)$$

Note $\tau = 2\eta\, {}^n/_d$ and $w_\tau, b_\tau$ the trained parameters.

Vincent et al., *Stacked denoising AEs:Learning useful representations in a deep net. with a local denoising criterion,* JMLR 2010
15

$t = 0$

$t = 1$

$X_0 \sim \mathcal{N}(0, \mathbb{I}_d)$

$X_1 \sim \rho$

$$\frac{d}{dt} X_t = \left( \dot{\beta_t} - \frac{\dot{\alpha_t}}{\alpha_t} \beta_t + \epsilon_t \frac{\beta_t}{\alpha_t^2} \right) f(t, X_t) + \left( \frac{\dot{\alpha_t}}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2} \right) X_t + \sqrt{2\epsilon_t} dW_t$$

16

$t = 0$

$t = 1$

$X_0 \sim \mathcal{N}(0, \mathbb{I}_d)$

$X_1 \sim \rho$    $X_1 \sim \hat{\rho}_\tau$

$$\frac{d}{dt} X_t = \left( \dot{\beta}_t - \frac{\dot{\alpha}_t}{\alpha_t} \beta_t + \epsilon_t \frac{\beta_t}{\alpha_t^2} \right) f(t, X_t) + \left( \frac{\dot{\alpha}_t}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2} \right) X_t + \sqrt{2\epsilon_t} dW_t$$

$$\frac{d}{dt} X_t = \left( \dot{\beta}_t - \frac{\dot{\alpha}_t}{\alpha_t} \beta_t + \epsilon_t \frac{\beta_t}{\alpha_t^2} \right) f_{b_\tau, w_\tau}(X_t) + \left( \frac{\dot{\alpha}_t}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2} \right) X_t + \sqrt{2\epsilon_t} dW_t$$

*Using the trained AE in the generative SDE*

$$\hat{\rho}_\tau(t) = \mathrm{Law}[X_t] \ ?$$

17

Gaussian mixture supported on *a low-dimensional latent manifold*

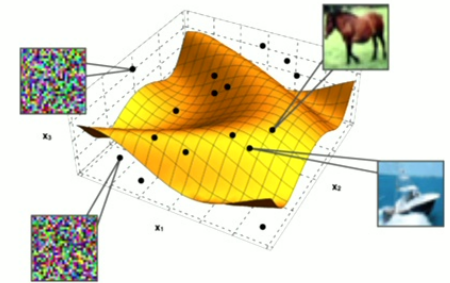$$\rho = \int_{\mathbb{R}^\kappa} d\pi(c)\mathcal{N}(\mu(c), \Sigma(c))$$

Figure from Goldt et al., *Modelling the influence of data structure in learning in neural networks: the hidden manifold model*, PRX 2020.

Tenenbaum., Silva and Langford, *A global geometric framework for nonlinear dimensionality reduction*. science, 2000
Weinberger and Saul, *Unsupervised learning of image manifolds by semidefinite programming*. Int. journal of computer vision, 2006.

18

Gaussian mixture supported on *a low-dimensional latent manifold*

$$\rho = \int_{\mathbb{R}^{\kappa}} d\pi(c)\mathcal{N}(\mu(c), \Sigma(c))$$



centroids      $\mu: \mathbb{R}^{\kappa} \to \mathbb{R}^{d}$       $\exists\, D > 0$, w.p. 1, $\|\mu(c)\| \leq D$

$K = \dim \mathrm{span}\{\mu(c)\}_c$ is low

covariances     $\Sigma: \mathbb{R}^{\kappa} \to \mathcal{S}^{d}(\mathbb{R})$     assumed jointly diagonalizable, with a well-defined joint limiting spectral density.
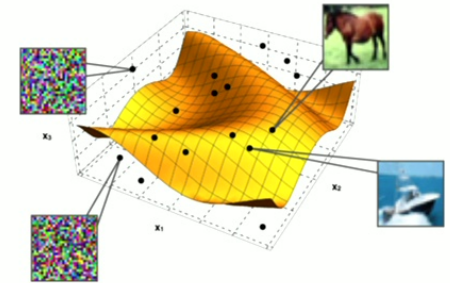
Figure from Goldt et al., *Modelling the influence of data structure in learning in neural networks: the hidden manifold model*, PRX 2020.

Tenenbaum., Silva and Langford, *A global geometric framework for nonlinear dimensionality reduction*. science, 2000
Weinberger and Saul, *Unsupervised learning of image manifolds by semidefinite programming*. Int. journal of computer vision, 2006.

19

Gaussian mixture supported on *a low-dimensional latent manifold*

$$\rho = \int_{\mathbb{R}^{\kappa}} d\pi(c)\mathcal{N}(\mu(c), \Sigma(c))$$

centroids        $\mu: \mathbb{R}^{\kappa} \to \mathbb{R}^d$      $\exists\, D > 0,\ \text{w.p. } 1,\ \|\mu(c)\| \le D$

$K = \dim \text{span}\{\mu(c)\}_c$ is low

covariances     $\Sigma: \mathbb{R}^{\kappa} \to S^d(\mathbb{R})$   assumed jointly diagonalizable, with a well-defined joint limiting spectral density.

Figure from Goldt et al., *Modelling the influence of data structure in learning in neural networks: the hidden manifold model*, PRX 2020.

*Average extension of the density*    $\Lambda = \int d\pi(c)\,\frac{1}{d}\text{Tr}[\Sigma(c)]$

Tenenbaum., Silva and Langford, *A global geometric framework for nonlinear dimensionality reduction*. science, 2000
Weinberger and Saul, *Unsupervised learning of image manifolds by semidefinite programming*. Int. journal of computer vision, 2006.

20

*Target density*
$$\rho = \int_{\mathbb{R}^\kappa} d\pi(c)\mathcal{N}(\mu(c), \Sigma(c))$$

Architecture
$$f_{b,w}(x) = bx + \frac{w}{\sqrt{d}}\sigma\left(\frac{w^\top x}{\sqrt{d}}\right)$$

Learning
$$b_{\mu+1} - b_\mu = -\frac{\eta}{d^2}\left(\partial_b\mathbb{E}_t \left\|x_1^\mu - f_{b_\mu,w_\mu}(\alpha_t x_0^\mu + \beta_t x_1^\mu)\right\|^2\right)$$
$$w_{\mu+1} - w_\mu = -\eta\left(\nabla_w\mathbb{E}_t \left\|x_1^\mu - f_{b_\mu,w_\mu}(\alpha_t x_0^\mu + \beta_t x_1^\mu)\right\|^2 + \lambda/dw_\mu\right)$$

for a time $\tau$

Sampling
$$\frac{d}{dt}X_t = \Gamma_t \frac{w_\tau}{\sqrt{d}}\sigma\left(\frac{w_\tau^\top X_t}{\sqrt{d}}\right) + \Delta_t^\tau X_t + \sqrt{2\epsilon_t}dW_t$$

for a time $t$

$$\text{with} \qquad \Gamma_t = \left(\dot{\beta}_t - \frac{\dot{\alpha}_t}{\alpha_t}\beta_t + \epsilon_t\frac{\beta_t}{\alpha_t^2}\right) \qquad \Delta_t^\tau = b_\tau\Gamma_t + \frac{\dot{\alpha}_t}{\alpha_t} - \frac{\epsilon_t}{\alpha_t^2}$$

*Generated density*   $\boxed{\hat{\rho}_\tau(t) = \text{Law}[X_t]}$

21

Asymptotic limit    $n, d \to \infty$ with $n/d, \kappa, r, D, K = \Theta_d(1)$

**Finite width, large amount of data**, large dimension

Saad and Solla, *Exact solution for on-line learning in multilayer neural networks*, PRL 1995,
...

Gabrié, Mean-Field inference methods for neural networks, J. Phys. A 2020.
**HC,** High-dimensional learning of narrow networks, J.Stat Mech 2025

22

*Tight characterization of low-dimensional projections of the generated density $\hat{\rho}_\tau(t)$*

Consider a low-dimensional subspace $\mathcal{E} \subset \mathbb{R}^d$ , with $\dim\mathcal{E} = R = \Theta_d(1)$. The distribution of the $R-$dimensional projection $\Pi_\mathcal{E} X_t$ is given by

$$\Pi_\mathcal{E} X_t =^d \Theta_\tau^\top Q_\tau^+ Z_t + Y_t$$

23

*Tight characterization of low-dimensional projections of the generated density* $\hat{\rho}_\tau(t)$

Consider a low-dimensional subspace $\mathcal{E} \subset \mathbb{R}^d$ , with $\dim\mathcal{E} = R = \Theta_d(1)$. The distribution of the $R$ −dimensional projection $\Pi_{\mathcal{E}} X_t$ is given by

$$\Pi_{\mathcal{E}} X_t =^d \Theta_\tau^\top Q_\tau^+ Z_t + Y_t$$
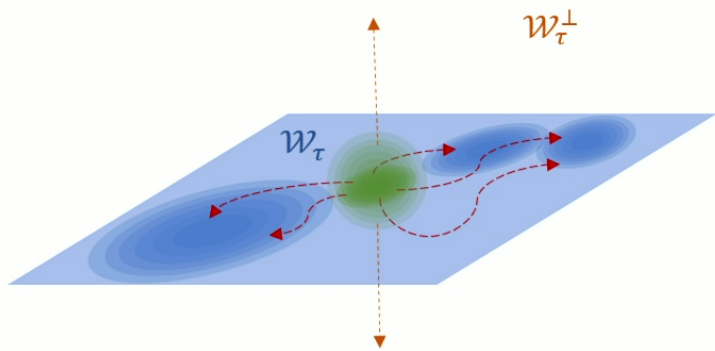
Where:

▶  $Y_t \in \mathbb{R}^R$ is Gaussian

$$Y_t \sim \mathcal{N}\left(0_R, e^{2\int_0^t \Delta_s^\tau ds}\left[1 + 2\int_0^t \epsilon_s e^{-2\int_0^s \Delta_z^\tau dz} ds\right](\mathbb{I}_R - \Theta_\tau^\top Q_\tau^+ \Theta_\tau)\right)$$

24

*Tight characterization of low-dimensional projections of the generated density $\hat{\rho}_\tau(t)$*

Consider a low-dimensional subspace $\mathcal{E} \subset \mathbb{R}^d$ , with $\dim\mathcal{E} = R = \Theta_d(1)$. The distribution of the $R$ −dimensional projection $\Pi_{\mathcal{E}} X_t$ is given by

$$\Pi_{\mathcal{E}} X_t =^d \Theta_\tau^\top Q_\tau^+ Z_t + Y_t$$

Where:

▶  $Y_t \in \mathbb{R}^R$ is Gaussian

$$Y_t \sim \mathcal{N}\left(0_R, e^{2\int_0^t \Delta_s^\tau ds}\left[1 + 2\int_0^t \epsilon_s e^{-2\int_0^s \Delta_z^\tau dz} ds\right](\mathbb{I}_R - \Theta_\tau^\top Q_\tau^+ \Theta_\tau)\right)$$

▶  $Z_t \in \mathbb{R}^r$ is distributed as the solution of the SDE

$$\frac{d}{dt}Z_t = \Delta_t^\tau Z_t + \Gamma_t Q_\tau \sigma(Z_t) + \sqrt{2\epsilon_t} Q_\tau^{1/2} W_t$$

From initialization $Z_0 \sim \mathcal{N}(0_r, Q_\tau)$

25

---

*Tight characterization of low-dimensional projections of the generated density $\hat{\rho}_\tau(t)$*

Consider a low-dimensional subspace $\mathcal{E} \subset \mathbb{R}^d$, with $\dim\mathcal{E} = R = \Theta_d(1)$. The distribution of the $R$-dimensional projection $\Pi_\mathcal{E}X_t$ is given by

$$\Pi_\mathcal{E}X_t =^d \Theta_\tau^\top Q_\tau^+ Z_t + Y_t$$

Where:

▶  $Y_t \in \mathbb{R}^R$ is Gaussian

$$Y_t \sim \mathcal{N}\left( 0_R, e^{2\int_0^t \Delta_s^\tau ds}\left[1 + 2\int_0^t \epsilon_s e^{-2\int_0^s \Delta_z^\tau dz}ds\right](\mathbb{I}_R - \Theta_\tau^\top Q_\tau^+ \Theta_\tau) \right)$$

▶  $Z_t \in \mathbb{R}^r$ is distributed as the solution of the SDE

$$\frac{d}{dt}Z_t = \Delta_t^\tau Z_t + \Gamma_t Q_\tau \sigma(Z_t) + \sqrt{2\epsilon_t}Q_\tau^{1/2}W_t$$

From initialization $Z_0 \sim \mathcal{N}(0_r, Q_\tau)$

▶  The parameters $\Theta_\tau \in \mathbb{R}^{r\times R}, Q_\tau \in \mathbb{R}^{r\times r}$ are the solutions of a set of 5 coupled **low-dimensional** deterministic **ODEs**.

26

$$\frac{d}{dt}X_t = \boxed{\Gamma_t \frac{w_\tau}{\sqrt{d}} \sigma\left(\frac{w_\tau{}^\top X_t}{\sqrt{d}}\right)} + \Delta_t^\tau X_t + \sqrt{2\epsilon_t}\,dW_t$$

Non-linear transport in       Linear in $\mathcal{W}_\tau^\perp$
$\mathcal{W}_\tau = \mathrm{span}(\{w_i\}_{i=1}^r)$

27

$$\frac{d}{dt}X_t = \boxed{\Gamma_t \frac{w_\tau}{\sqrt{d}} \sigma\left(\frac{w_\tau{}^\top X_t}{\sqrt{d}}\right)} + \Delta_t^\tau X_t + \sqrt{2\epsilon_t}\, dW_t$$

Non-linear transport in         Linear in $\mathcal{W}_\tau^\perp$
$\mathcal{W}_\tau = \mathrm{span}(\{w_i\}_{i=1}^r)$

Dynamics of $Z_t = \frac{w_\tau{}^\top X_t}{\sqrt{d}}$

$$\frac{d}{dt}Z_t = \Delta_t^\tau Z_t + \Gamma_t \frac{w_\tau^\top w_\tau}{d} \sigma(Z_t) + \sqrt{2\epsilon_t}\left(\frac{w_\tau^\top w_\tau}{d}\right)^{\frac{1}{2}} W_t$$

28

$$\frac{d}{dt}X_t = \boxed{\Gamma_t \frac{w_\tau}{\sqrt{d}} \sigma\left(\frac{w_\tau^\top X_t}{\sqrt{d}}\right)} + \Delta_t^\tau X_t + \sqrt{2\epsilon_t}dW_t$$

Non-linear transport in              Linear in $\mathcal{W}_\tau^\perp$
$\mathcal{W}_\tau = \mathrm{span}(\{w_i\}_{i=1}^r)$

Dynamics of $Z_t = \frac{w_\tau^\top X_t}{\sqrt{d}}$

$$\frac{d}{dt}Z_t = \Delta_t^\tau Z_t + \Gamma_t \frac{w_\tau^\top w_\tau}{d} \sigma(Z_t) + \sqrt{2\epsilon_t}\left(\frac{w_\tau^\top w_\tau}{d}\right)^{\frac{1}{2}} W_t$$

The SGD dynamics of the summary statistic $Q_\tau = \frac{w_\tau^\top w_\tau}{d}$ (and others) self-average and can be characterized in closed-form by a set of low-dimensional ODEs.

Saad and Solla, *Exact solution for on-line learningin multilayer neural networks*, PRL 1995

29

Intuition:

• The network **identifies** a $r-$dimensional $\mathcal{W}_\tau$ subspace where the target $\rho$ has important structure, and implements a non-linear transport.

• It approximates $\rho$ in the orthogonal space by an **isotropic Gaussian**, whose variance is tuned by the skip connection strength.

Special case: linear networks $\sigma(x) = x$

• Linear networks approximately learn $\approx$ **principal components** $\quad \mathcal{W}_\tau \approx \mathrm{PCA}_r\left[\{x_1^\mu\}_\mu\right]$

　Pretorius et al,. *Learning dynamics of linear denoising autoencoders*. ICML, 2018. , …..

• The linear diffusion model does a Gaussian approximation in the principal space. In the orthogonal space, approximates by an isotropic Gaussian.

30

$\sigma = $ ReLU activation,
$r = 4$ hidden units
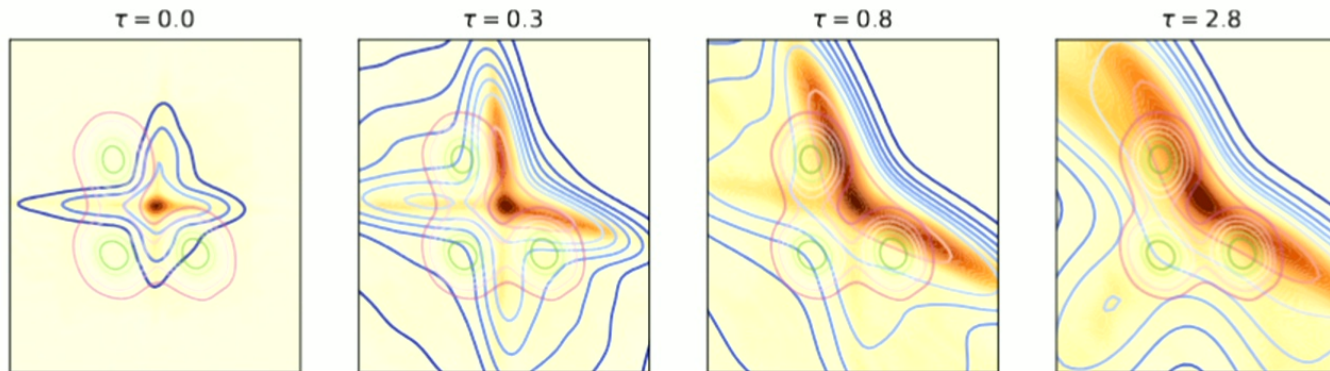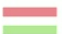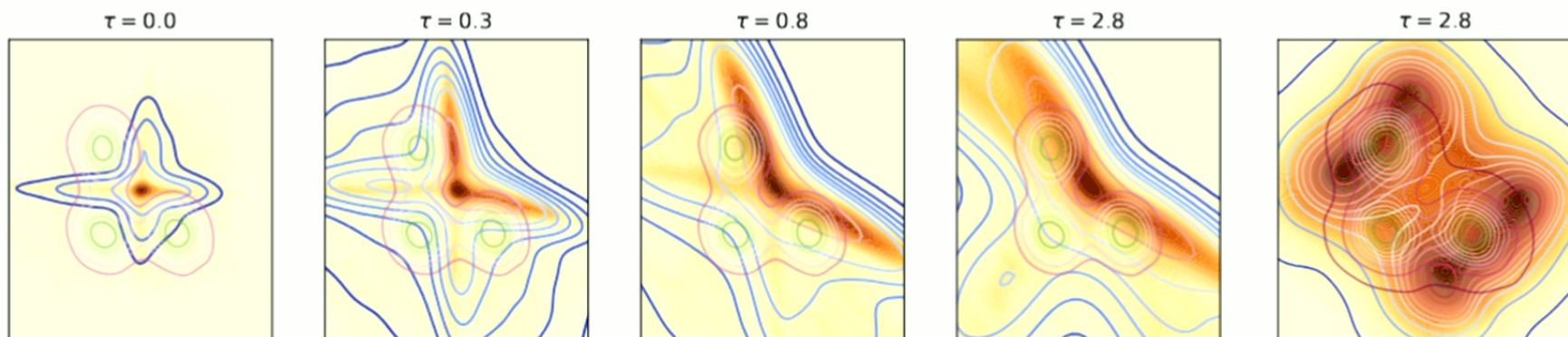
$\tau = 0.0$



━━  Target density $\rho$          ━━  (**theory**) Generated density $\hat{\rho}$          ▉  (**exp**) Generated density $\hat{\rho}$

31

$\sigma = $ ReLU activation,
$r = 4$ hidden units

$\tau = 0.0$　　　　$\tau = 0.3$

| ▬ | Target density $\rho$ | ▬ | (**theory**) Generated density $\hat{\rho}$ | ▬ | (**exp**) Generated density $\hat{\rho}$ |

32

$\sigma$ = ReLU activation,
$r = 4$ hidden units
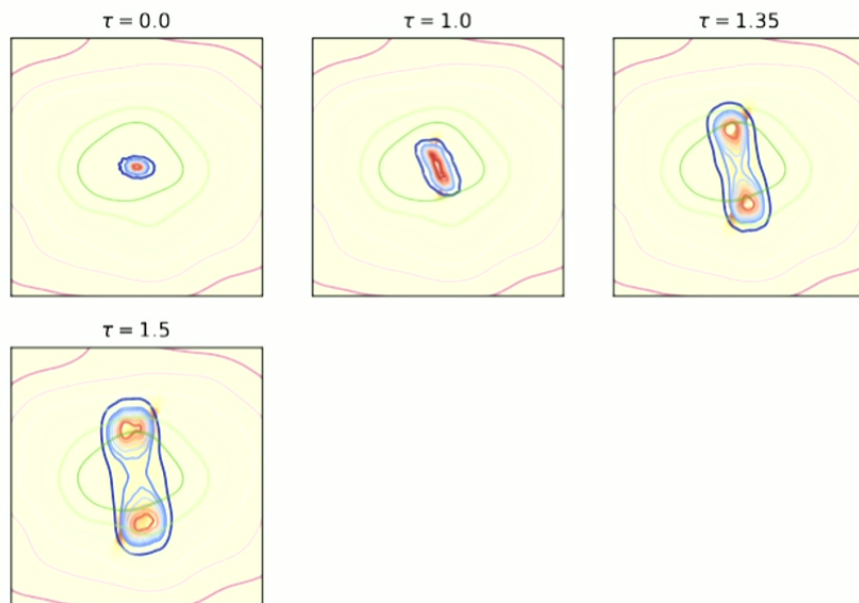


$\tau = 0.0$    $\tau = 0.3$    $\tau = 0.8$

Target density $\rho$    (**theory**) Generated density $\hat{\rho}$    (**exp**) Generated density $\hat{\rho}$

33

$\sigma = \text{ReLU activation,}$
$r = 4$ hidden units

$\tau = 0.0$      $\tau = 0.3$      $\tau = 0.8$      $\tau = 2.8$



— Target density $\rho$     — **(theory)** Generated density $\hat{\rho}$     ▇ **(exp)** Generated density $\hat{\rho}$

34

$\sigma = $ ReLU activation,
$r = 4$ hidden units

$\sigma = $ tanh activation,
$r = 2$ hidden units



— Target density $\rho$     — (**theory**) Generated density $\hat{\rho}$     ■ (**exp**) Generated density $\hat{\rho}$

35

$\tau = 0.0$    $\tau = 1.0$    $\tau = 1.35$
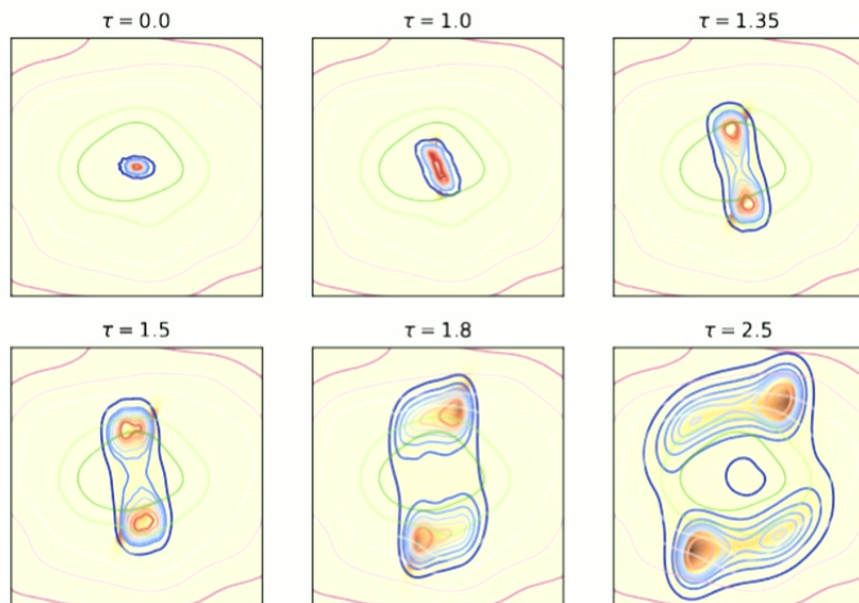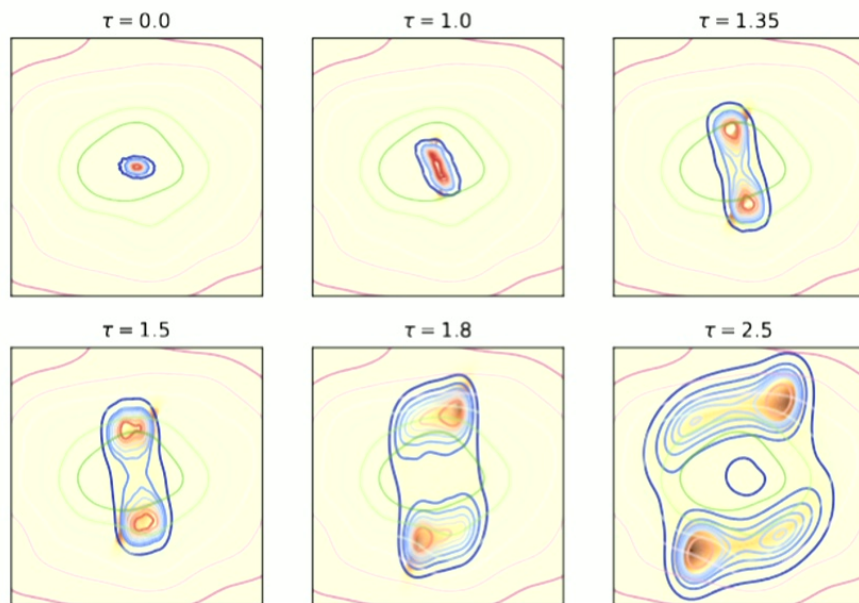
$\tau = 1.5$

$\sigma = $ tanh activation,
$r = 2$ hidden units
Gaussian $\rho$ with MNIST covariance

━━  Target density $\rho$    ━━  (**theory**) Generated density $\hat{\rho}$    ▮  (**exp**) Generated density $\hat{\rho}$

37

$\tau = 0.0$

$\sigma = $ tanh activation,
$r = 2$ hidden units
Gaussian $\rho$ with MNIST covariance

▬ ▬ Target density $\rho$        ▬ ▬ (**theory**) Generated density $\hat{\rho}$        ▮ (**exp**) Generated density $\hat{\rho}$

36

$\sigma = $ tanh activation,
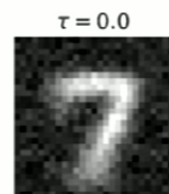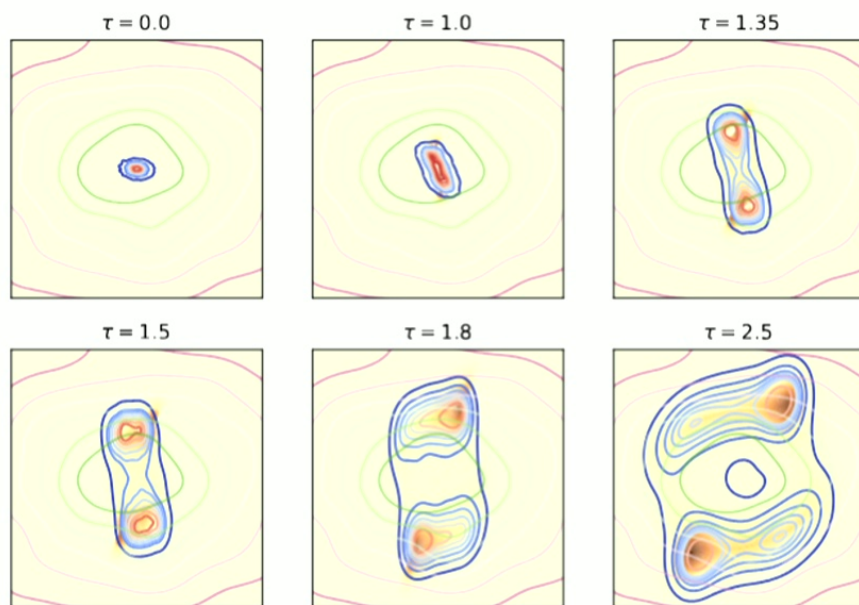$r = 2$ hidden units
Gaussian $\rho$ with MNIST covariance

— Target density $\rho$ 　　 — (**theory**) Generated density $\hat{\rho}$ 　　 ■ (**exp**) Generated density $\hat{\rho}$

38

$\sigma =$ tanh activation,
$r = 2$ hidden units
Gaussian $\rho$ with MNIST covariance

$\tau = 0.0$

— Target density $\rho$    — (**theory**) Generated density $\hat\rho$    ▮ (**exp**) Generated density $\hat\rho$
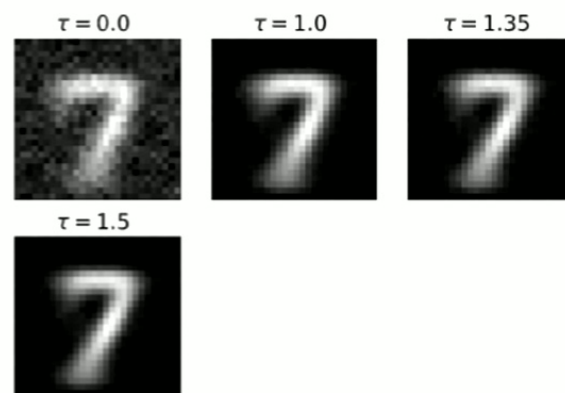
39

$\sigma$ = tanh activation,
$r$ = 2 hidden units
Gaussian $\rho$ with MNIST covariance



—— Target density $\rho$   —— (**theory**) Generated density $\hat{\rho}$   ▮ (**exp**) Generated density $\hat{\rho}$
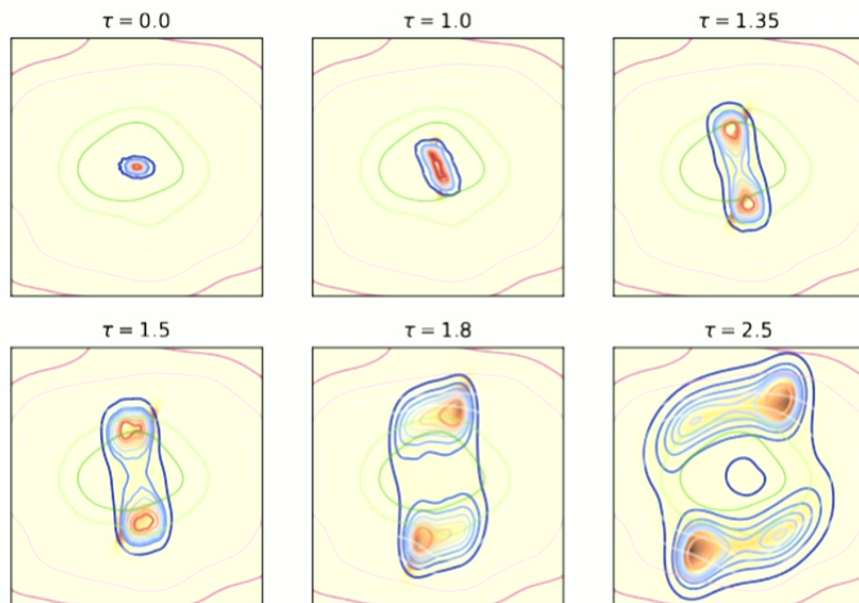
40

$\sigma = $ tanh activation,
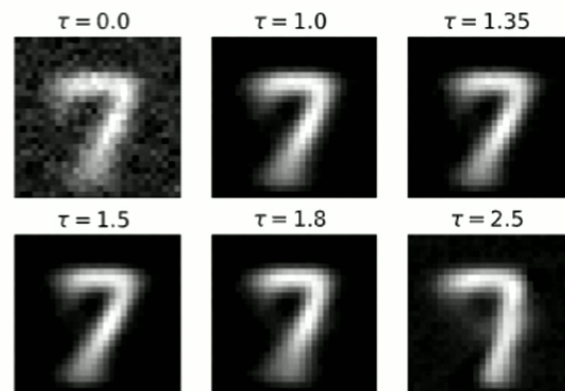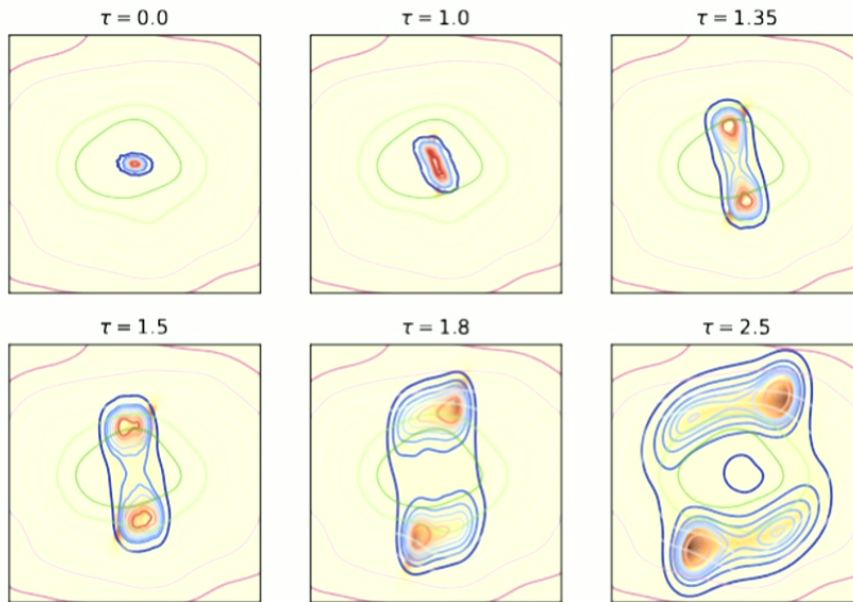$r = 2$ hidden units
Gaussian $\rho$ with MNIST covariance



— Target density $\rho$    — (**theory**) Generated density $\hat{\rho}$    ▮ (**exp**) Generated density $\hat{\rho}$

41

$\tau = 0.0$  $\tau = 1.0$  $\tau = 1.35$
$\tau = 1.5$  $\tau = 1.8$  $\tau = 2.5$

Closed-form expression for the trained skip connection

$$b_\tau = \frac{\Lambda \mathbb{E}_t[\beta_t]\left[1 - (1 - b_0)e^{-\left(\Lambda \mathbb{E}_t[\beta_t^2] + \mathbb{E}_t[\alpha_t^2]\right)\tau}\right]}{\Lambda \mathbb{E}_t[\beta_t^2] + \mathbb{E}_t[\alpha_t^2]}$$

Average cov. eigenvalue $\quad \Lambda = \int d\pi(c) \frac{1}{d} \mathrm{Tr}[\Sigma(c)]$

is typically **small** in real datasets, causing
$\approx$ ***mode collapse***

Goodfellow et al., *Generative adversarial nets.* NeurIPS 2014.

—— Target density $\rho$     —— (**theory**) Generated density $\hat{\rho}$     ■ (**exp**) Generated density $\hat{\rho}$
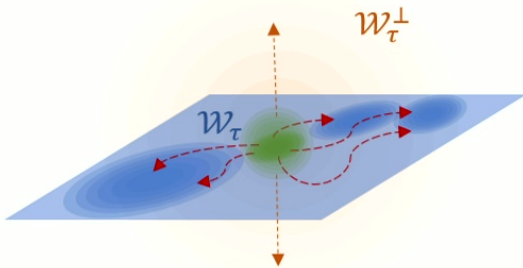
42

Can this bias be ***aggravated*** when using synthetic data to train a new generative model ?

$$\rho \rightarrow \hat{\rho}^{(1)} \rightarrow \hat{\rho}^{(2)} \rightarrow \cdots \rightarrow \hat{\rho}^{(g)}$$
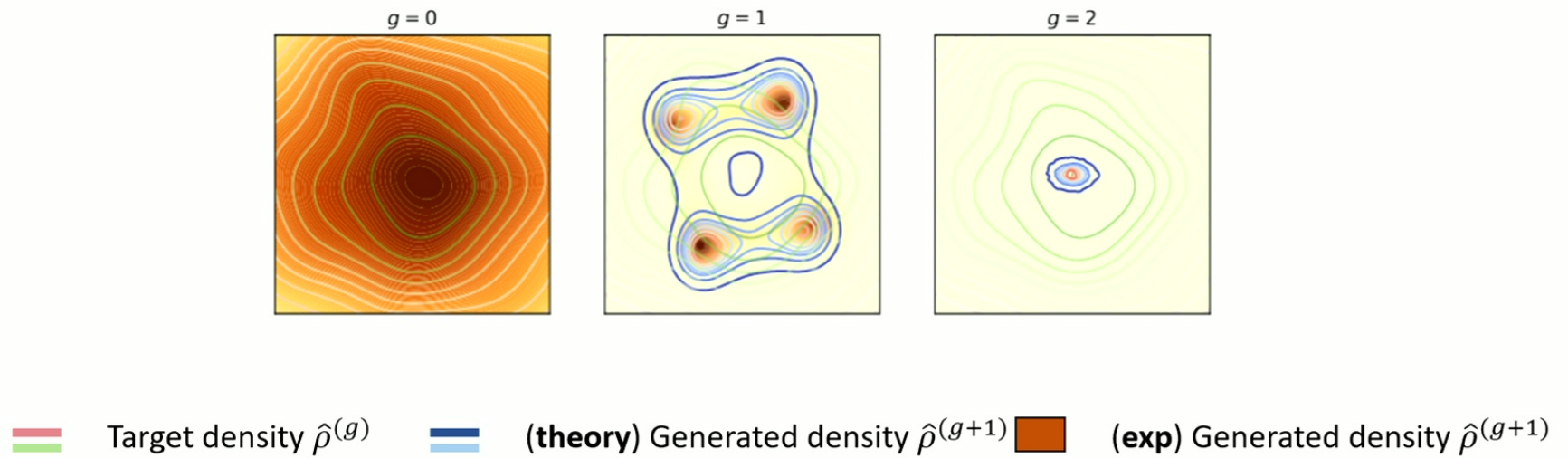
43

Can this bias be ***aggravated*** when using synthetic data to train a new generative model ?

$$\rho \to \hat{\rho}^{(1)} \to \hat{\rho}^{(2)} \to \cdots \to \hat{\rho}^{(g)}$$



**Remark**: *Manifold form of the generated density*

$\hat{\rho}^{(1)}$ is still of the form $\int d\pi(c)\mathcal{N}\big(\mu(c), \Sigma(c)\big)$, with $\mu(c) = c$ and

$$\pi = \Pi_{\mathcal{W}_\tau}\,\hat{\rho}^{(1)}$$

$$\Sigma(c) = e^{2\int_0^t \Delta_s^\tau ds}\left[1 + 2\int_0^t \epsilon_s e^{-2\int_0^s \Delta_z^\tau dz}\,ds\right]\Pi_{\mathcal{W}_\tau^\perp}$$

Thus the analysis *carries over iteratively* to generations $\hat{\rho}^{(2)}$, ...

44

g = 0          g = 1          g = 2

Target density $\hat{\rho}^{(g)}$     (**theory**) Generated density $\hat{\rho}^{(g+1)}$     (**exp**) Generated density $\hat{\rho}^{(g+1)}$

Shumailov et al., *Ai models collapse when trained onrecursively generated data.* Nature, 2024

45

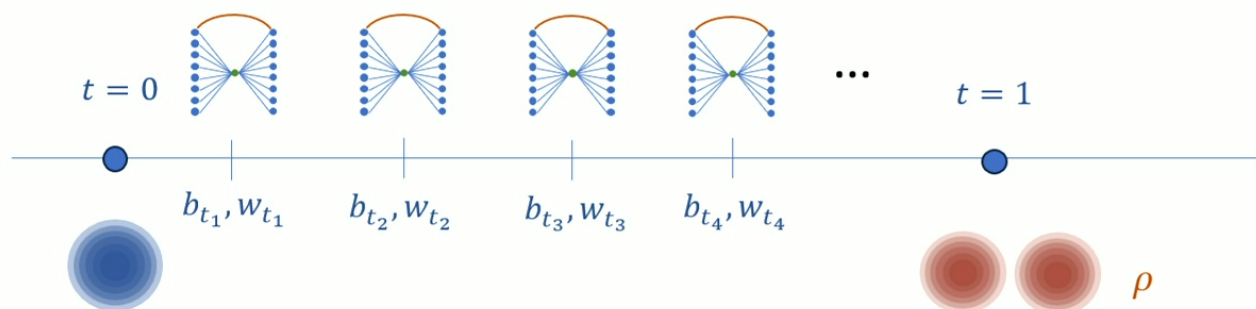Binary, isotropic Gaussian mixture    $\rho = {}^{1}/_{2}\,\mathcal{N}\!\left(-\mu, \sigma^2\,\mathbb{I}_d\right) + {}^{1}/_{2}\,\mathcal{N}\!\left(+\mu, \sigma^2\mathbb{I}_d\right)$

$t = 0$                      $t = 1$

$\rho$

46

Binary, isotropic Gaussian mixture　　　$\rho = {}^{1}/_{2}\,\mathcal{N}\left(-\mu, \sigma^2\,\mathbb{I}_d\right) + {}^{1}/_{2}\,\mathcal{N}\left(+\mu, \sigma^2\mathbb{I}_d\right)$

At *each sampling time*, train a **separate** AE with $r = 1$ hidden unit and $\sigma =$ sign activation

$$b_t, w_t = \operatorname{argmin}_{\theta \in \mathbb{R}^{d \times r}} \sum_{\mu=1}^{n} \left\| f_{b,w}\left(\alpha_t x_0^{\mu} + \beta_t x_1^{\mu}\right) - x_1^{\mu} \right\|_2^2 + \lambda \|w\|^2$$



47

Binary, isotropic Gaussian mixture　　　$\rho = {}^{1}\!/_{2}\,\mathcal{N}\big(-\mu, \sigma^{2}\,\mathbb{I}_{d}\big) + {}^{1}\!/_{2}\,\mathcal{N}\big(+\mu, \sigma^{2}\mathbb{I}_{d}\big)$

At *each sampling time*, train a **separate** AE with $r = 1$ hidden unit and $\sigma =$ sign activation

$$b_t, w_t = \mathrm{argmin}_{\theta \in \mathbb{R}^{d \times r}} \sum_{\mu=1}^{n} \left\| f_{b,w}\big(\alpha_t x_0^{\mu} + \beta_t x_1^{\mu}\big) - x_1^{\mu} \right\|_2^2 + \lambda \|w\|^2$$

Sampling :　　　$\dfrac{d}{dt} X_t = \left( \dot{\beta}_t - \dfrac{\dot{\alpha}_t}{\alpha_t} \beta_t \right) f_{b_t, w_t}(X_t) + \dfrac{\dot{\alpha}_t}{\alpha_t} X_t$

## *Closed form characterization of the dynamics*

In the asymptotic limit $d \to \infty$ with $n = \Theta_d(1)$, $\|\mu\| = \Theta_d(\sqrt{d})$, the sampling dynamic is non-linear in $\text{span}(\mu, \xi, \eta)$ where

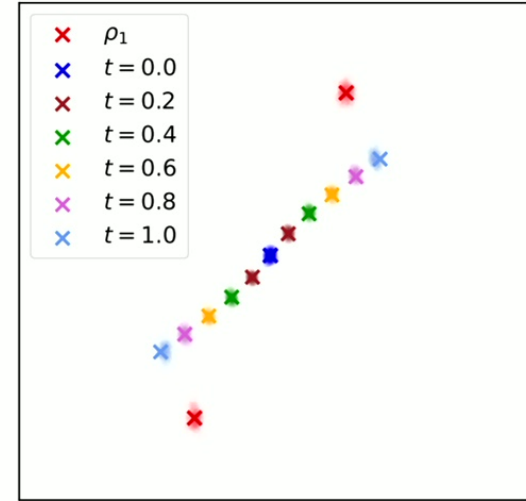$$\xi \equiv \sum_{\mu=1}^{n} s^\mu x_0^\mu, \qquad\qquad \eta \equiv \sum_{\mu=1}^{n} s^\mu (x_1^\mu - s^\mu \mu),$$

The coordinates $M_t, Q_t^\xi, Q_t^\eta$ of a sample $X_t$ follow the ODEs

$$\begin{cases} \frac{d}{dt} M_t = \frac{\left(\dot{\beta}(t)\beta(t)(\lambda(1+\sigma^2)+(n-1)\sigma^2)+\dot{\alpha}(t)\alpha(t)(\lambda+n-1)\right)M_t + \left(\alpha(t)\dot{\beta}(t)-\dot{\alpha}(t)\beta(t)\right)\frac{n\alpha(t)(\lambda+n-1)}{\lambda+n}}{\alpha(t)^2(\lambda+n-1)+\beta(t)^2(\lambda(1+\sigma^2)+(n-1)\sigma^2)} \\ \frac{d}{dt} Q_t^\xi = \frac{\left(\dot{\beta}(t)\beta(t)(\lambda(1+\sigma^2)+(n-1)\sigma^2)+\dot{\alpha}(t)\alpha(t)(\lambda+n-1)\right)Q_t^\xi - \left(\alpha(t)\dot{\beta}(t)-\dot{\alpha}(t)\beta(t)\right)\frac{\beta(t)(\lambda(1+\sigma^2)+(n-1)\sigma^2)}{\lambda+n}}{\alpha(t)^2(\lambda+n-1)+\beta(t)^2(\lambda(1+\sigma^2)+(n-1)\sigma^2)} \\ \frac{d}{dt} Q_t^\eta = \frac{\left(\dot{\beta}(t)\beta(t)(\lambda(1+\sigma^2)+(n-1)\sigma^2)+\dot{\alpha}(t)\alpha(t)(\lambda+n-1)\right)Q_t^\eta + \left(\alpha(t)\dot{\beta}(t)-\dot{\alpha}(t)\beta(t)\right)\frac{\alpha(t)(\lambda+n-1)}{\lambda+n}}{\alpha(t)^2(\lambda+n-1)+\beta(t)^2(\lambda(1+\sigma^2)+(n-1)\sigma^2)} \end{cases}$$

The component $X_t^\perp$ orthogonal to $\text{span}(\mu, \xi, \eta)$ evolves linearly

$$\frac{d}{dt} X_t^\perp = \frac{\left(\dot{\beta}(t)\beta(t)(\lambda(1+\sigma^2)+(n-1)\sigma^2)+\dot{\alpha}(t)\alpha(t)(\lambda+n-1)\right)}{\alpha(t)^2(\lambda+n-1)+\beta(t)^2(\lambda(1+\sigma^2)+(n-1)\sigma^2)} X_t^\perp$$



**HC**, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity*, **ICLR 2024**     49
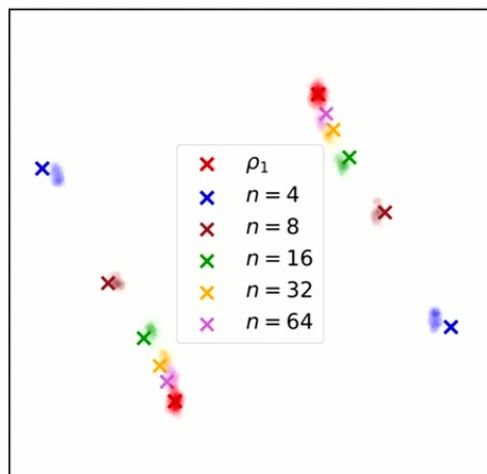
*Corollary*

The *mixture Wasserstein distance* between the target $\rho$ and the generated density $\hat{\rho}$ decays as

$$\mathrm{M}\mathcal{W}_2[\rho, \hat{\rho}] = O\left(\frac{1}{n}\right)$$

**HC**, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity*, **ICLR 2024**    50

*Corollary*

The *mixture Wasserstein distance* between the target $\rho$ and the generated density $\hat{\rho}$ decays as

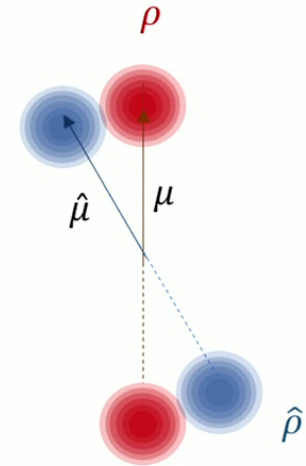$$\mathrm{M}\mathcal{W}_2[\rho, \hat{\rho}] = O\left(\frac{1}{n}\right)$$



HC, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity,* **ICLR 2024**　　　51

## *Corollary*

The *mixture Wasserstein distance* between the target $\rho$ and the generated density $\hat{\rho}$ decays as

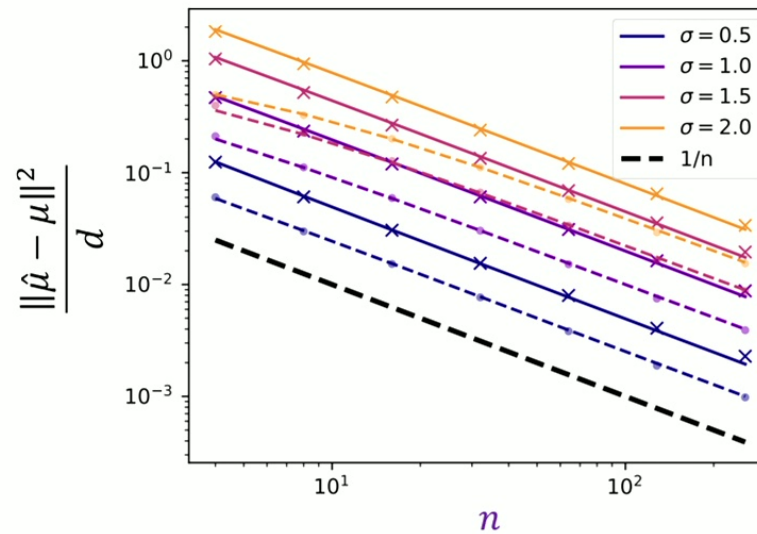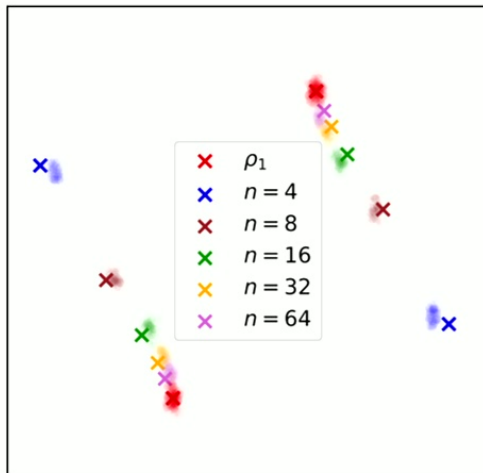$$\mathsf{M}\mathcal{W}_2[\rho, \hat{\rho}] = O\left(\frac{1}{n}\right)$$



**HC**, Krzakala, Vanden-Eijnden, Zdeborová, *Analysis of a learning a flow-based generative model from finite sample complexity,* **ICLR 2024**              52

**Intuition** : The optimal denoising function follows from *Tweedie's formula* (Empirical Bayes) and is of the **same functional form** as the AE

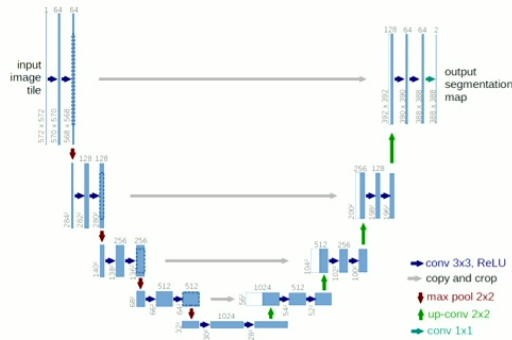$$f_t^\star(x) = \frac{\beta(t)\sigma^2}{\alpha(t)^2 + \beta(t)^2\sigma^2}x + \frac{\alpha(t)^2}{\alpha(t)^2 + \beta(t)^2\sigma^2}\mu \times \tanh\left(\frac{\beta(t)}{\alpha(t)^2 + \beta(t)^2\sigma^2}\mu^\top x\right)$$

→ The architectural bias *is **aligned** with the target distribution.*

Bradley Efron. *Tweedie's formula and selection bias*. Journal of the American Statistical Association, 2011
Robbins, Proceedings of the Third Berkeley Symposium on Mathematical Statistics and Probability, vol. 1: Contributions to the Theory of Statistics.
Koichi Miyasawa. *An empirical Bayes estimator of the mean of a normal population*. Bulletin of the International Statistical Institute, 1961

53

## Perspectives



Ronneberger, Fischer, and Brox *U-net: Convolutional networks for biomedical image segmentation*. MICCAI 2015
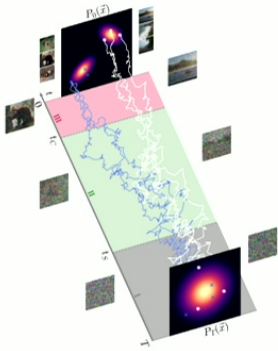
Inductive bias of **Unets**?

U-nets are suited to data with a *hierarchical structure*

Kadkhodaie et al., *Generalization in diffusion models arises from geometry-adaptive harmonic representation*, ICLR 2024
Mei, S. *U-nets as belief propagation: Efficient classification,denoising, and diffusion in generative hierarchical models*, arXiv:2404.18444, 2024.

(Recall also Alessandro's talk!)



For infinitely expressive networks who can perfectly overfit the data, dynamical transitions in the sampling process.

Biroli et al, *Dynamical Regimes of Diffusion Models,* Nature Comm. 2024

How are they altered for networks with finite expressivity?

54

# Collaborators

**Yue M Lu**
(*Harvard*)

**Cengiz Pehlevan**
(*Harvard*)

**Lenka Zdeborová**
(*EPFL*)

**Florent Krzakala**
(*EPFL*)

**Eric Vanden-Eijnden**
(*NYU*)

55

Thank you for your attention !