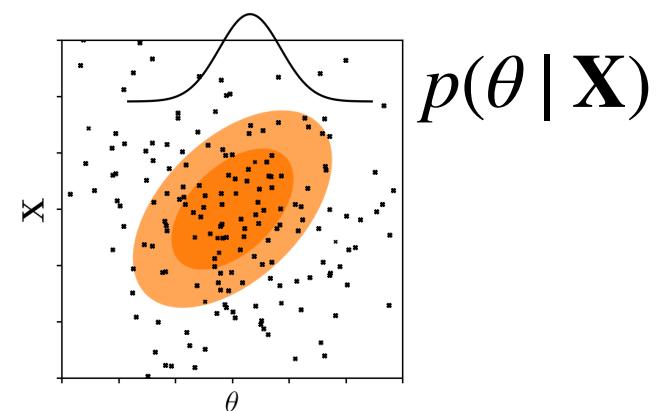


summary of first lecture

SBI is a density estimation problem -

allows to approximate a posterior distribution without an explicit likelihood using a forward model



deep learning enables density estimation at very high dimension, enabling fast inference

$$(x_0, x_1, \dots, x_N) \in X \longrightarrow p_w(X)$$

(evaluation + sampling)

density estimation with neural networks

 = neural network

$$(x_0, x_1, \dots, x_N) \in X \longrightarrow p_w(X) \longrightarrow$$

Find w so that:
 $\max \log p_w(X)$

density estimation with neural networks

■ = neural network

$$(x_0, x_1, \dots, x_N) \in X \longrightarrow p_w(X) \longrightarrow$$

Find w so that:
 $\max \log p_w(X)$

MDN

$$p_w(y | x) = L(w) = \prod_{i=1}^N \sum_{k=1}^2 \pi_w^{(k)}(x_i) \frac{1}{\sqrt{2\pi} \sigma_w^{(k)}(x_i)} \exp\left(-\frac{(y_i - f_w^{(k)}(x_i))^2}{2 [\sigma_w^{(k)}(x_i)]^2}\right)$$

a general probability function:

$$p_{\theta}(x) = \frac{\exp(f_{\theta}(x))}{Z(\theta)},$$
$$Z(\theta) = \int_{\mathcal{X}} \exp(f_{\theta}(x)) dx,$$
$$\log p_{\theta}(x) = f_{\theta}(x) - \log Z(\theta).$$

Neural network
 θ =weights

a general probability function:

$$p_\theta(x) = \frac{\exp(f_\theta(x))}{Z(\theta)},$$
$$Z(\theta) = \int_{\mathcal{X}} \exp(f_\theta(x)) dx,$$
$$\log p_\theta(x) = f_\theta(x) - \log Z(\theta).$$

Neural network
 θ =weights

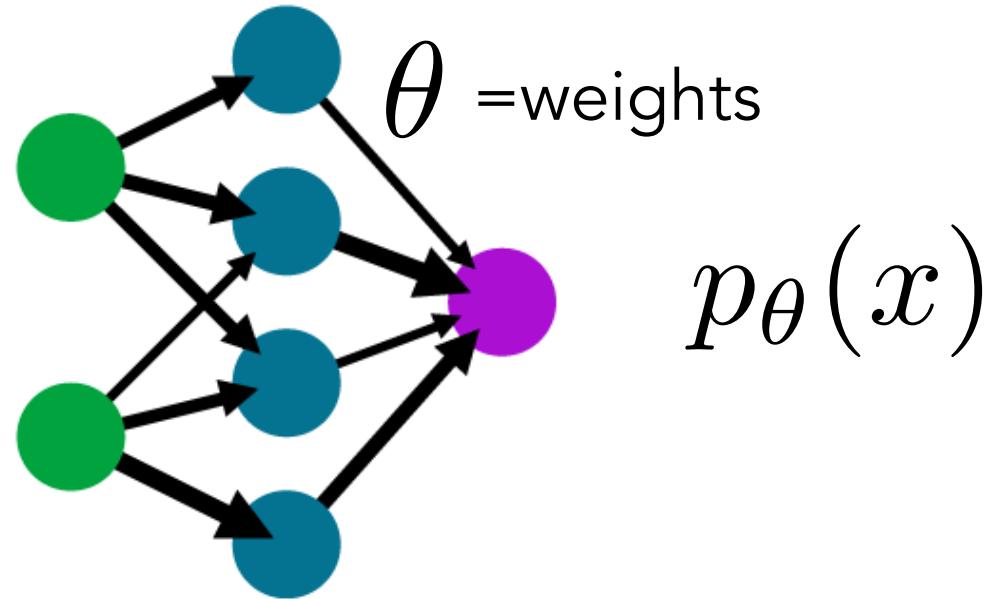
however, computing the gradient is intractable (norm. constant)

$$\begin{aligned}\nabla_\theta \log Z(\theta) &= \frac{1}{Z(\theta)} \int \exp(f_\theta(x)) \nabla_\theta f_\theta(x) dx \\ &= \int \frac{\exp(f_\theta(x))}{Z(\theta)} \nabla_\theta f_\theta(x) dx \\ &= \mathbb{E}_{x \sim p_\theta} [\nabla_\theta f_\theta(x)]\end{aligned}$$

It's an expectation under the **current model (neural network)** p_θ which itself depends on $Z(\theta)$

generative model

$$z \sim \mathcal{N}(0, 1)$$
$$p(z)$$



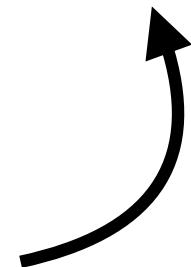
generative models are proposed solutions to make this optimization problem practical by avoiding to compute the normalization constant $Z(\theta)$.

<u>Types of generative models</u>			<u>loss</u>
Likelihood based models	directly learn the distribution's probability density function with approximations	VAEs, Normalizing Flows	$\max_{\theta} \sum_{i=1}^N \log p_{\theta}(x_i)$
Implicit generative models	the probability distribution is implicitly represented by a model of its sampling process	GANs	$\mathcal{L}_D = -\frac{1}{m} \sum_{i=1}^m [\log D(x_i) + \log (1 - D(\tilde{x}_i))], \quad \mathcal{L}_G^{NS} = -\frac{1}{m} \sum_{i=1}^m \log D(\tilde{x}_i)$
Score based models	model the gradient of the log probability density function	Score based diffusion	$\nabla_x \log p_{\theta}(x)$

High
dimensional data
(e.g. images)

Evaluation

$$p_w(X)$$



Likely not
Gaussian

$$z \sim \mathcal{N}(0, 1)$$

(Easy to sample, easy to evaluate)

Sampling

Evaluation

High
dimensional data
(e.g. images)

$$p_w(X)$$

Likely not
Gaussian

NNs == universal
approximators

$$z \sim \mathcal{N}(0, 1)$$

(Easy to sample, easy to evaluate)

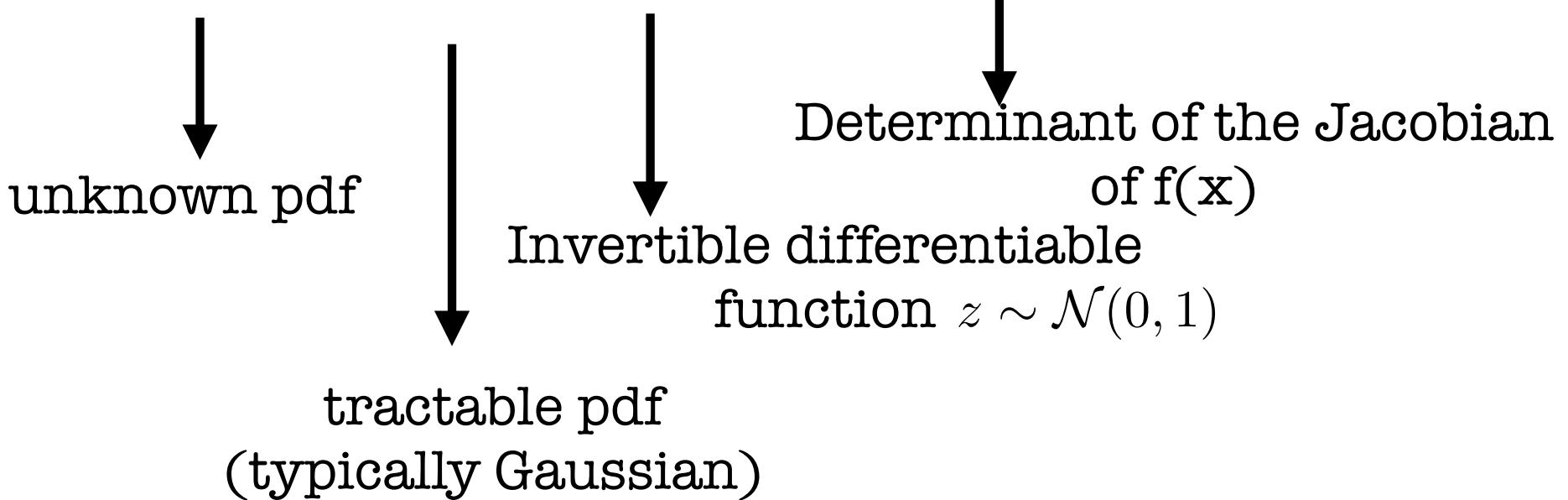
(Optimization problem)

Sampling

Normalizing Flows

Based on change of variables:

$$p_X(x) = p_Z(f(x)) |\det J_f(x)|, \quad J_f(x) = \frac{\partial f(x)}{\partial x^\top}$$

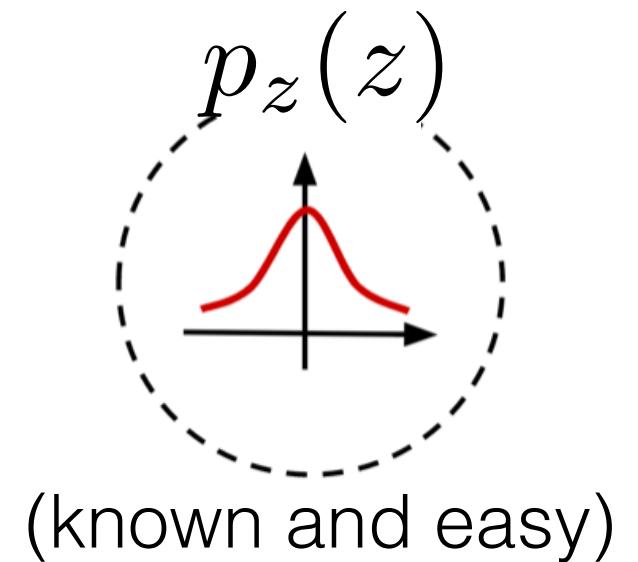
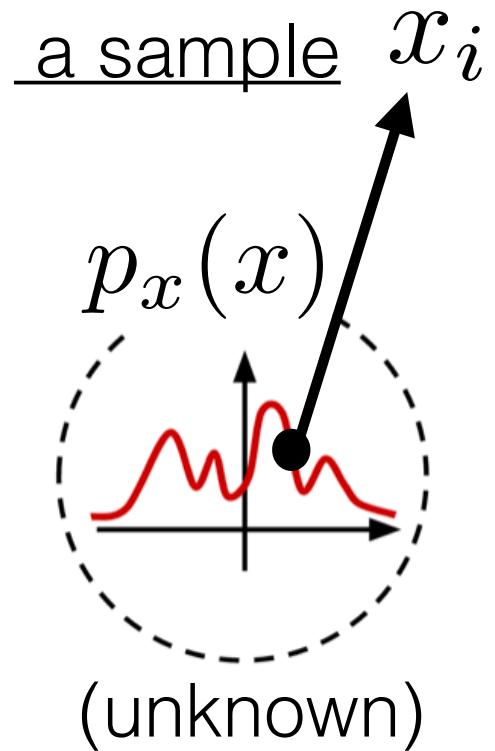


the neural network parametrizes the function f : $f_\theta(x)$

and is optimized so that:

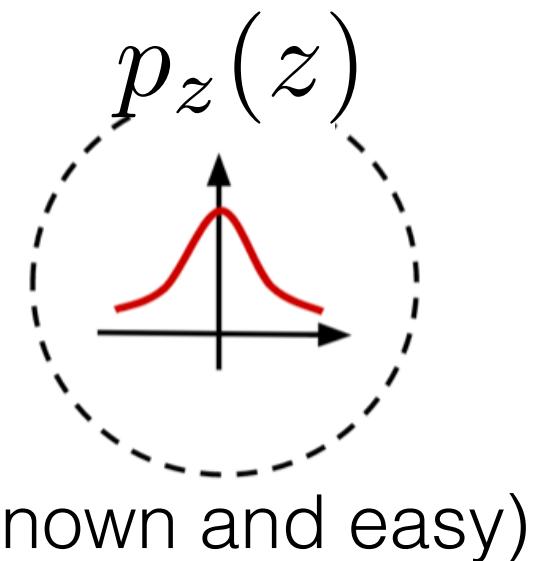
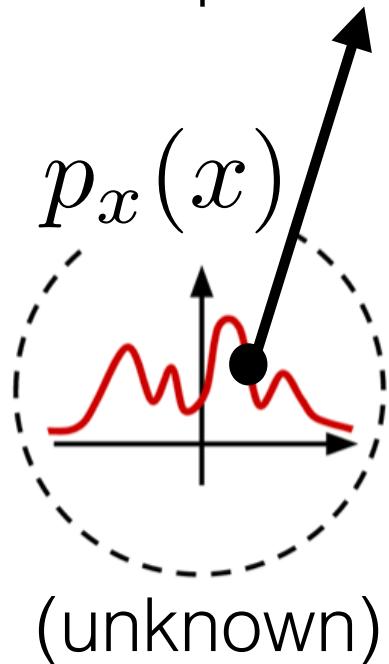
$$p_Z(f_\theta(x)) \cdot |\det J_f(x)| \quad \text{is maximum}$$
$$x \in X$$
$$\sim p_X(x)$$

training procedure

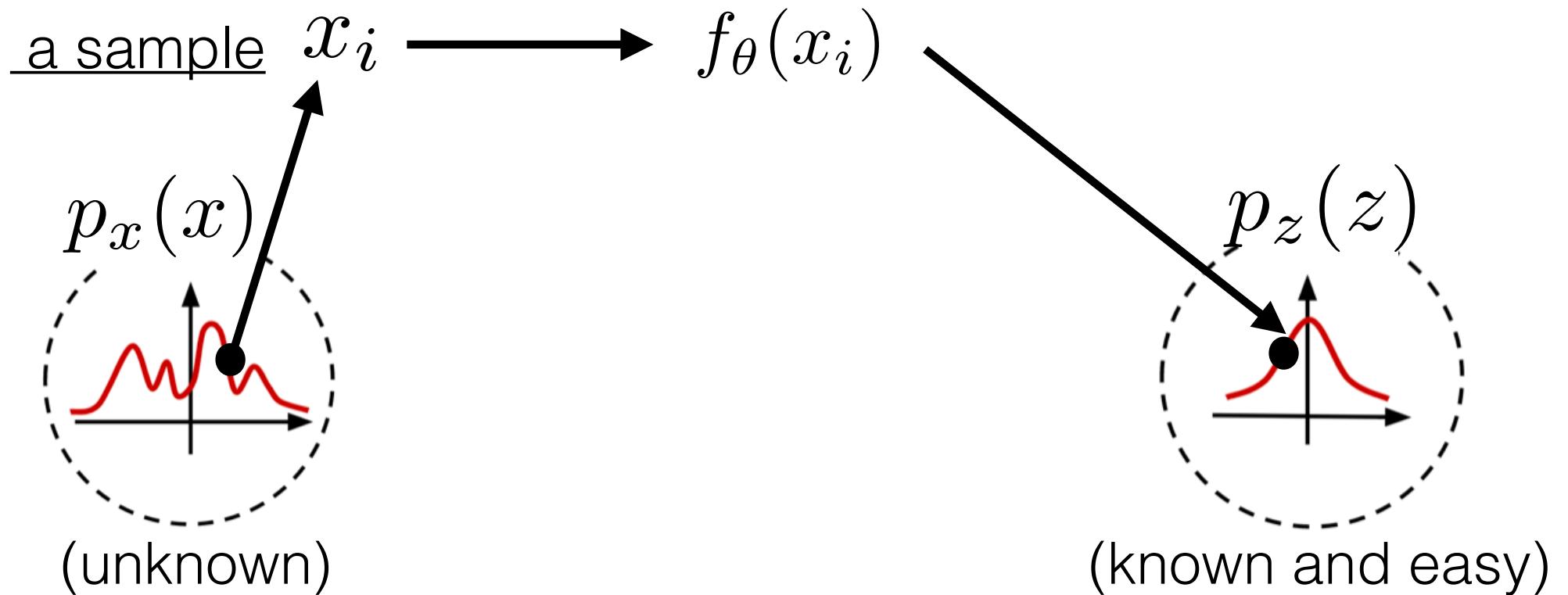


training procedure

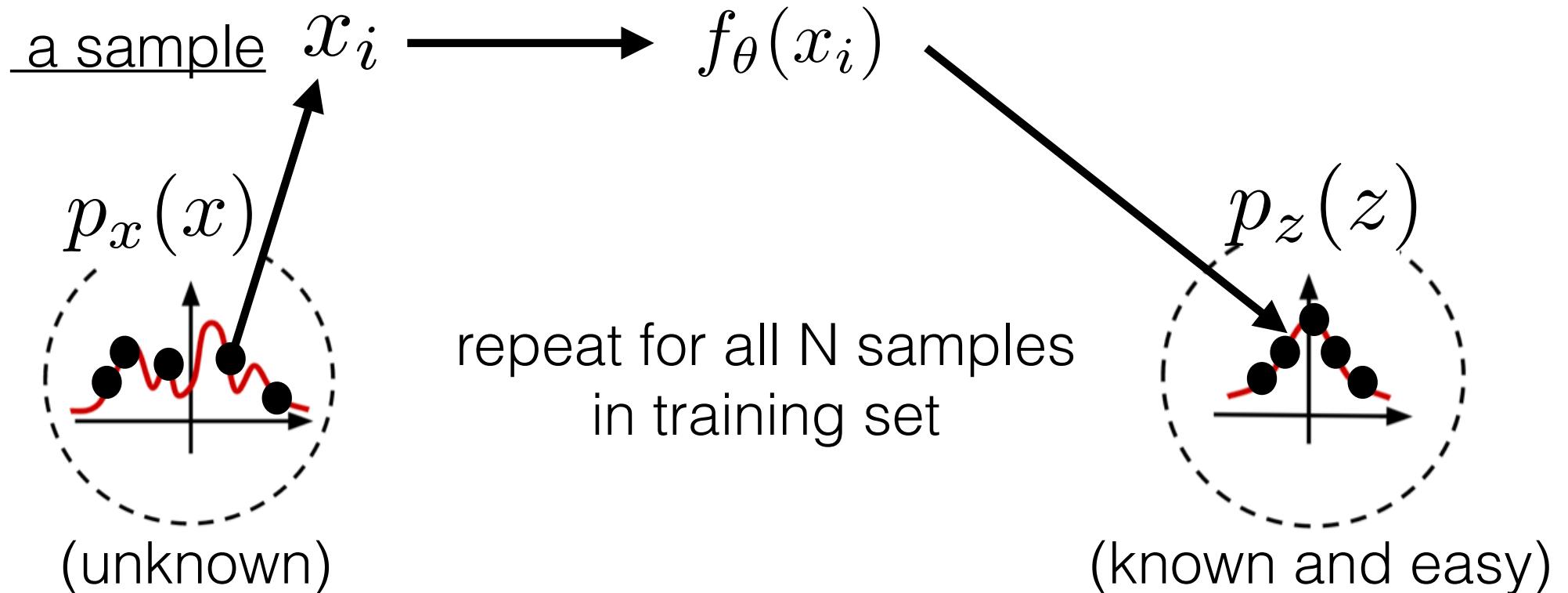
a sample $x_i \longrightarrow f_\theta(x_i)$



training procedure



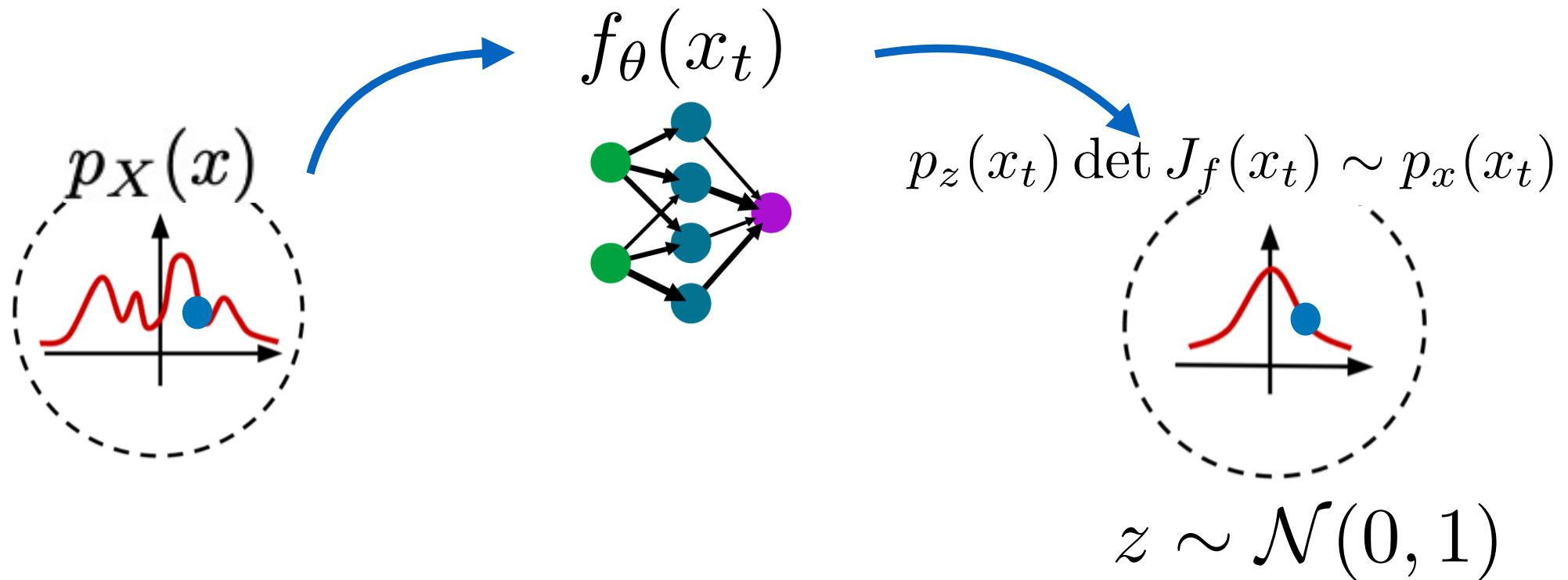
training procedure



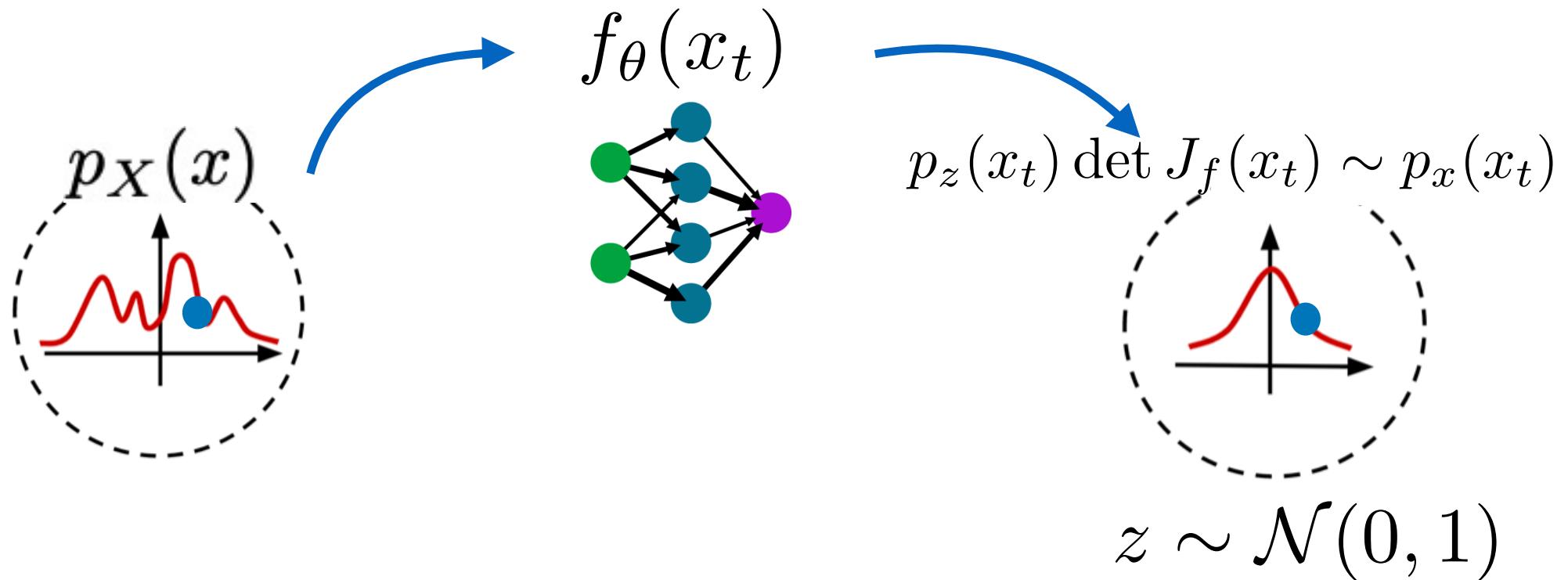
achieved by maximizing this loss function

$$\sum_{i=1}^N \log p_z(x_i) + \log \det J_f(x_i)$$

Likelihood evaluation

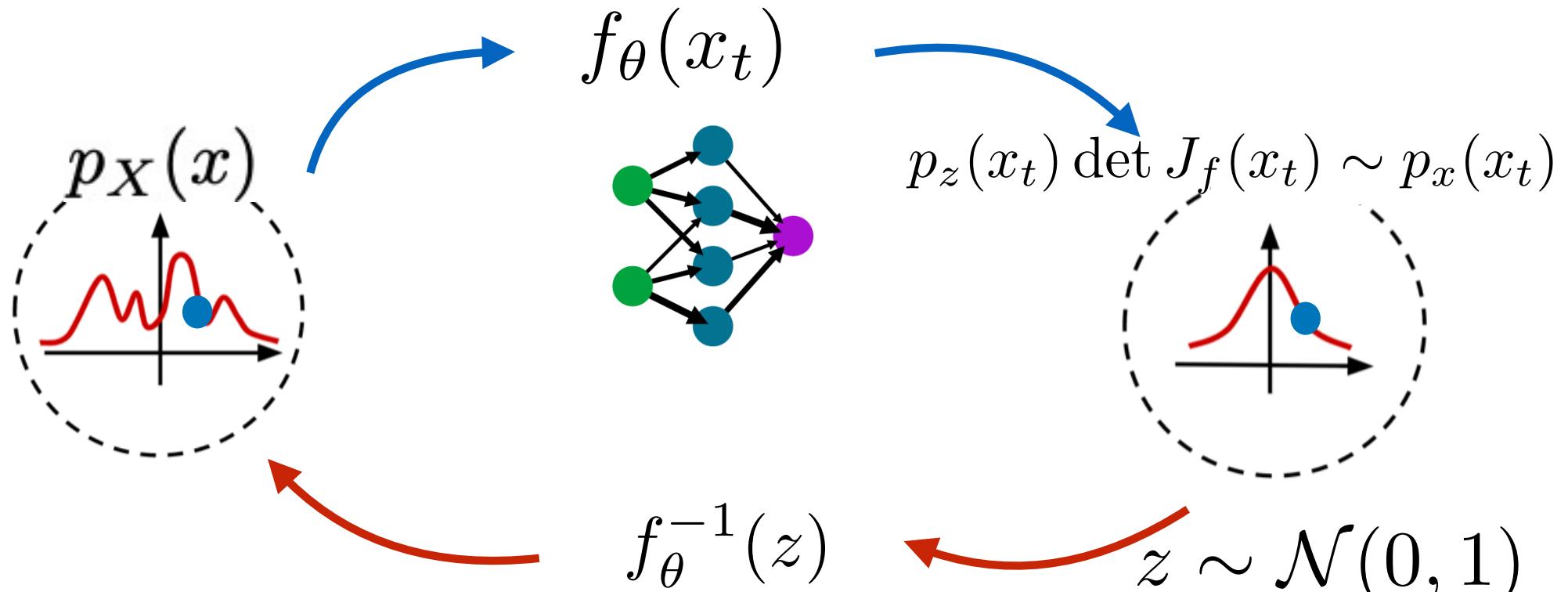


Likelihood evaluation

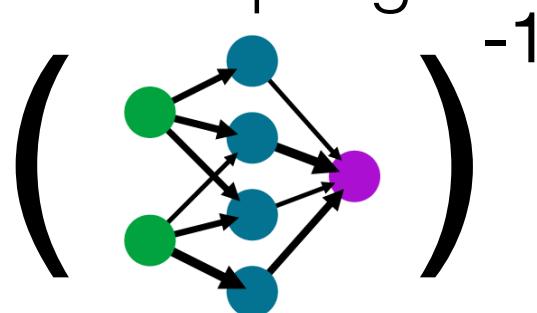


Sampling?

Likelihood evaluation



Sampling



The price to pay is that we need to put some restrictions on the NN architecture

- 1. Easily invertible - to enable sampling
- 2. Jacobian tractable - to optimize the network

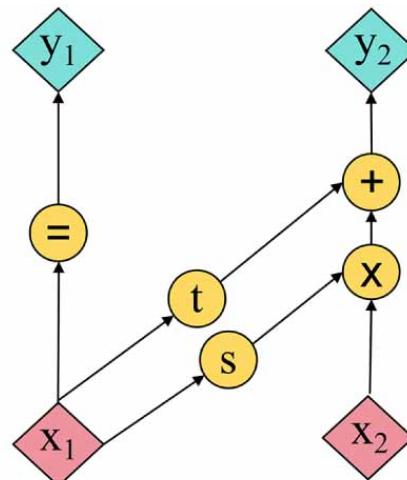
What functions satisfy these conditions?

An example are: RealNVPs (Real-valued Non-Volume Preserving)

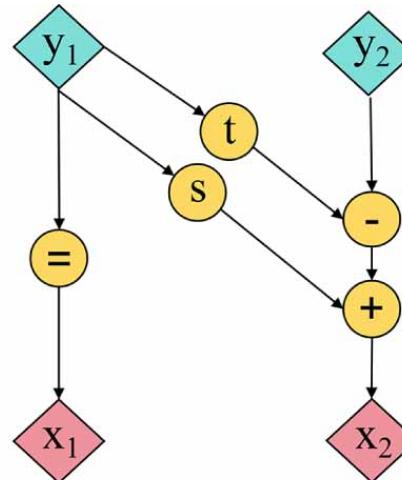
$$\mathbf{y}_{1:d} = \mathbf{x}_{1:d}$$

$$\mathbf{y}_{d+1:D} = \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d})$$

NN



(a) Forward transformation



(b) Inverse transformation

- Easily invertible:

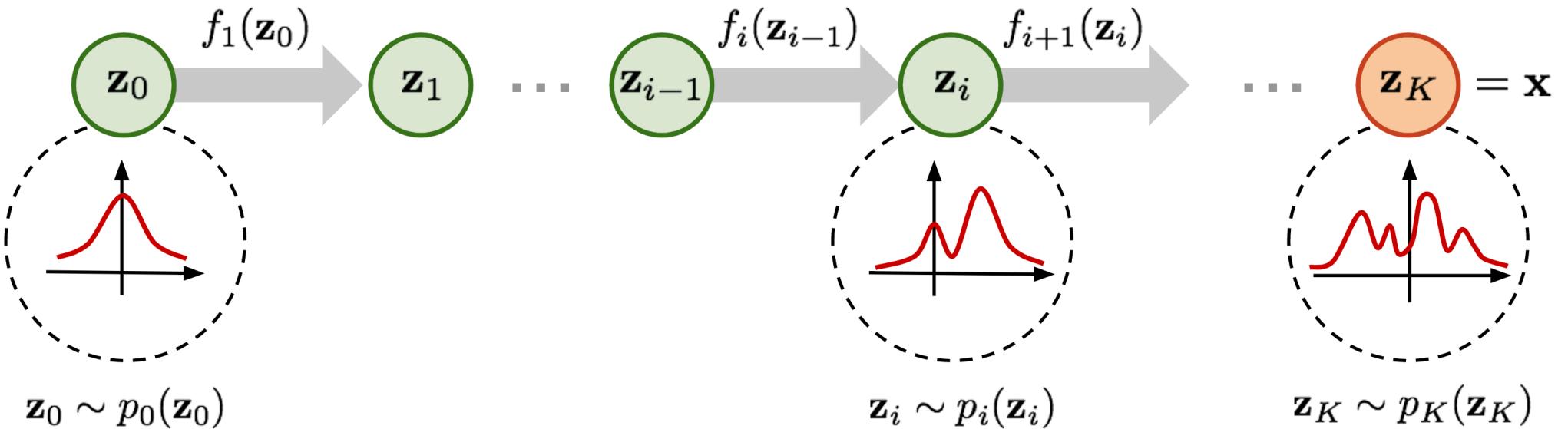
$$\begin{cases} \mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} &= \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} &= (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot \exp(-s(\mathbf{y}_{1:d})) \end{cases}$$

- Jacobian:

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix}$$

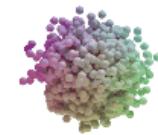
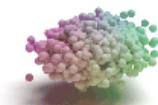
$$\det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp(s(\mathbf{x}_{1:d}))_j = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_{1:d})_j\right)$$

Normalizing flow: in practice we use a concatenation of neural networks (called bijectors)

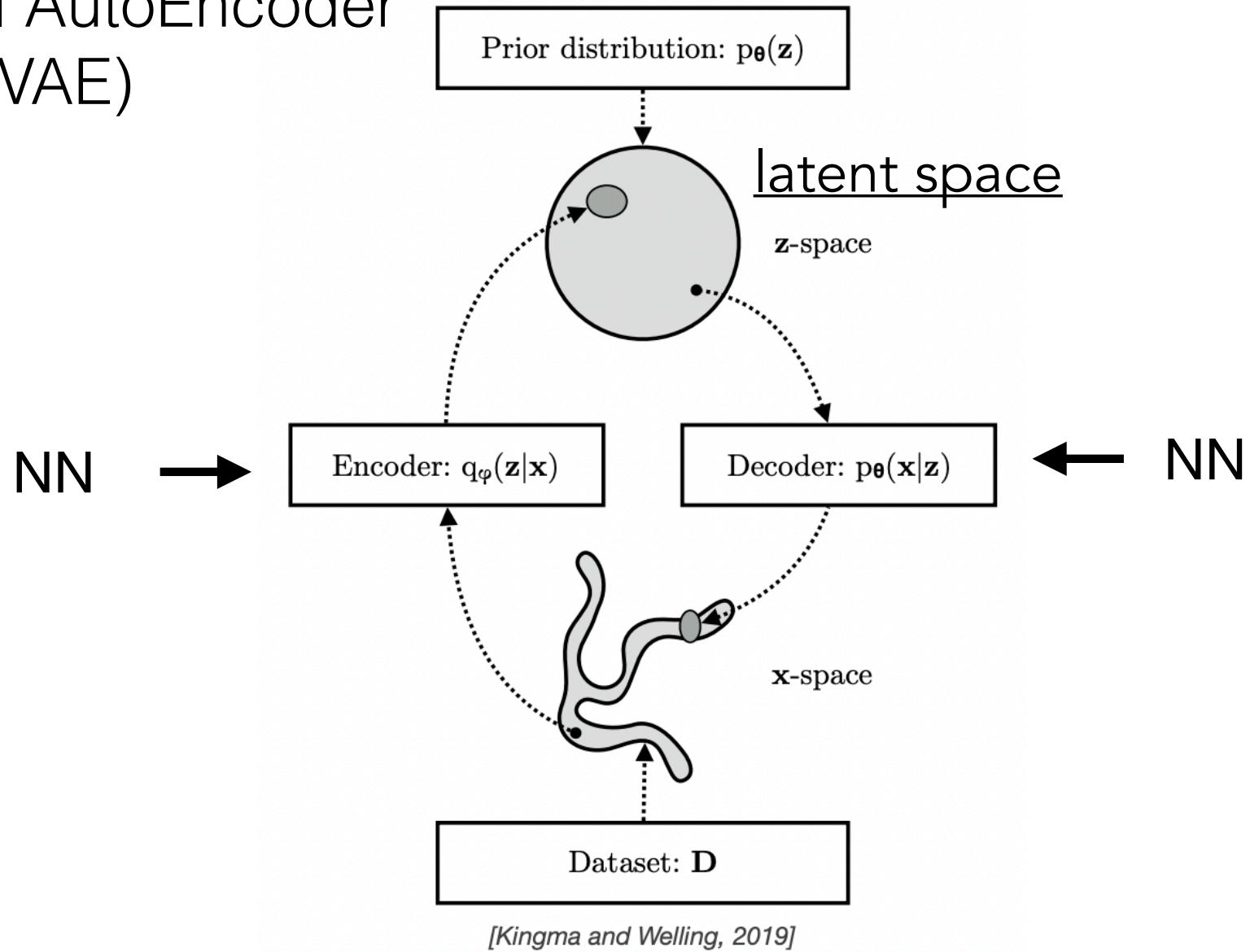


$z_i = f_i(z_{i-1})$ are **invertible** and **differentiable** transformations

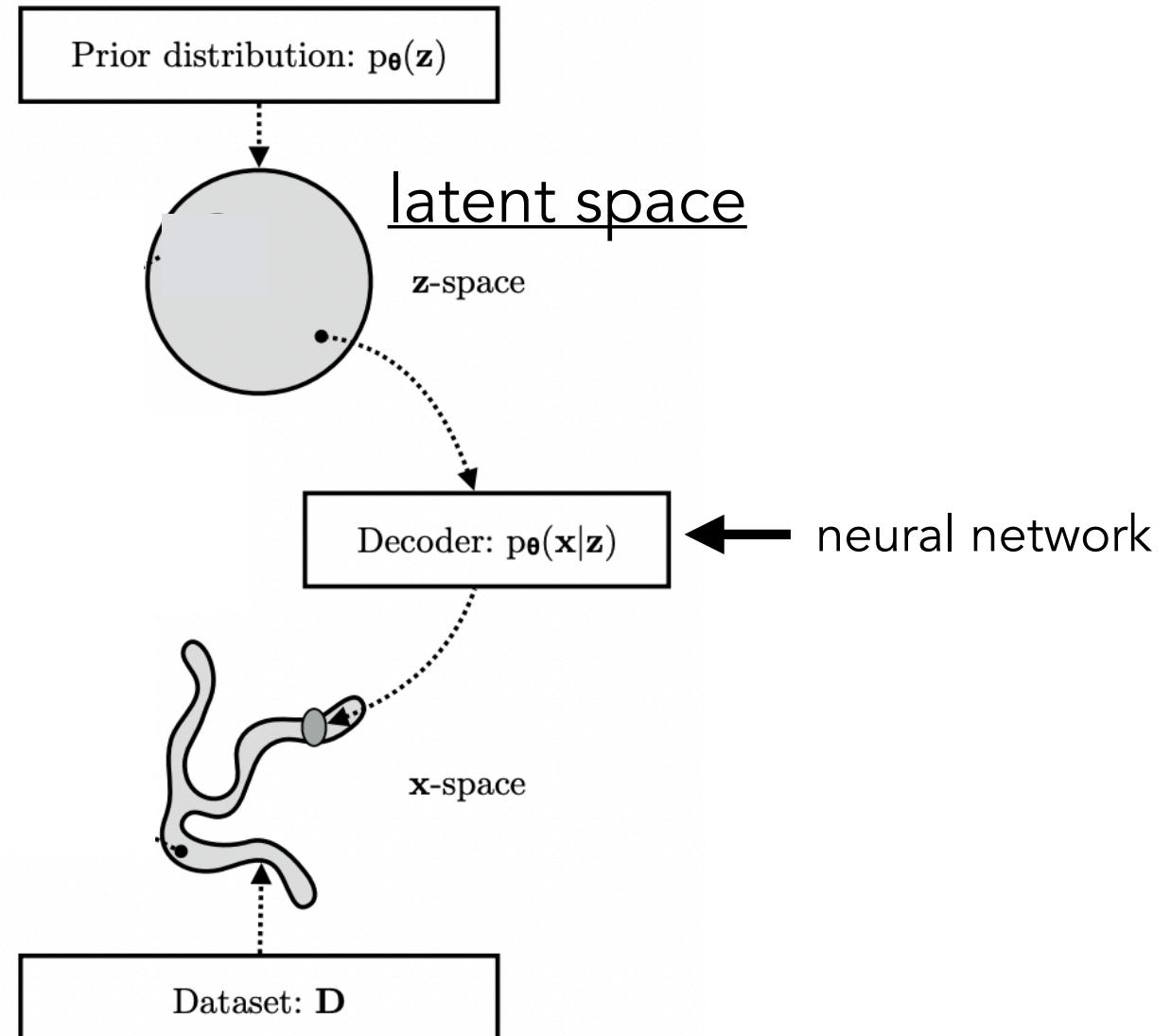
$f = f_1 \circ f_2 \dots \circ f_{k-1} \circ f_k$ is also invertible and differentiable



Variational AutoEncoder (VAE)

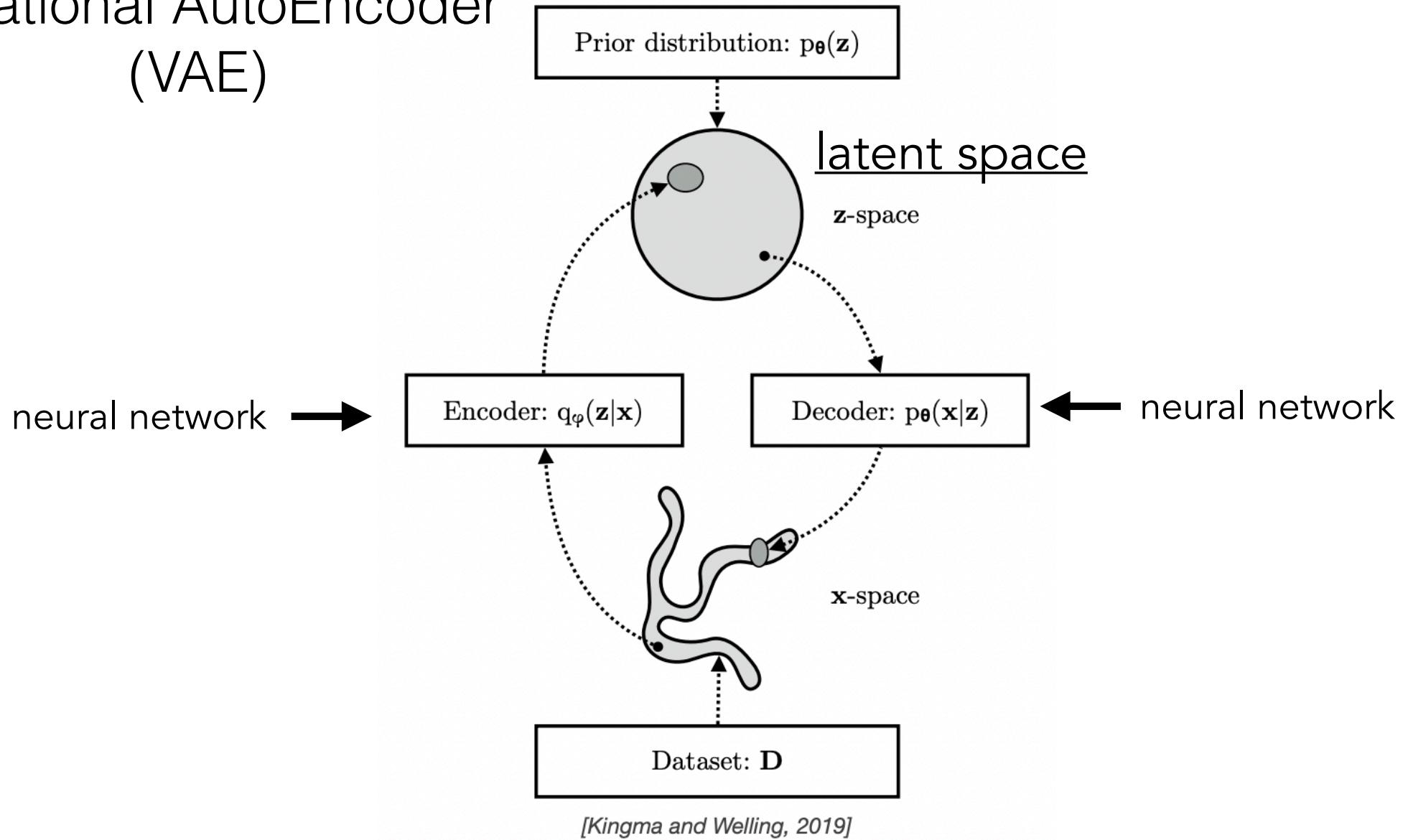


Variational AutoEncoder (VAE)

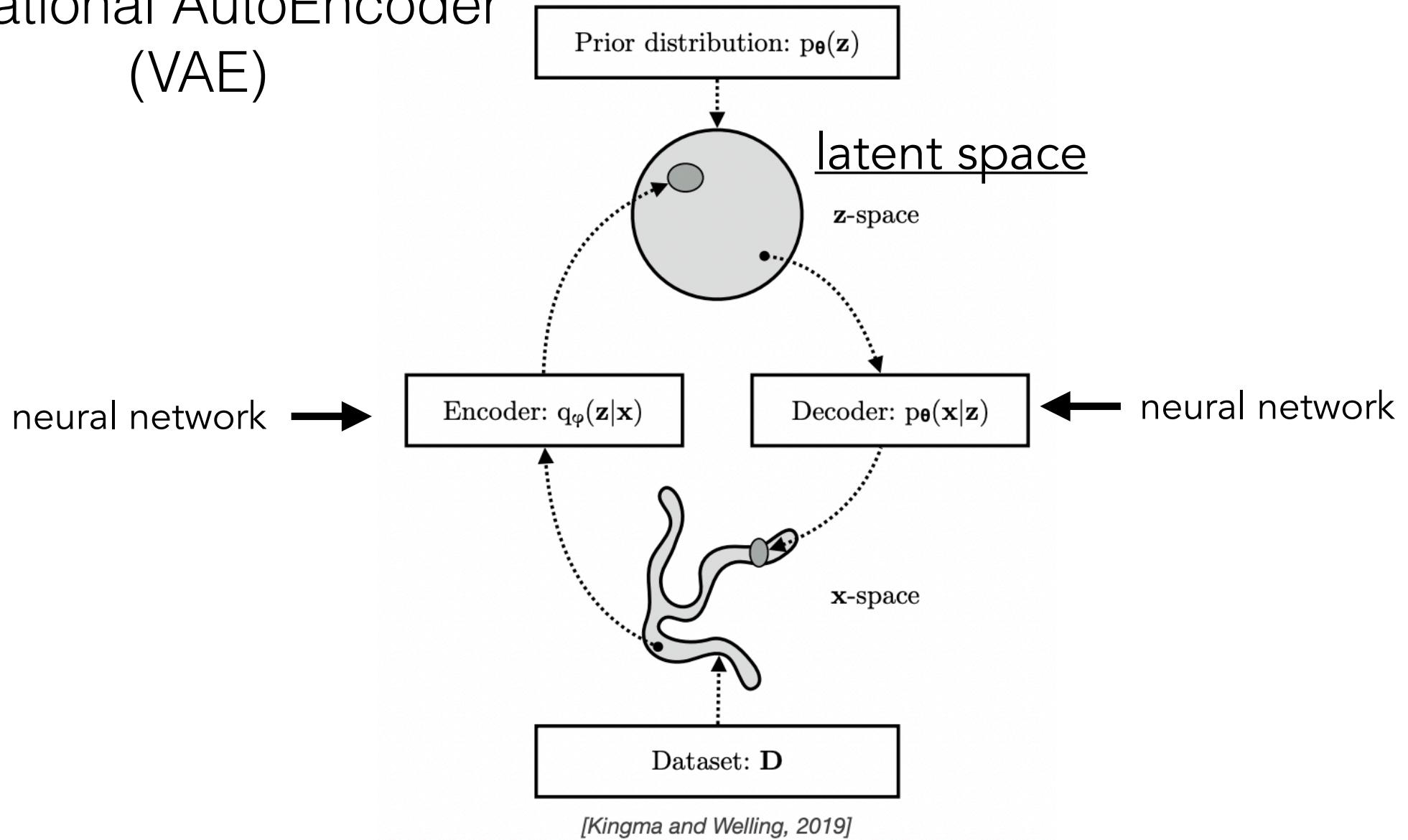


[Kingma and Welling, 2019]

Variational AutoEncoder (VAE)



Variational AutoEncoder (VAE)



Reconstruction term

Regularization term

$$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_\phi(z|x)} [\log p_\theta(x | z)] - \text{KL}(q_\phi(z | x) \| p_\theta(z))$$

we want to maximize the likelihood of our data:

$$\log p_{\theta}(x) = \log \int p_{\theta}(x, z) dz = \log \int p_{\theta}(x | z) p_{\theta}(z) dz$$

(Bayes Rule)

$$\log p_{\theta}(x) = \log \int q_{\phi}(z | x) \frac{p_{\theta}(x, z)}{q_{\phi}(z | x)} dz = \log \mathbb{E}_{q_{\phi}(z|x)} \left[\frac{p_{\theta}(x, z)}{q_{\phi}(z | x)} \right]$$

(Trick to introduce q)

$$\log p_{\theta}(x) \geq \mathbb{E}_{q_{\phi}(z|x)} [\log p_{\theta}(x, z) - \log q_{\phi}(z | x)] \equiv \mathcal{L}(\theta, \phi; x)$$

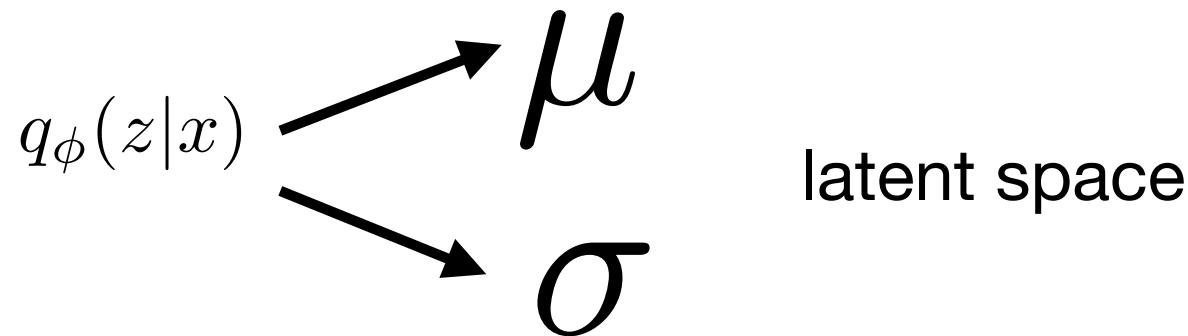
(Jensen's inequality) ELBO

Reconstruction term	Regularization term
$\mathcal{L}(\theta, \phi; x) = \mathbb{E}_{q_{\phi}(z x)} [\log p_{\theta}(x z)] - \text{KL}(q_{\phi}(z x) \ p_{\theta}(z))$	

VAE loss

regularization term (multivariate gaussian and gaussian prior)

$$q_\phi(z | x) = \mathcal{N}(\mu_\phi(x), \text{diag}(\sigma_\phi^2(x))), \quad p_\theta(z) = \mathcal{N}(0, I) \Rightarrow \text{KL} = \frac{1}{2} \sum_{j=1}^d (\mu_j^2 + \sigma_j^2 - \log \sigma_j^2 - 1)$$



reconstruction term (multivariate gaussian and gaussian prior)

$$p_\theta(x | z) = \mathcal{N}(f_\theta(z), \sigma^2 I) \Rightarrow \mathbb{E}_{q_\phi}[\log p_\theta(x | z)] = -\frac{1}{2\sigma^2} \mathbb{E}_{q_\phi}[\|x - f_\theta(z)\|^2] - \frac{D}{2} \log(2\pi\sigma^2)$$

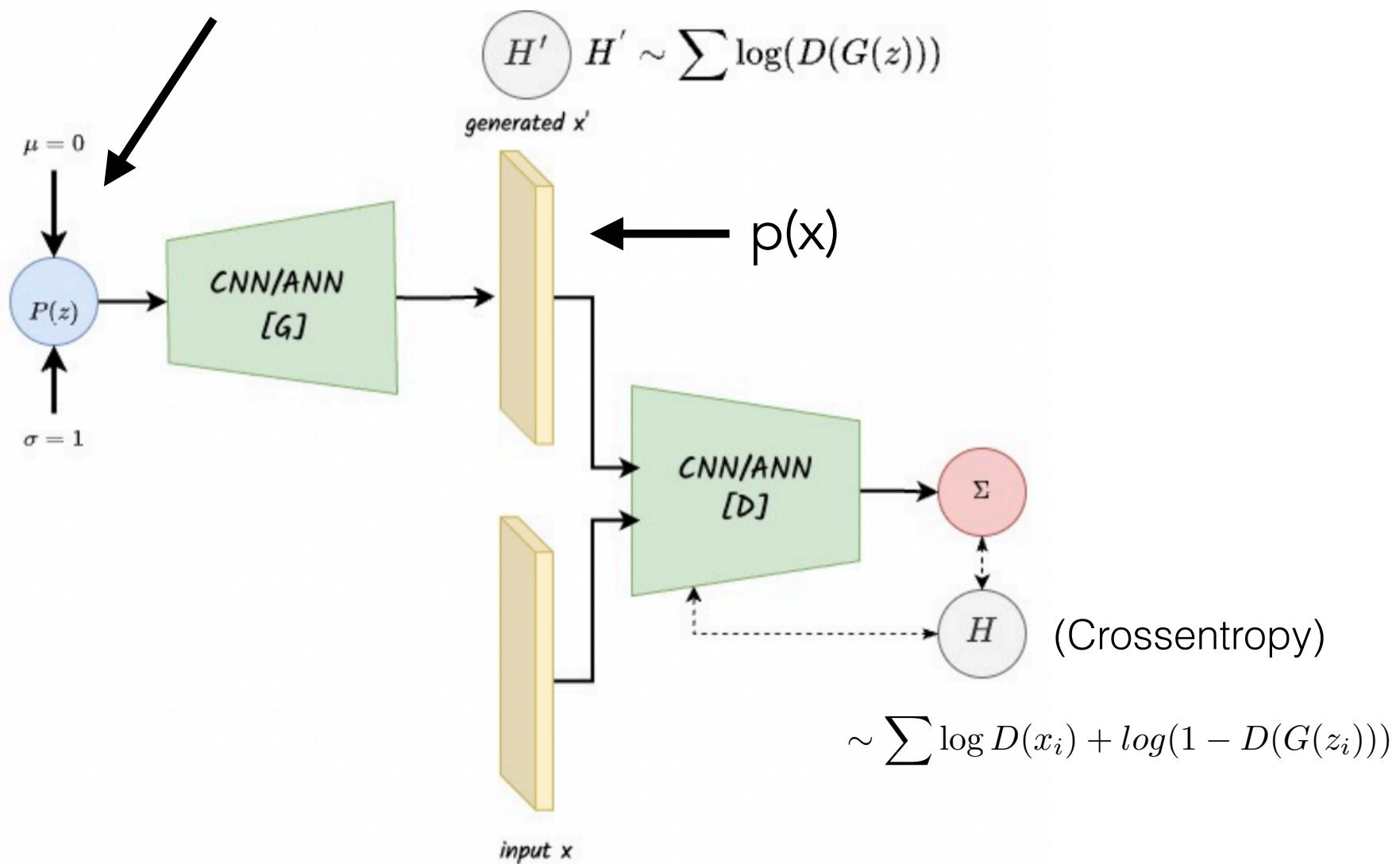
data space

Generative Adversarial Networks

GANs also transform a probability distribution $p(z)$ into $p(x)$

They convert the problem into a “supervised approach” by using two competing neural networks

The latent variable is
here



Generative Adversarial Networks learn the mapping between $p(z)$ and $p(x)$ using an adversarial approach.

discriminator loss

$$\mathcal{L}_D = -\mathbb{E}_{x \sim p_{\text{data}}} [\log D(x)] - \mathbb{E}_{z \sim p(z)} [\log (1 - D(G(z)))].$$

D maximizes the log-likelihood of the correct label (real/fake).

generator loss

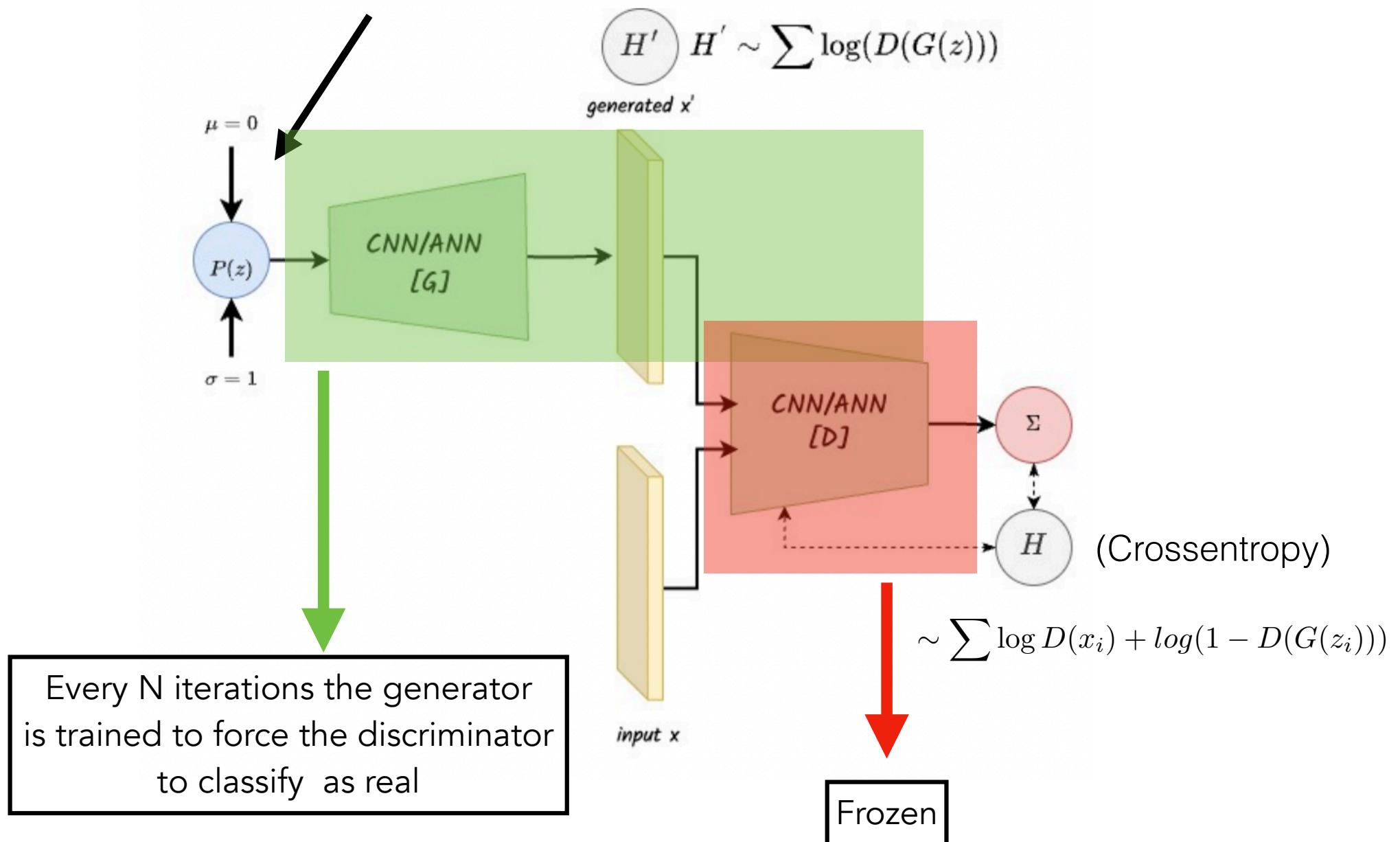
$$\mathcal{L}_G^{\text{NS}} = -\mathbb{E}_{z \sim p(z)} [\log D(G(z))]$$

G tries to make fakes that D labels as real.

$$\mathcal{L}_D = -\frac{1}{m} \sum_{i=1}^m [\log D(x_i) + \log (1 - D(\tilde{x}_i))], \quad \mathcal{L}_G^{\text{NS}} = -\frac{1}{m} \sum_{i=1}^m \log D(\tilde{x}_i)$$

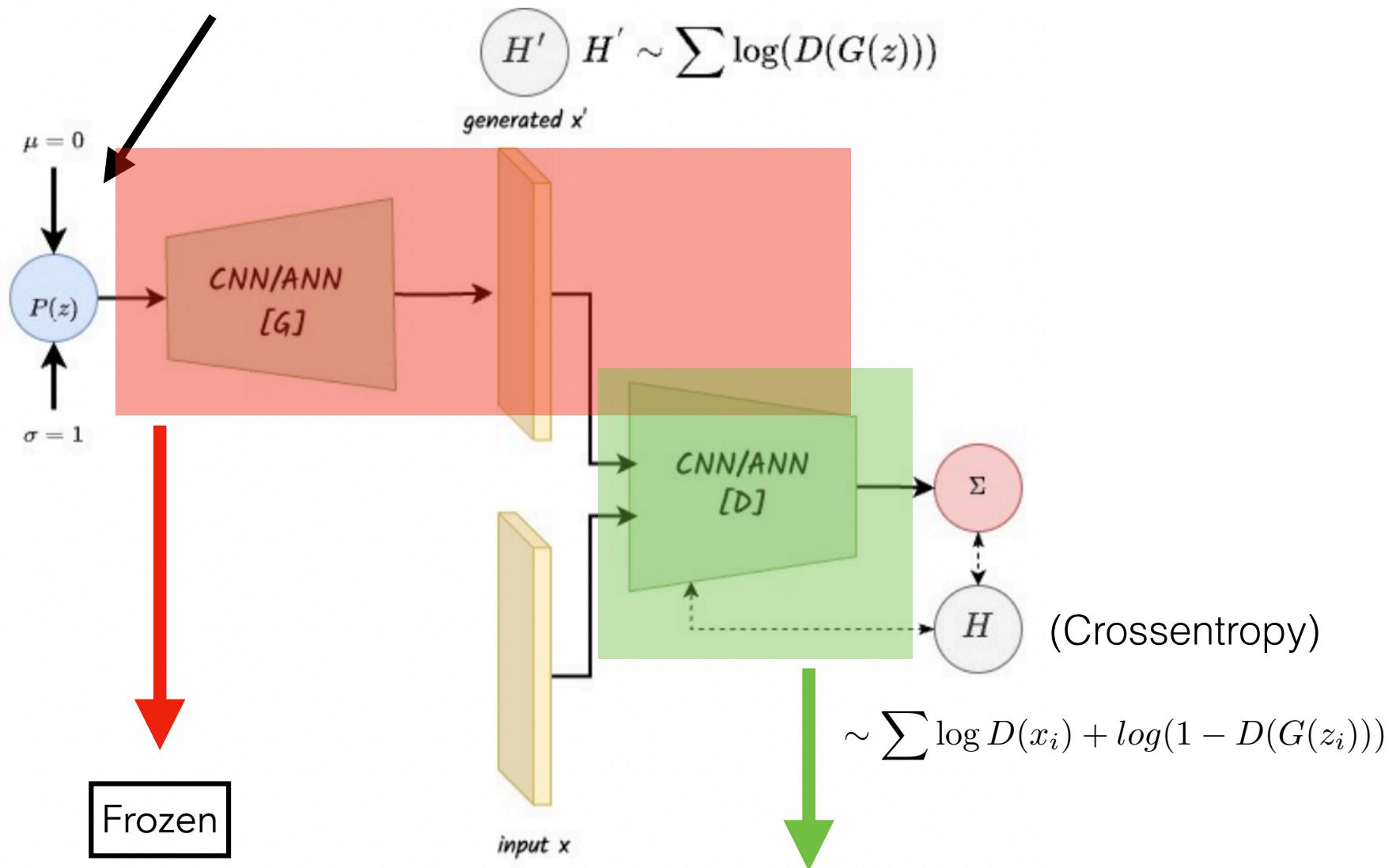
adversarial training

The latent variable is
here



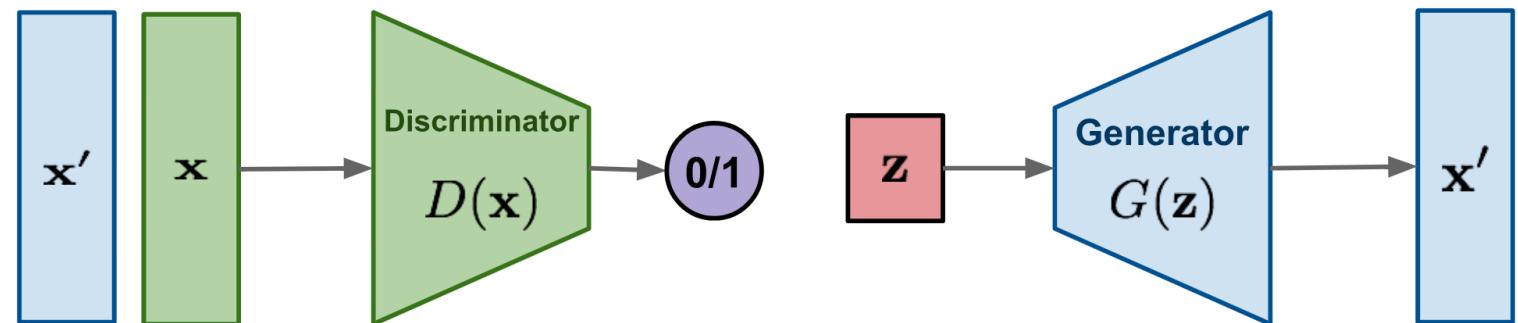
adversarial training

The latent variable is
here

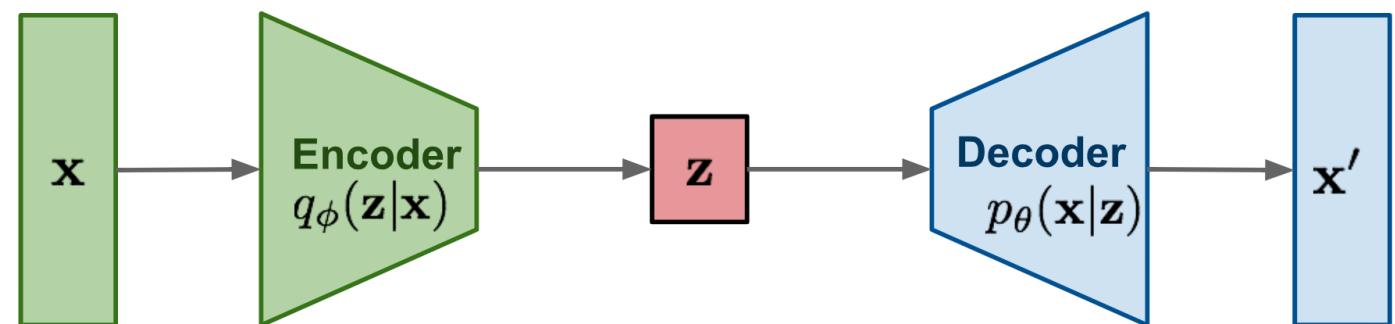


Every N iterations the discriminator
is trained to force to distinguish between
real and fake

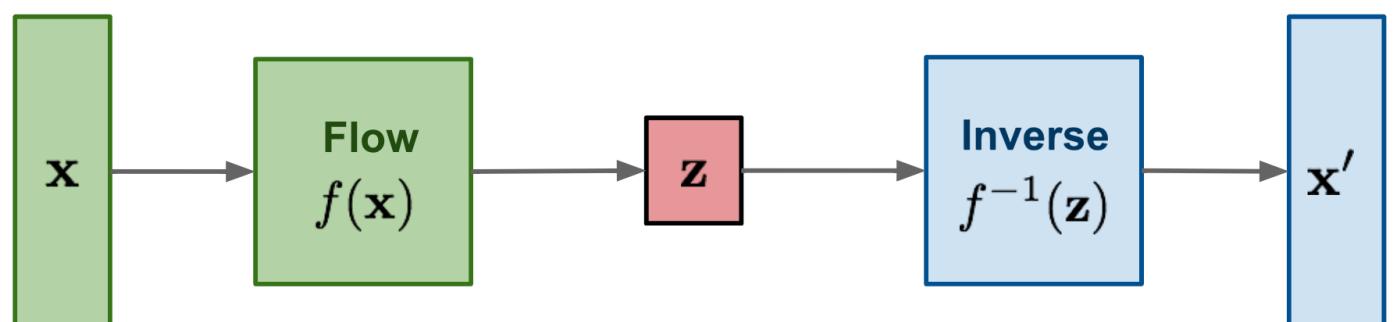
GAN: minimax the classification error loss.



VAE: maximize ELBO.



Flow-based generative models: minimize the negative log-likelihood



Likelihood-based models estimate either a lower bound (~ELBO) or require restrictions in the NN architectures (~Flows). GANs kind of work around these limitations but adversarial training is unstable.

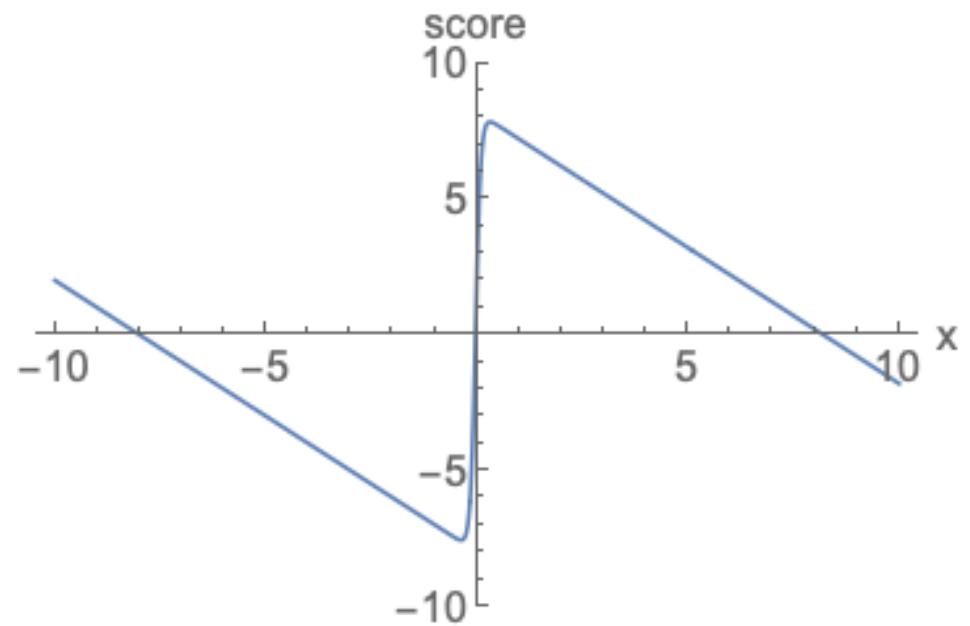
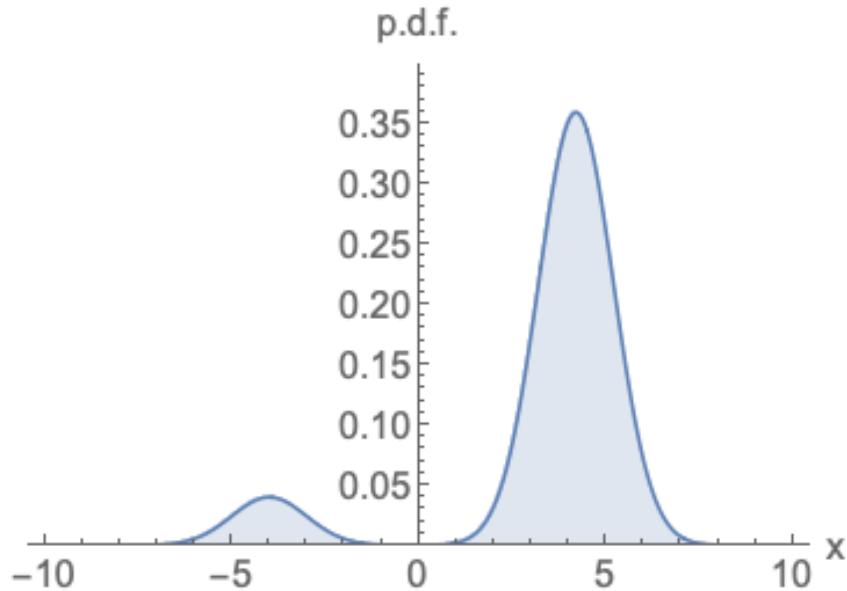
Likelihood-based models estimate either a lower bound (~ELBO) or require restrictions in the NN architectures (~Flows). GANs kind of work around these limitations but adversarial training is unstable.

$$p(x) \longrightarrow \nabla_x \log p(x)$$

$$S_\theta$$

: Score network, i.e. Neural Network that approximates the gradient of p with respect to the data

why is the score interesting?



Independent of normalization as compared to p

$$p_{\theta}(x) = \frac{\exp(f_{\theta}(x))}{Z(\theta)},$$

$$\nabla_x \log p_{\theta}(x) = \nabla_x f_{\theta}(x) - \nabla_x Z(\theta)$$

$$\nabla_x \log Z(\theta) = 0$$

One key ingredient is Langevin Dynamics which allows one to sample from a probability distribution using only the score

start with random numbers:

$$z_t \sim \mathcal{N}(0, 1) \quad \tilde{x}_0 \sim \pi(x) \quad (\text{arbitrary distribution})$$

run iteratively for $t < T$:

$$\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{\epsilon}{2} \nabla_{\mathbf{x}} \log p(\tilde{\mathbf{x}}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t, \quad \epsilon \ll$$

for T large:

\tilde{x}_T is exact sample of $p(x)$

One key ingredient is Langevin Dynamics which allows one to sample from a probability distribution using only the score

start with random numbers:

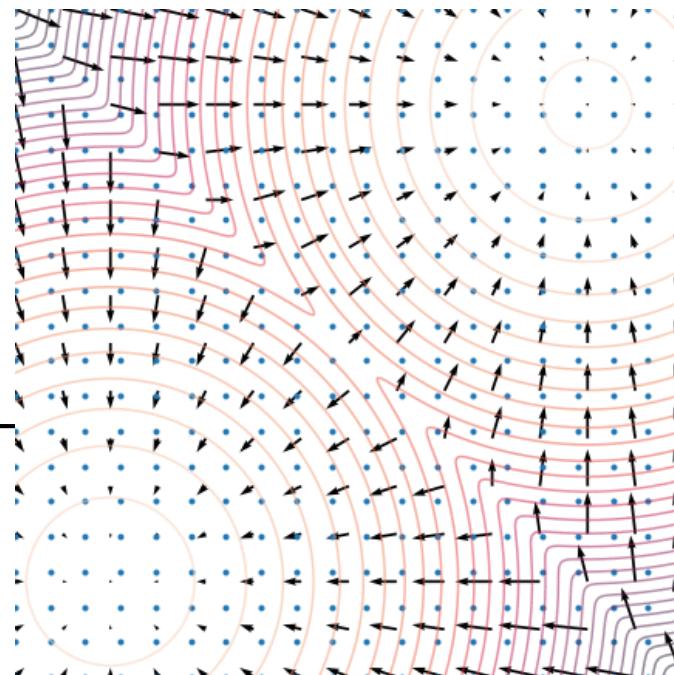
$$z_t \sim \mathcal{N}(0, 1) \quad \tilde{x}_0 \sim \pi(x) \quad (\text{arbitrary distribution})$$

run iteratively for $t < T$:

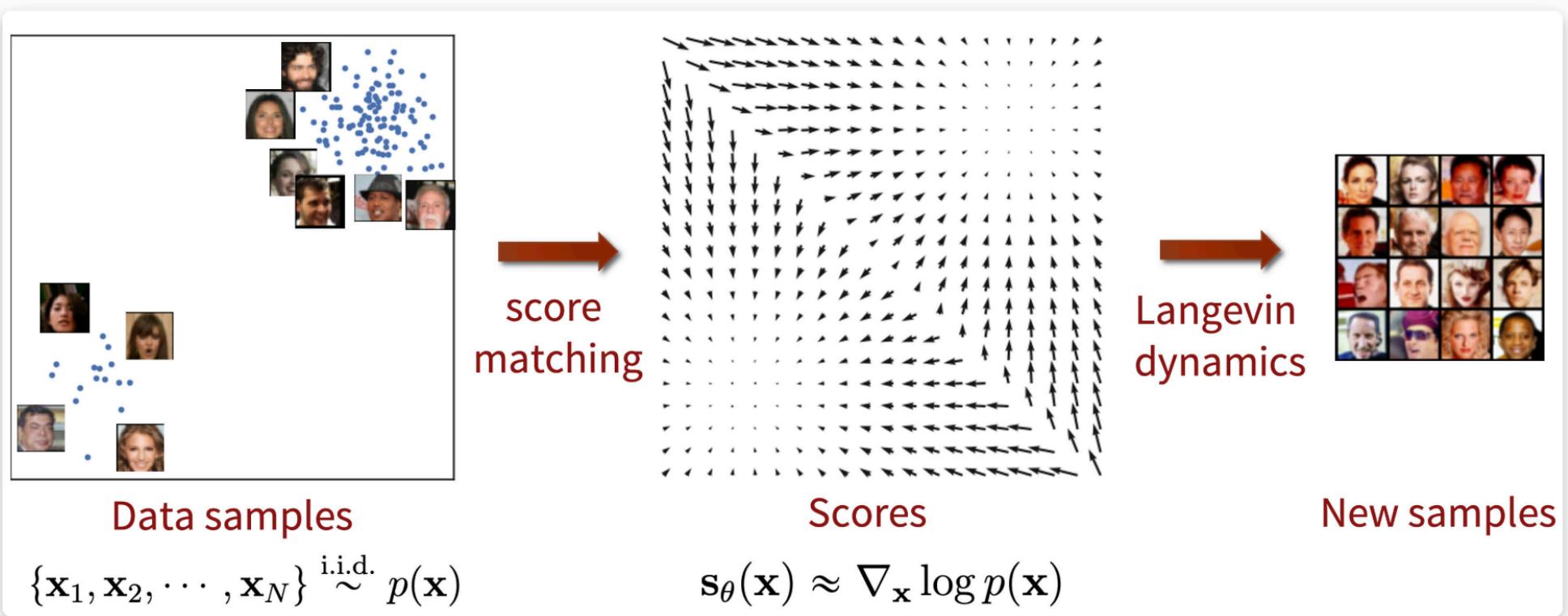
$$\tilde{\mathbf{x}}_t = \tilde{\mathbf{x}}_{t-1} + \frac{\epsilon}{2} \nabla_{\mathbf{x}} \log p(\tilde{\mathbf{x}}_{t-1}) + \sqrt{\epsilon} \mathbf{z}_t, \quad \epsilon \ll$$

for T large:

\tilde{x}_T is exact sample of $p(x)$



score based diffusion models



<https://yang-song.net/blog/2021/score/>

the neural network

We need to estimate the score (with an optimization problem): denoising score matching (DSM)

$$\tilde{x} \mid x \sim \mathcal{N}(x, \sigma^2 I), \quad \tilde{x} = x + \sigma \varepsilon, \quad \varepsilon \sim \mathcal{N}(0, I), q_\sigma(\tilde{x}) = \int p_{\text{data}}(x) \mathcal{N}(\tilde{x}; x, \sigma^2 I) dx.$$

$$\frac{1}{2} \mathbb{E}_{q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})} p_{\text{data}}(\mathbf{x}) \left[\|\mathbf{s}_\theta(\tilde{\mathbf{x}}) - \nabla_{\tilde{\mathbf{x}}} \log q_\sigma(\tilde{\mathbf{x}}|\mathbf{x})\|_2^2 \right].$$

$$s_\theta^*(x) = \nabla_x \log q_\sigma(x) \simeq \nabla_x \log p_{\text{data}}(x)$$

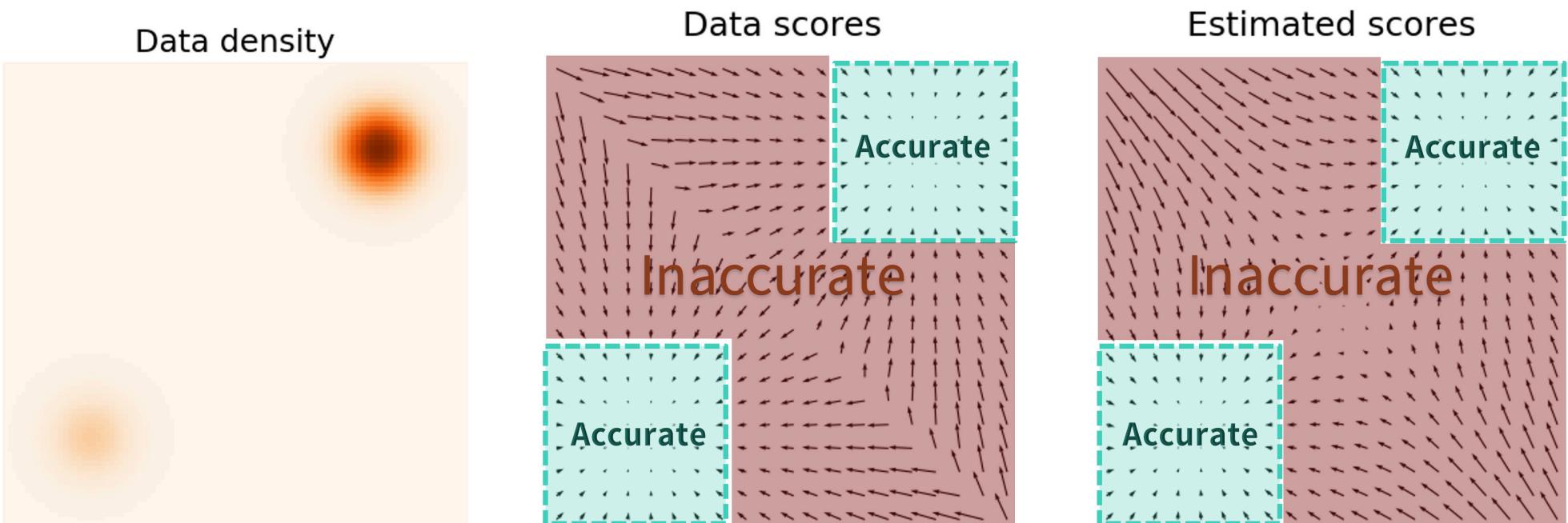
if noise is small

$$\mathcal{J}_{\text{DSM}}(\theta; \sigma) = \mathbb{E}_{x \sim p_{\text{data}}} \mathbb{E}_{\tilde{x} \sim \mathcal{N}(x, \sigma^2 I)} \left[\frac{1}{2} \left\| s_\theta(\tilde{x}, \sigma) + \frac{\tilde{x} - x}{\sigma^2} \right\|_2^2 \right].$$

$$\mathcal{J}_{\text{DSM}}(\theta; \sigma) = \mathbb{E}_{x \sim p_{\text{data}}} \mathbb{E}_{\varepsilon \sim \mathcal{N}(0, I)} \left[\frac{1}{2} \left\| s_\theta(x + \sigma \varepsilon, \sigma) + \frac{\varepsilon}{\sigma} \right\|_2^2 \right]$$

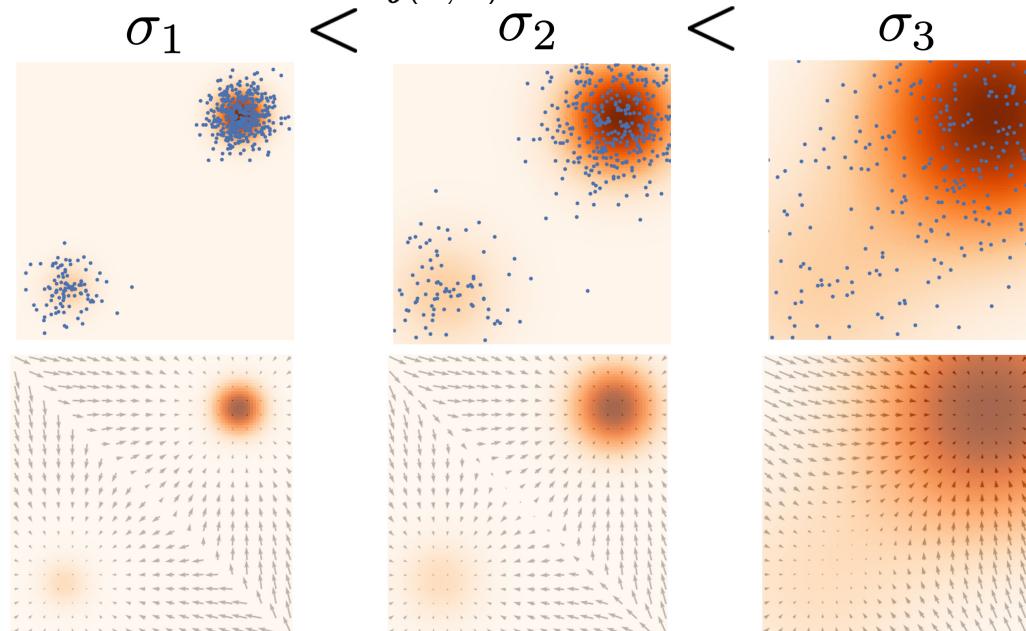
Estimating the score is difficult

In regions of low data density, score matching may not have enough evidence to estimate score functions accurately



noise conditional score network

Let $\{\sigma_i\}_{i=1}^L$ be a positive geometric sequence that satisfies $\frac{\sigma_1}{\sigma_2} = \dots = \frac{\sigma_{L-1}}{\sigma_L} > 1$. Let $q_\sigma(\mathbf{x}) \triangleq \int p_{\text{data}}(\mathbf{t}) \mathcal{N}(\mathbf{x} \mid \mathbf{t}, \sigma^2 I) d\mathbf{t}$ denote the perturbed data distribution. We choose the noise levels $\{\sigma_i\}_{i=1}^L$ such that σ_1 is large enough to mitigate the difficulties discussed in Section 3, and σ_L is small enough to minimize the effect on data. We aim to train a conditional score network to jointly estimate the scores of all perturbed data distributions, *i.e.*, $\forall \sigma \in \{\sigma_i\}_{i=1}^L : \mathbf{s}_\theta(\mathbf{x}, \sigma) \approx \nabla_{\mathbf{x}} \log q_\sigma(\mathbf{x})$. Note that $\mathbf{s}_\theta(\mathbf{x}, \sigma) \in \mathbb{R}^D$ when $\mathbf{x} \in \mathbb{R}^D$. We call $\mathbf{s}_\theta(\mathbf{x}, \sigma)$ a *Noise Conditional Score Network (NCSN)*.



The loss function for a given sigma is:

$$\ell(\boldsymbol{\theta}; \sigma) \triangleq \frac{1}{2} \mathbb{E}_{p_{\text{data}}(\mathbf{x})} \mathbb{E}_{\tilde{\mathbf{x}} \sim \mathcal{N}(\mathbf{x}, \sigma^2 I)} \left[\left\| \mathbf{s}_\theta(\tilde{\mathbf{x}}, \sigma) + \frac{\tilde{\mathbf{x}} - \mathbf{x}}{\sigma^2} \right\|_2^2 \right].$$

And then combined:

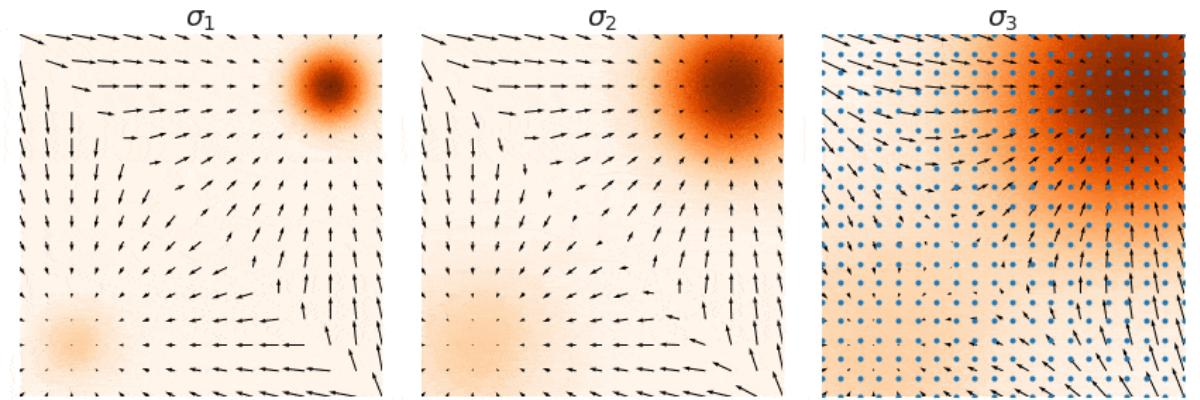
$$\mathcal{L}(\boldsymbol{\theta}; \{\sigma_i\}_{i=1}^L) \triangleq \frac{1}{L} \sum_{i=1}^L \lambda(\sigma_i) \ell(\boldsymbol{\theta}; \sigma_i),$$

annealed Langevin dynamics

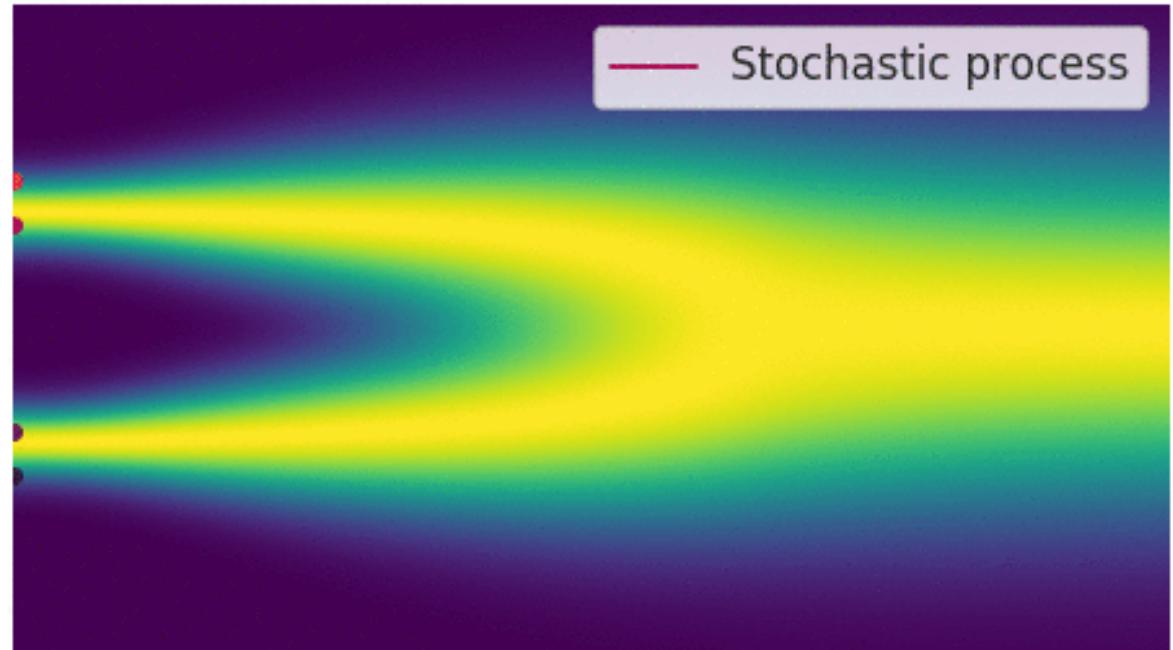
Algorithm 1 Annealed Langevin dynamics.

Require: $\{\sigma_i\}_{i=1}^L, \epsilon, T$.

- 1: Initialize $\tilde{\mathbf{x}}_0$
- 2: **for** $i \leftarrow 1$ to L **do**
- 3: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ ▷ α_i is the step size.
- 4: **for** $t \leftarrow 1$ to T **do**
- 5: Draw $\mathbf{z}_t \sim \mathcal{N}(0, I)$
- 6: $\tilde{\mathbf{x}}_t \leftarrow \tilde{\mathbf{x}}_{t-1} + \frac{\alpha_i}{2} \mathbf{s}_\theta(\tilde{\mathbf{x}}_{t-1}, \sigma_i) + \sqrt{\alpha_i} \mathbf{z}_t$
- 7: **end for**
- 8: $\tilde{\mathbf{x}}_0 \leftarrow \tilde{\mathbf{x}}_T$
- 9: **end for**
- return** $\tilde{\mathbf{x}}_T$



Since the distributions $\{q_{\sigma_i}\}_{i=1}^L$ are all perturbed by Gaussian noise, their supports span the whole space and their scores are well-defined, avoiding difficulties from the manifold hypothesis. When σ_1 is sufficiently large, the low density regions of $q_{\sigma_1}(\mathbf{x})$ become small and the modes become less isolated. As discussed previously, this can make score estimation more accurate, and the mixing of Langevin dynamics faster. We can therefore assume that Langevin dynamics produce good samples for $q_{\sigma_1}(\mathbf{x})$. These samples are likely to come from high density regions of $q_{\sigma_1}(\mathbf{x})$, which means they are also likely to reside in the high density regions of $q_{\sigma_2}(\mathbf{x})$, given that $q_{\sigma_1}(\mathbf{x})$ and $q_{\sigma_2}(\mathbf{x})$ only slightly differ from each other. As score estimation and Langevin dynamics perform better in high density regions, samples from $q_{\sigma_1}(\mathbf{x})$ will serve as good initial samples for Langevin dynamics of $q_{\sigma_2}(\mathbf{x})$. Similarly, $q_{\sigma_{i-1}}(\mathbf{x})$ provides good initial samples for $q_{\sigma_i}(\mathbf{x})$, and finally we obtain samples of good quality from $q_{\sigma_L}(\mathbf{x})$.



$$dx(t) = f(x(t), t) dt + g(t) dW_t, \quad t \in [0, 1], \quad x(0) \sim p_{\text{data}}$$

for example

$$\text{VP (DDPM): } dx = -\frac{1}{2} \beta(t) x dt + \sqrt{\beta(t)} dW_t,$$

- **SDE solution = process:** a continuous family of random variables $\{x(t)\}_{t \in [0,T]}$ tracing trajectories as time goes from $p_t(x) = p_{x(t)}(x)$.
- **Link to discrete noise scales:** continuous time $t \in [0, T] \leftrightarrow$ discrete index $i = 1, 2, \dots, L$; likewise densities correspond as $p_t(x) \leftrightarrow p_{\sigma_i}(x)$
- **Start (no noise):** $p_0(x) = p(x)$ (data distribution)
- **End (max noise / prior):** after long enough time, $p_T(x) \approx \pi(x)$ (tractable prior).
- **Finite-step analogy at the end:** $p_T(x) \leftrightarrow p_{\sigma_L}(x)$ with the largest noise level

$$\sigma_L = \max_{i \in \{1, \dots, L\}} \sigma_i$$

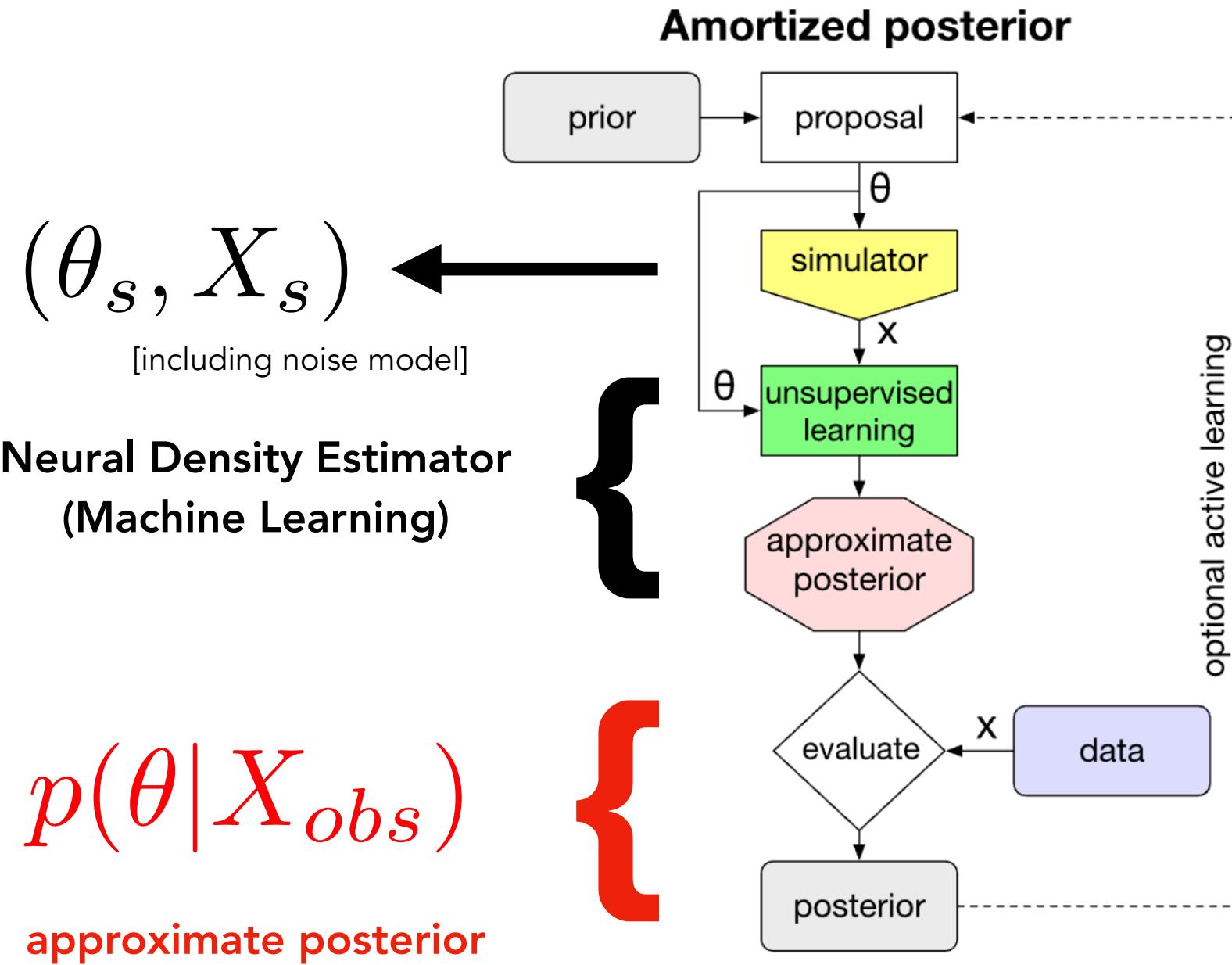
$$s_\theta(x,t) \; \approx \; \nabla_x \log p_t(x), \qquad \min_{\theta} \; \mathbb{E}_{t \sim \mathcal{U}(0,1)} \left[\lambda(t) \, \mathbb{E}_{x_0 \sim p_{\text{data}}, \, \varepsilon \sim \mathcal{N}(0,I)} \Big\| s_\theta(x_t,t) - \nabla_x \log q_t(x_t \mid x_0) \Big\|_2^2 \right].$$

$$\text{VP: } \; x_t = \sqrt{\bar{\alpha}(t)} x_0 + \sqrt{1-\bar{\alpha}(t)} \, \varepsilon \; \Rightarrow \; \nabla_x \log q_t(x_t \mid x_0) = - \frac{\varepsilon}{\sqrt{1-\bar{\alpha}(t)}}.$$

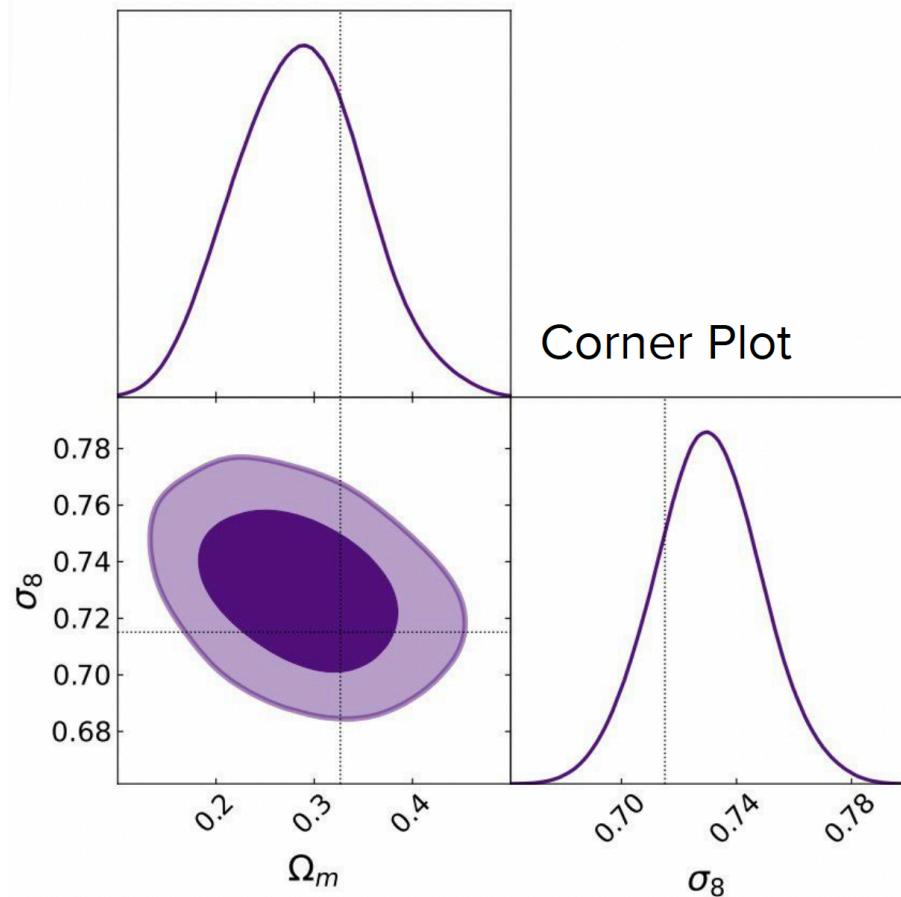
$$x \leftarrow x + \eta \, s_\theta(x,t) + \sqrt{2\eta} \, \xi, \qquad \xi \sim \mathcal{N}(0,I).$$

Back to SBI

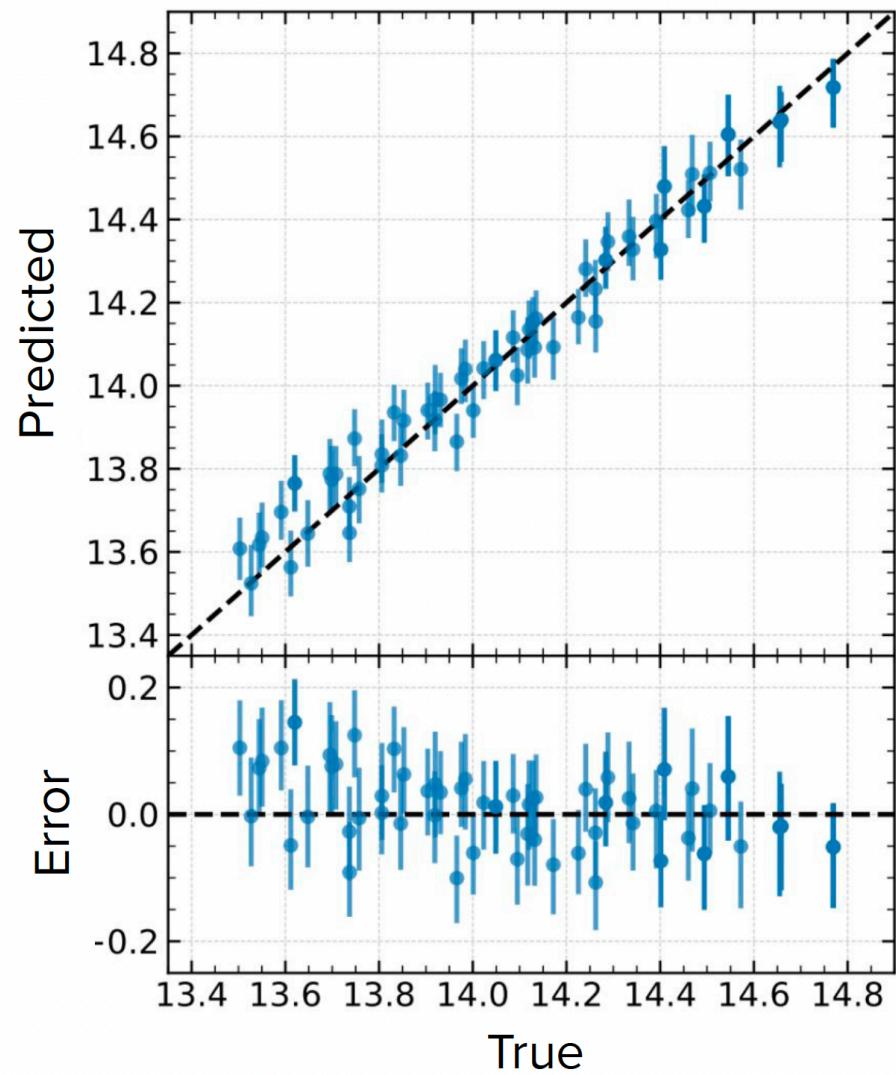
SBI: All you need is a ****good**** forward model



Validation



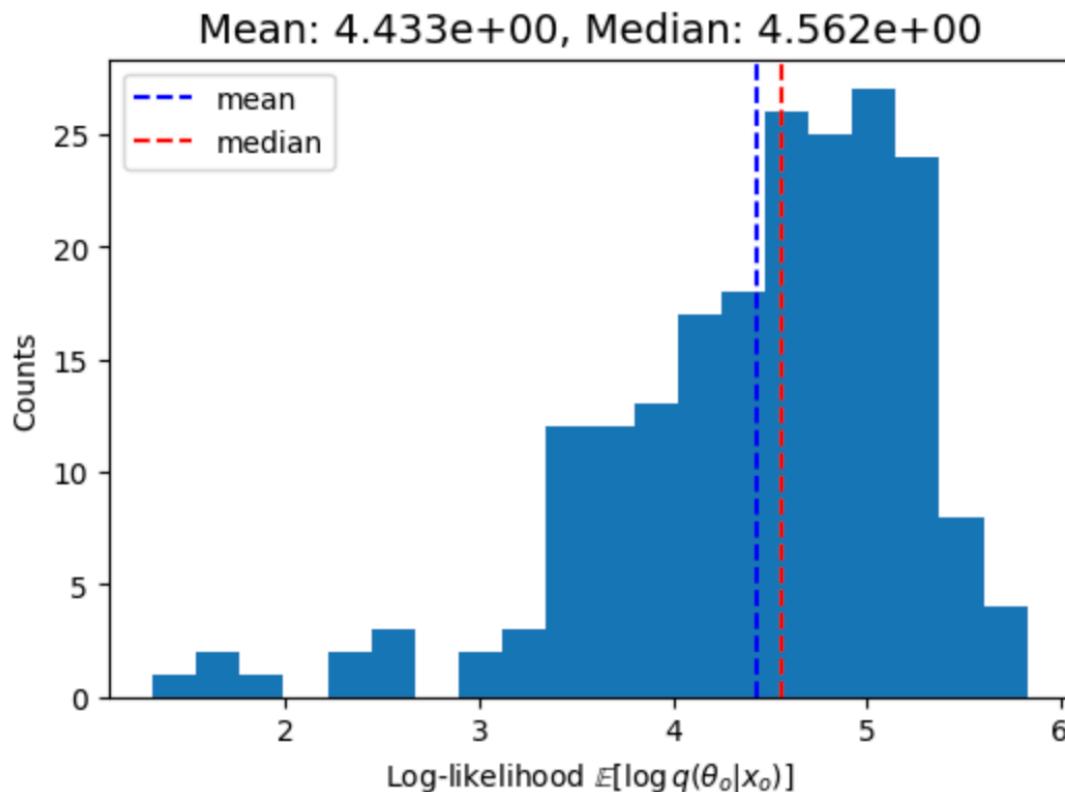
Corner Plot



*adapted from M. Ho (SBI conference)

Cumulative Likelihood of the Test Dataset

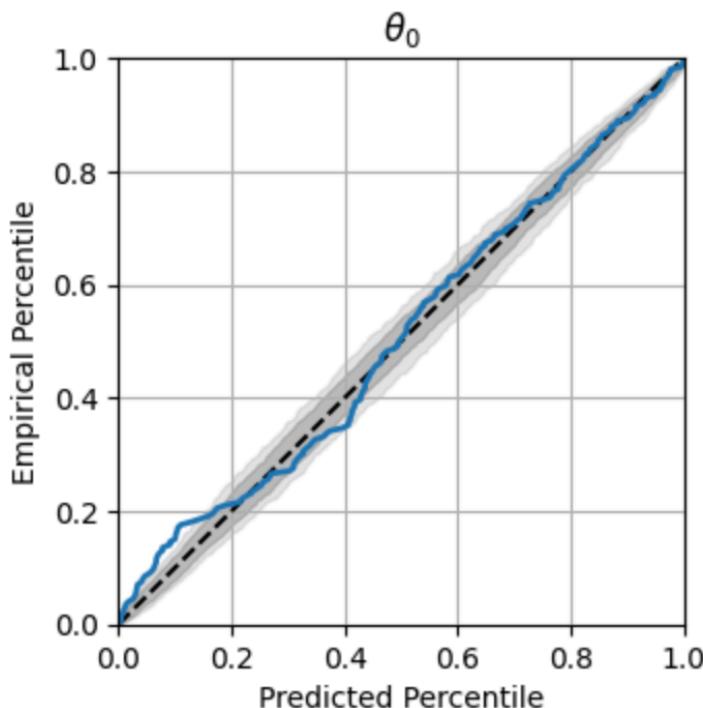
$$N_{test} \prod_{i=1}^{} \hat{P}(\theta_i | x_i)$$



Large values
mean good
predictive power

P-P plots

$$\text{PIT}(\theta; x_o) = \int_{-\infty}^{\theta} d\theta \hat{P}(\theta | x_o).$$



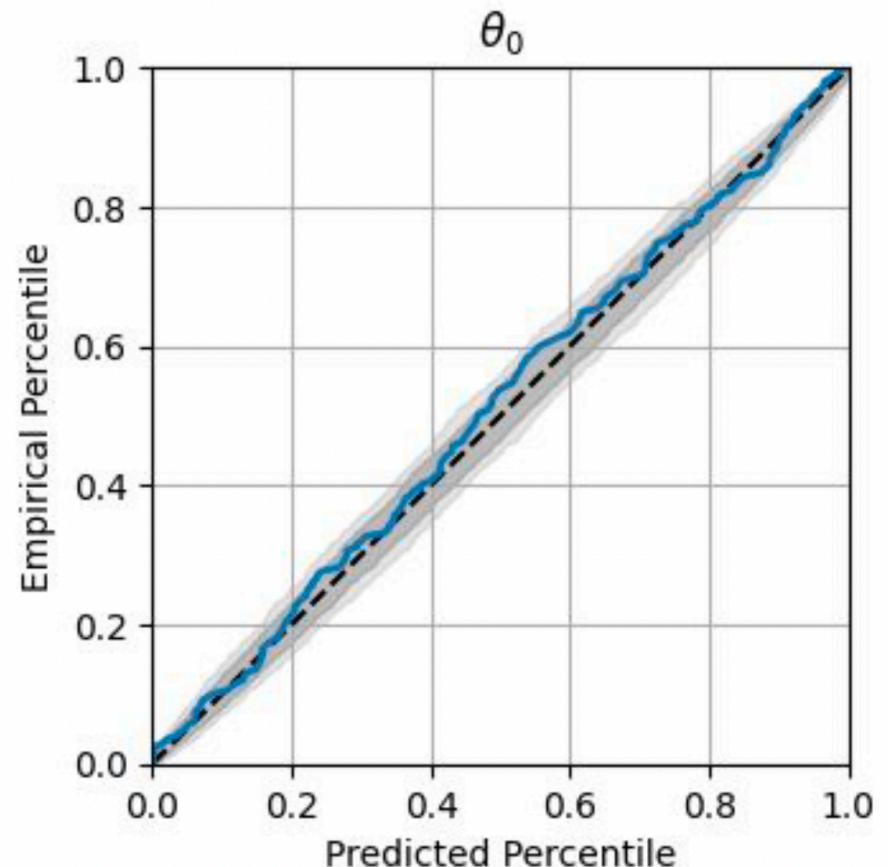
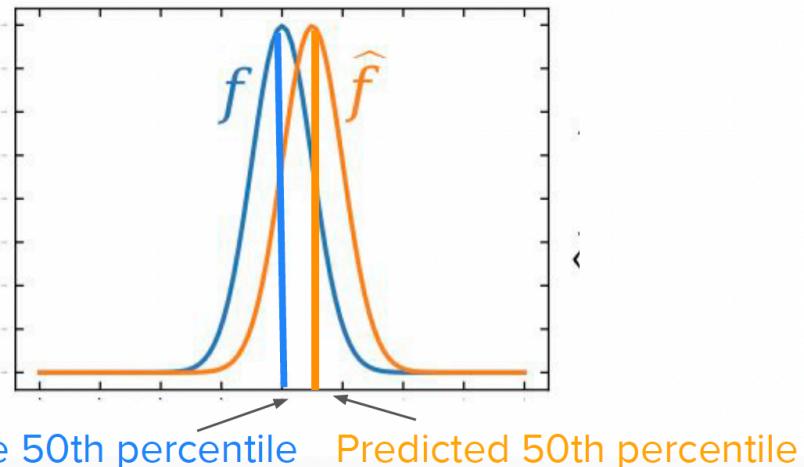
“What percentile level is my posterior model assigning to the true value, and with what frequency does this occur in the test set, e.g. **we should predict the true value below the 50-th percentile 50% of the time**”

Ho+24

Validation

PIT tests / P-P plots

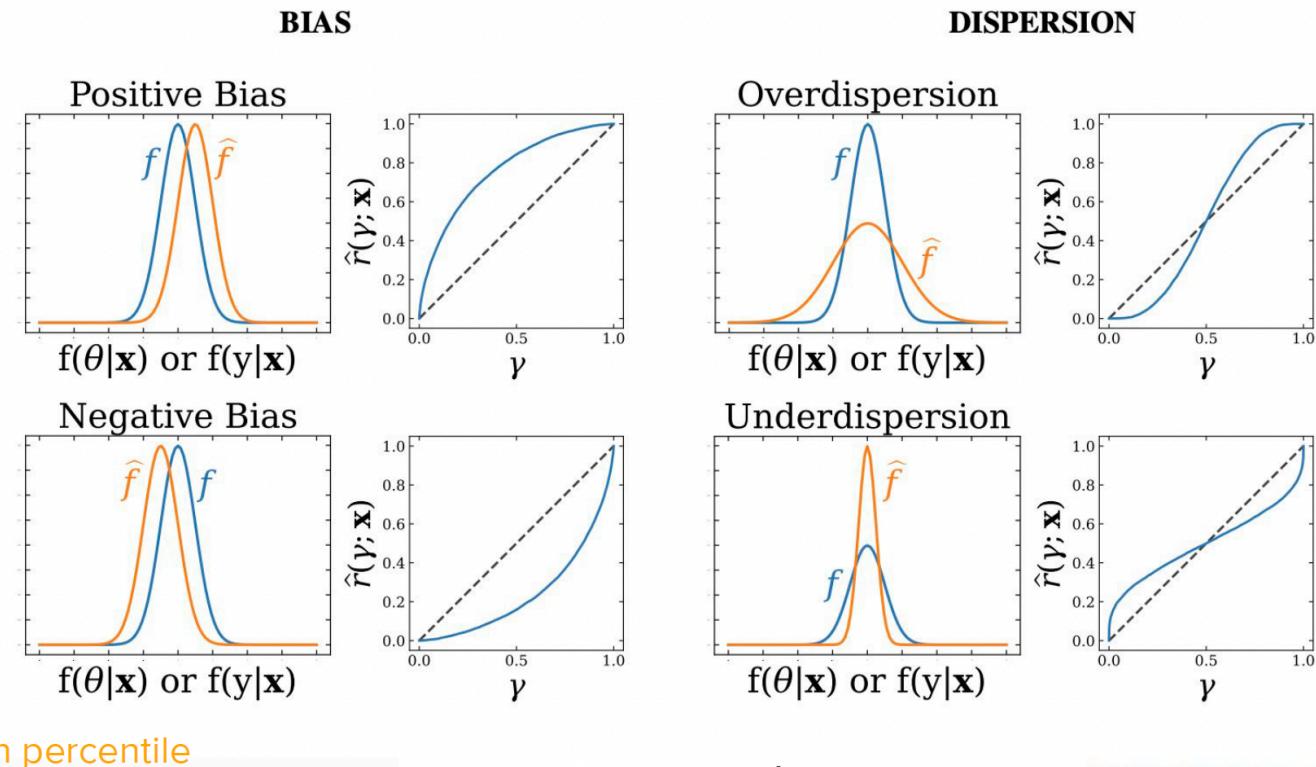
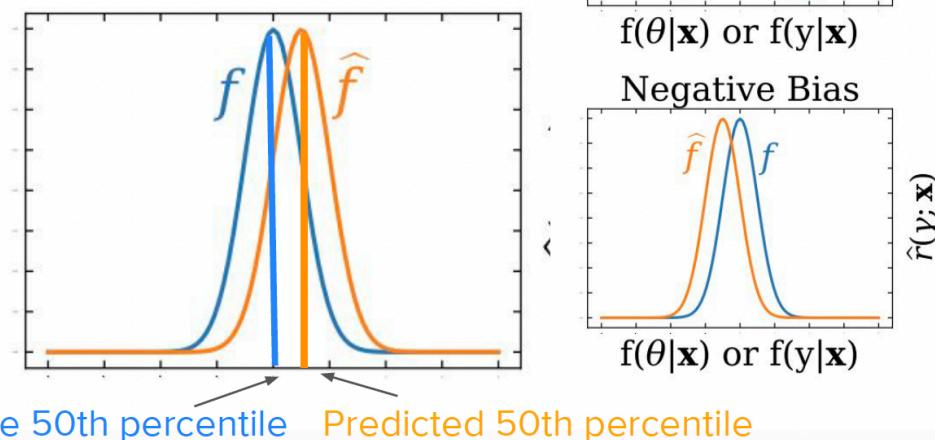
- Tests of marginal coverage



*adapted from M. Ho (SBI conference)

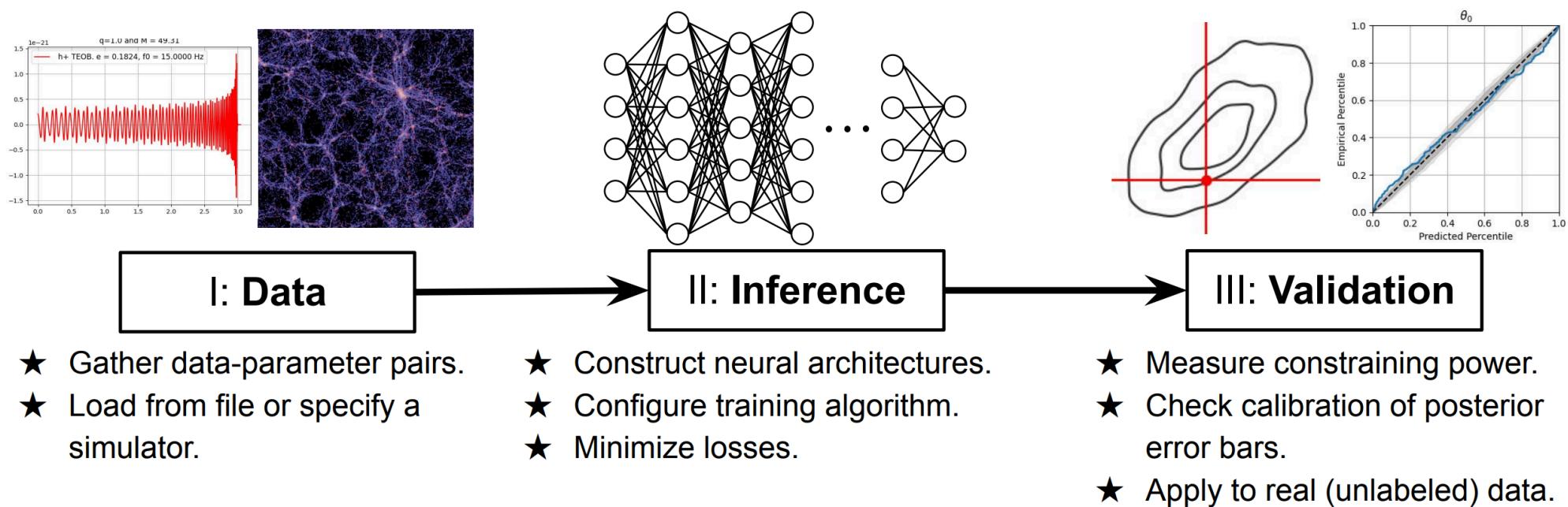
PIT tests / P-P plots

- Tests of marginal coverage



*adapted from M. Ho (SBI conference)

SBI pipeline



Ho+24

Examples of applications

Forward Model

θ

[SFR, M*, Av, SFH, Z]

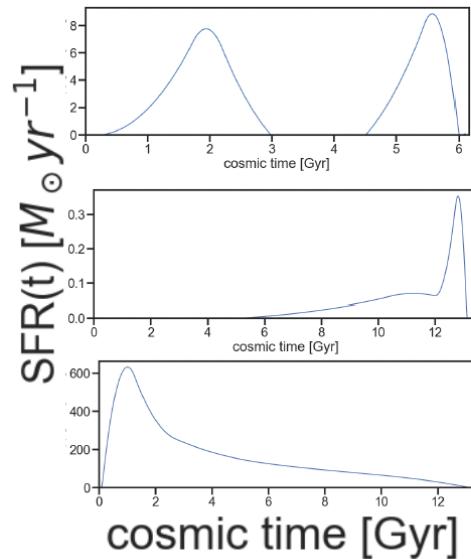


X_S

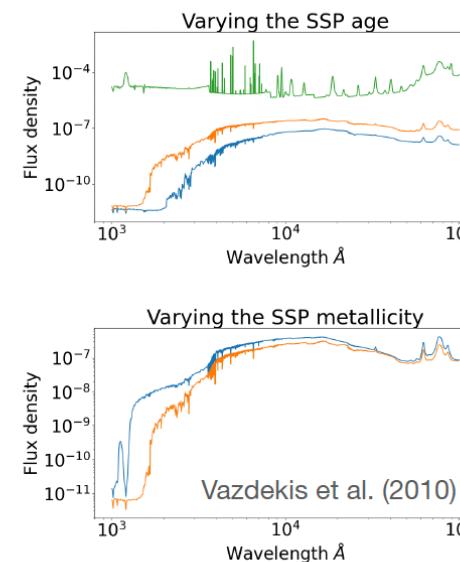
fluxes + noise model

$$F_\lambda(t) = \int_0^t \psi(t - t') SSP_\lambda(t', Z(t - t')) e^{-\tau(t')} dt'$$

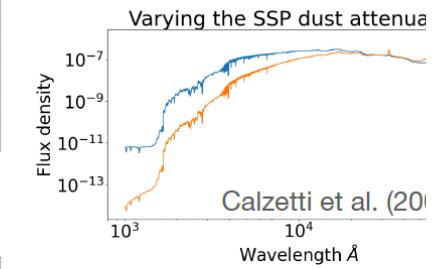
Non-parametric SFHs with a Dirichlet prior



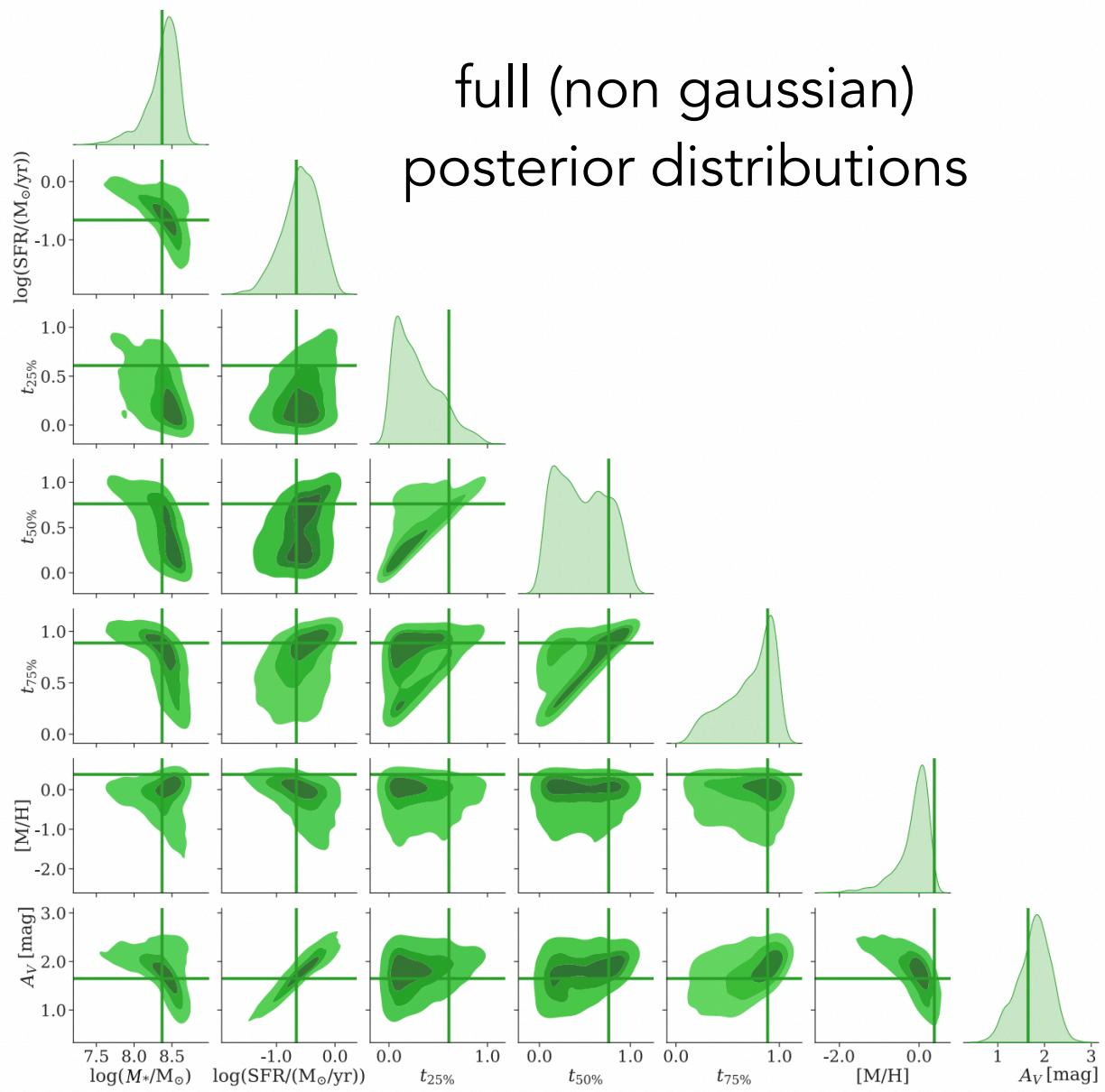
MILES SSP models

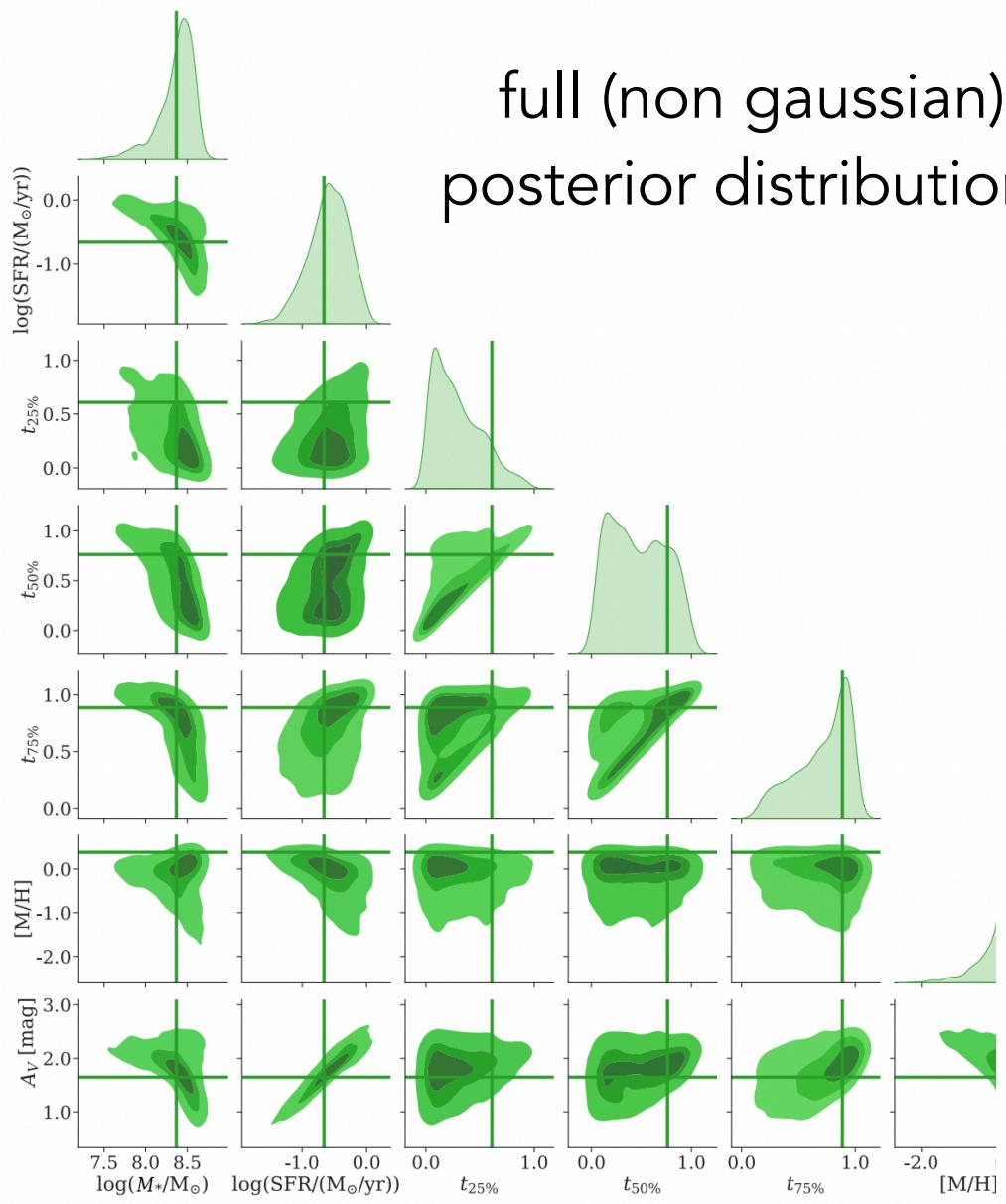


Dust attenuation model

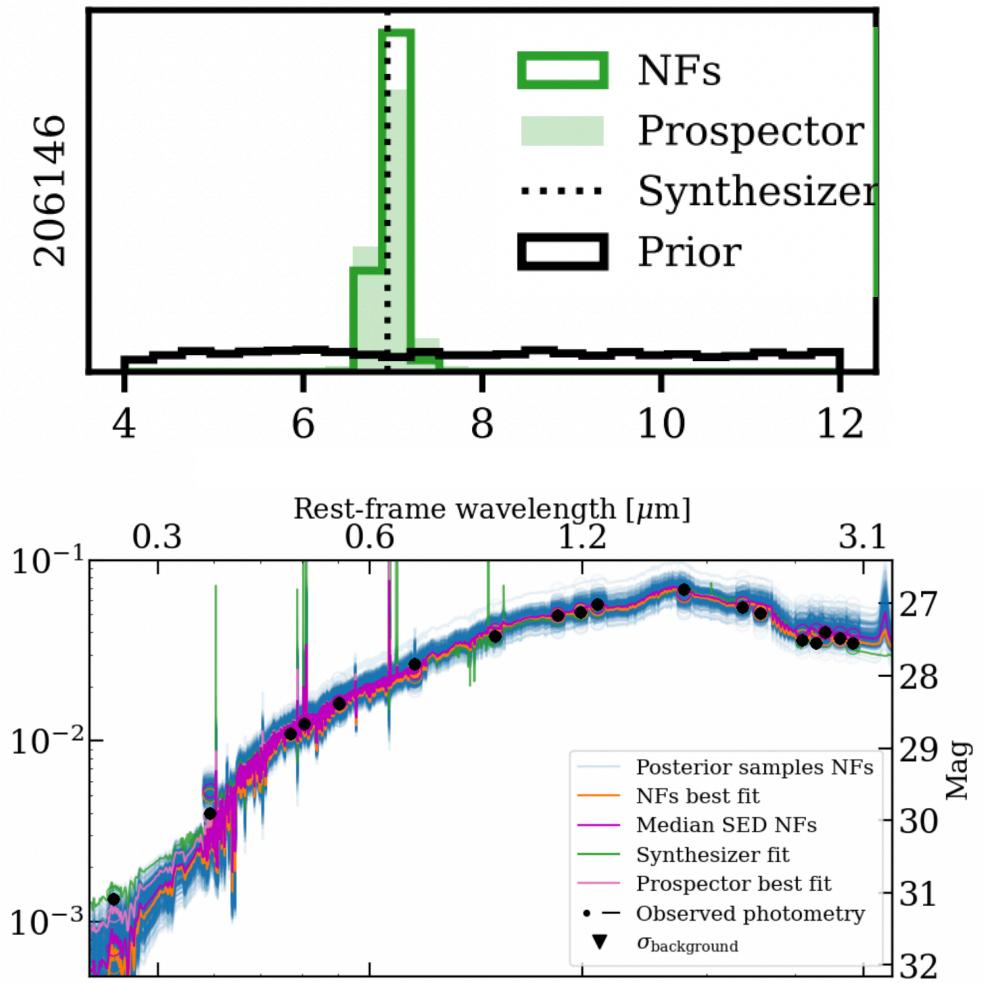


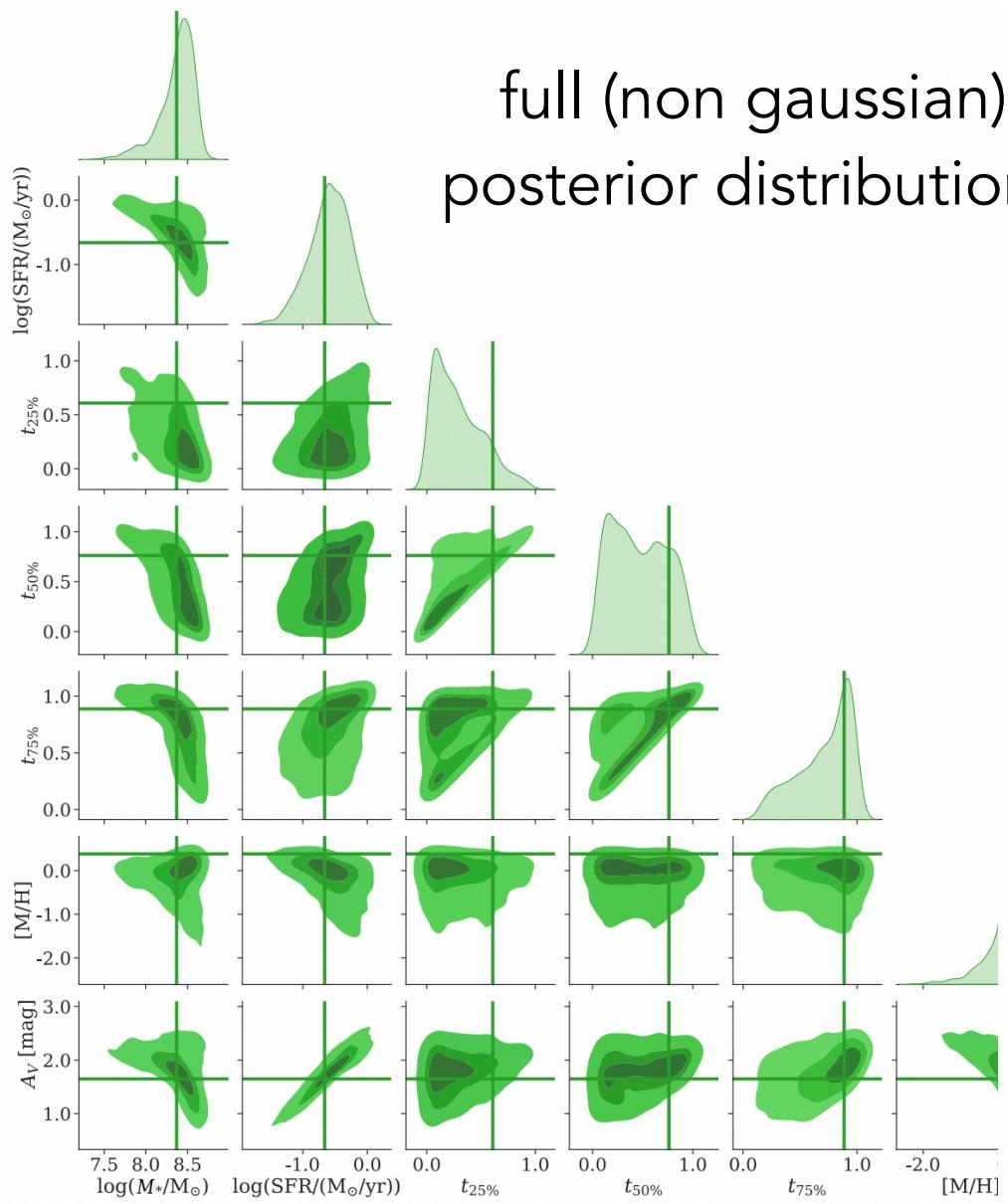
full (non gaussian) posterior distributions





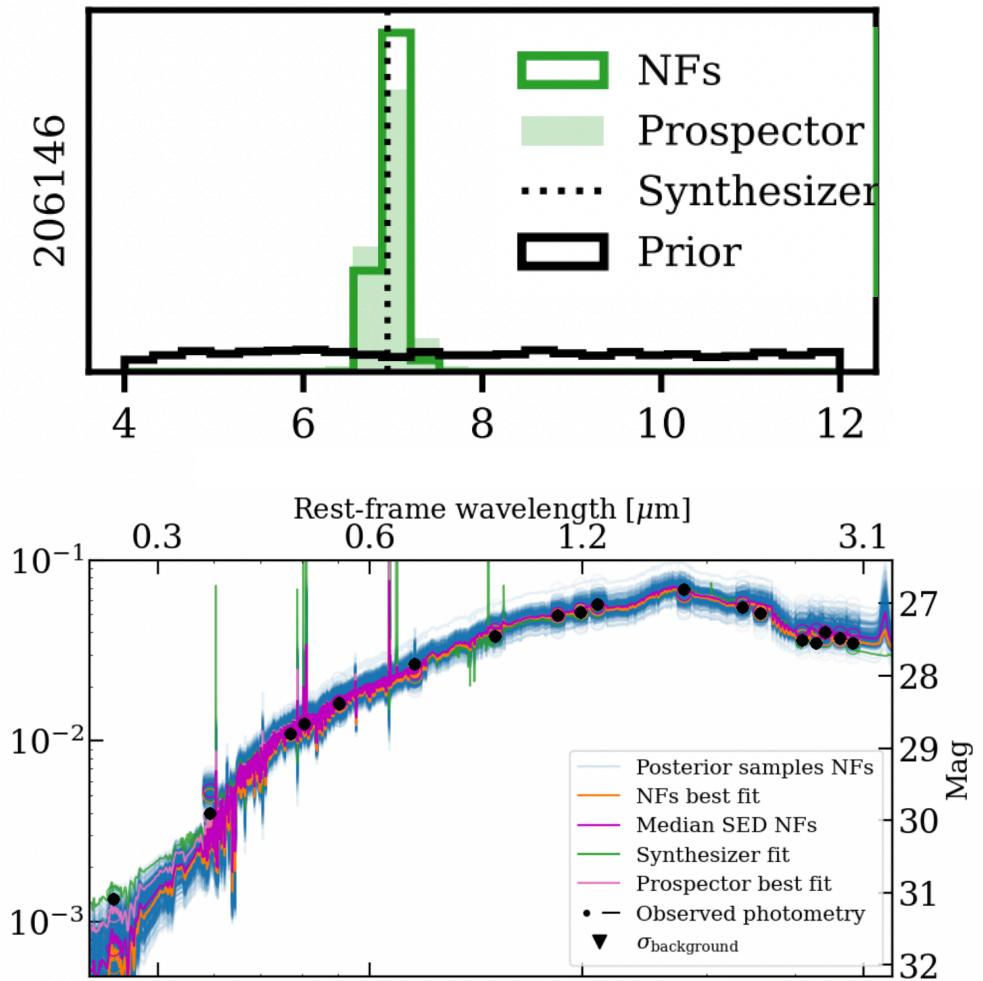
in good agreement with
MCMC



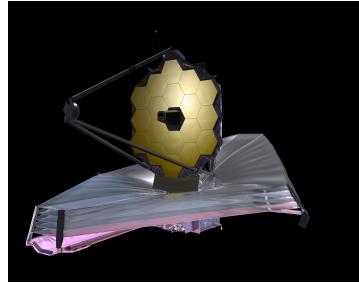


full (non gaussian)
posterior distributions

in good agreement with
MCMC



~ 3 orders of magnitude faster than Bayesian SED fitting



JWST + SBI enable
resolved stellar populations
of galaxies up to very high redshift

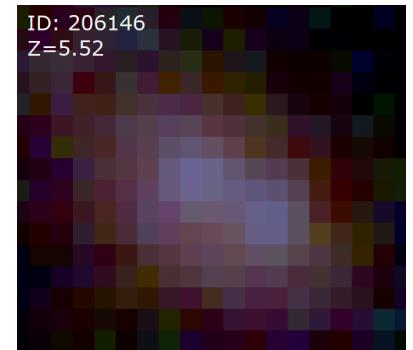
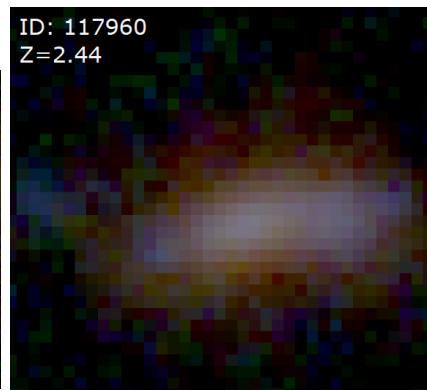
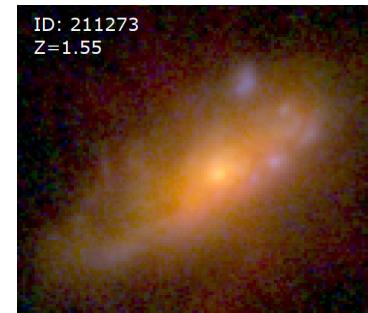
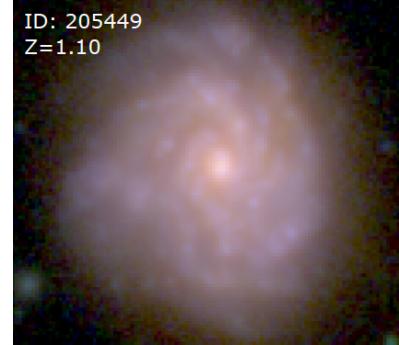
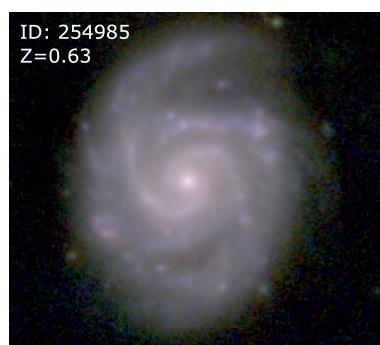
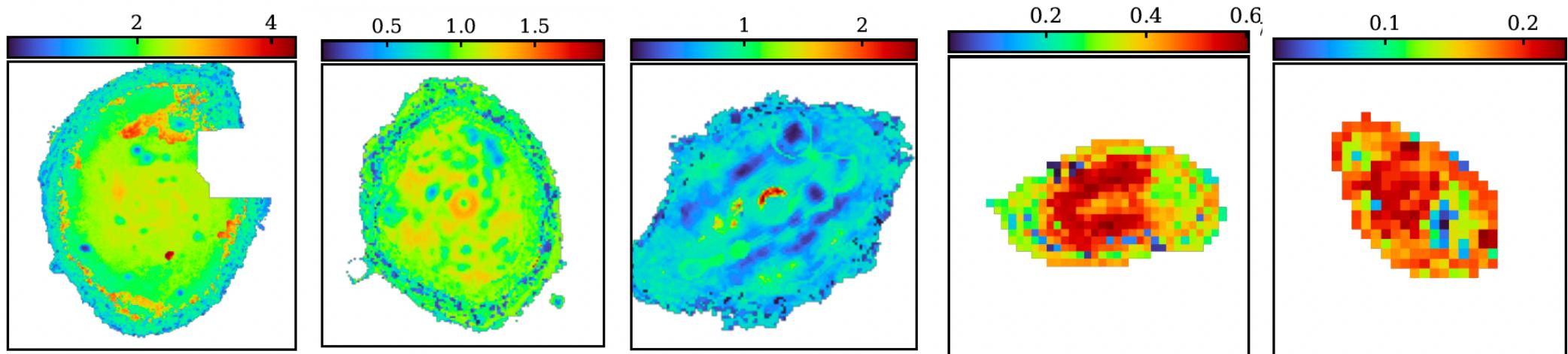


Fit every pixel
(~10k)

ML based SBI of resolved stellar populations



Mass-weighted
age [Gyr]



$z=0.63$

$z=1.10$

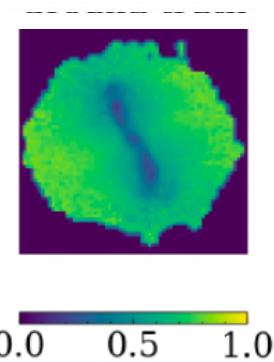
$z=1.55$

$z=2.44$

$z=5.52$

Forward Model: Cosmological simulations

θ

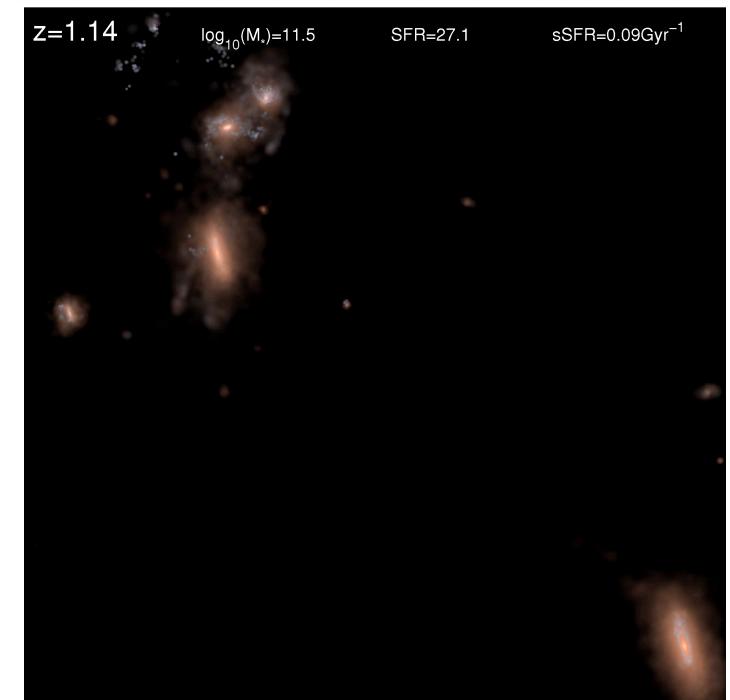
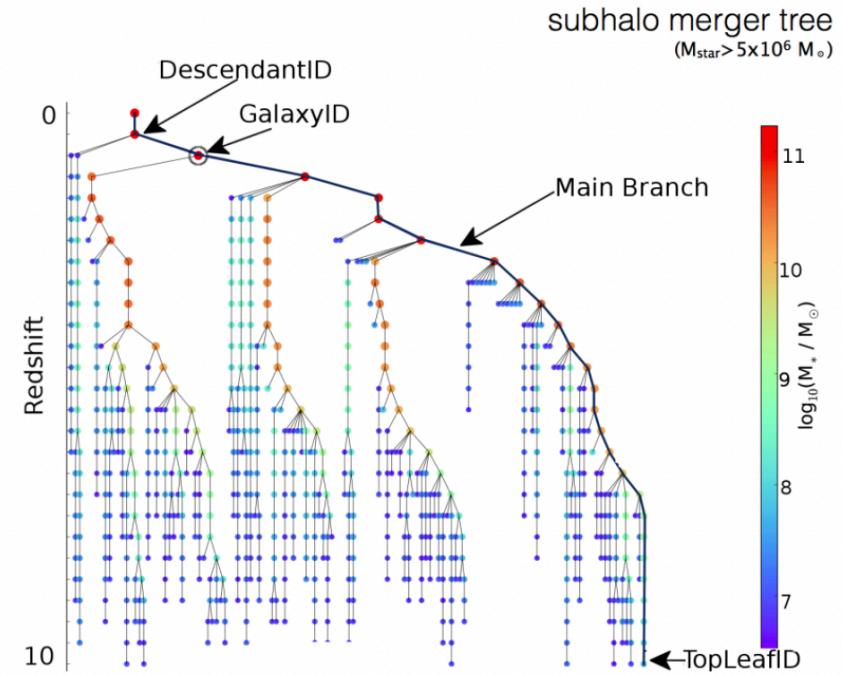


(nD)



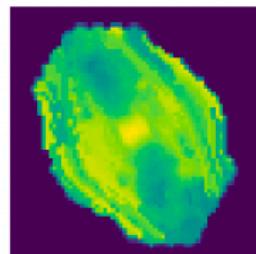
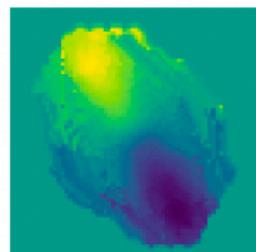
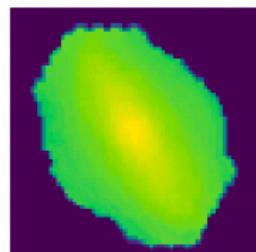
[kinematic and
stellar population
maps]

(nD)



Resolved accretion map of local galaxies via diffusion

Observable maps



Conditioning

Normalizing

Adding FF

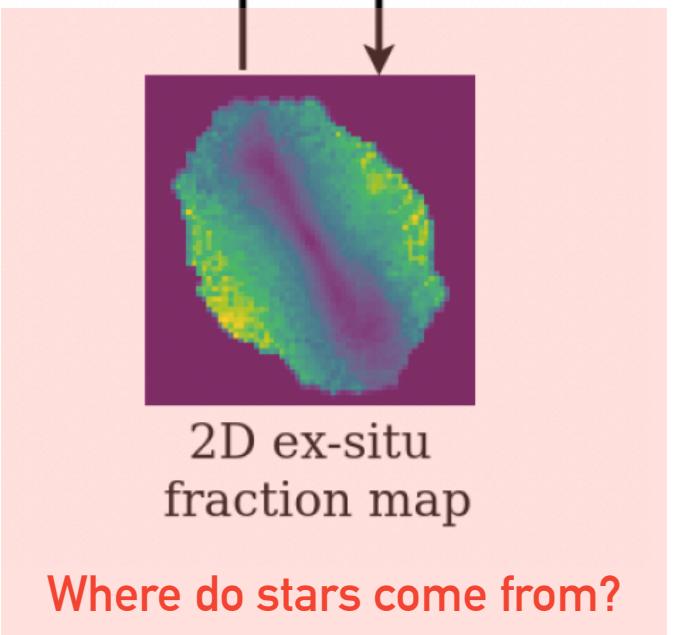
Adding gaussian noise

Data augmentations

Diffusion model

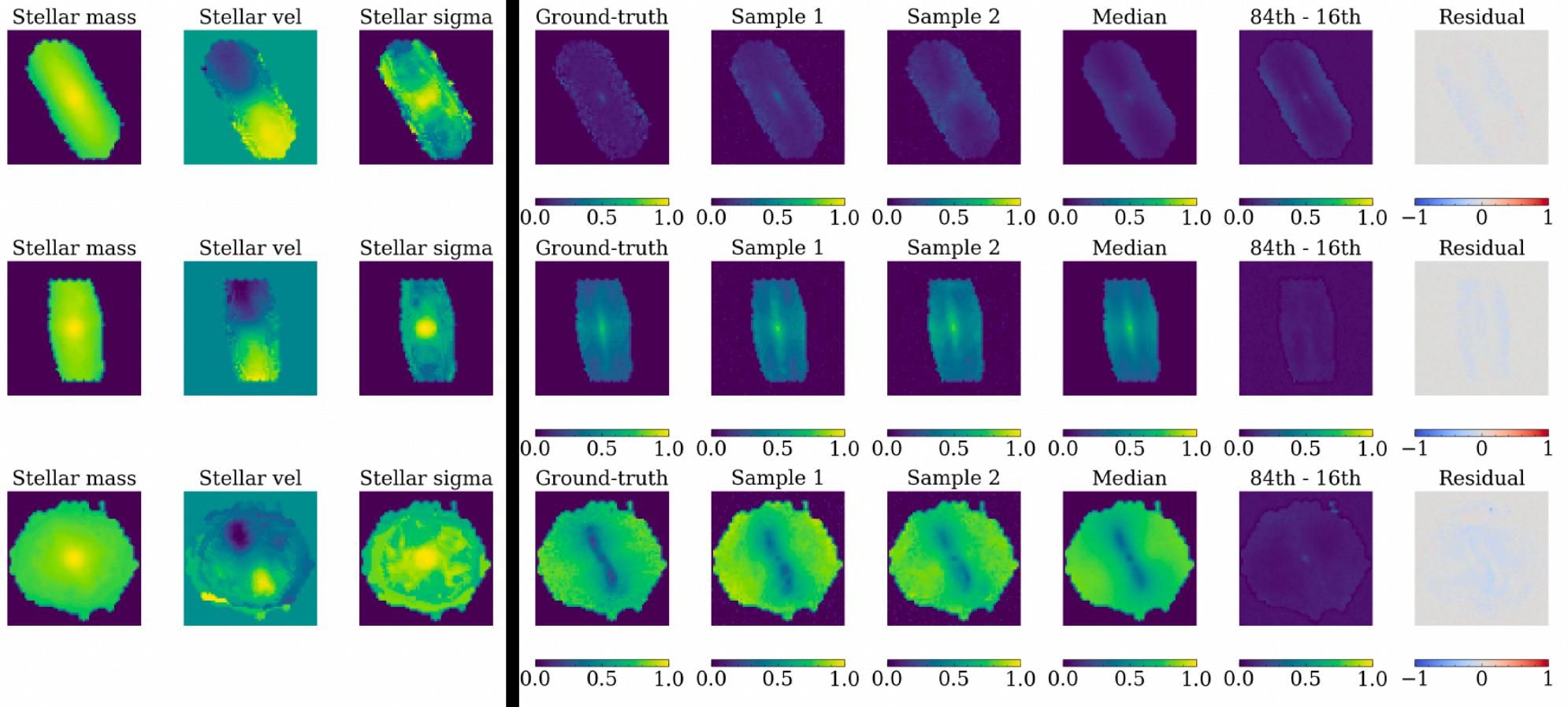
Training

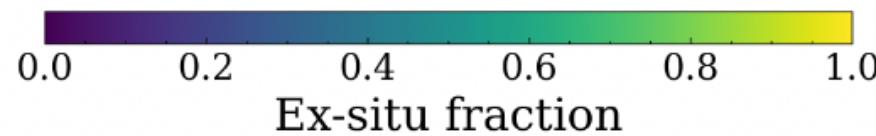
Sampling



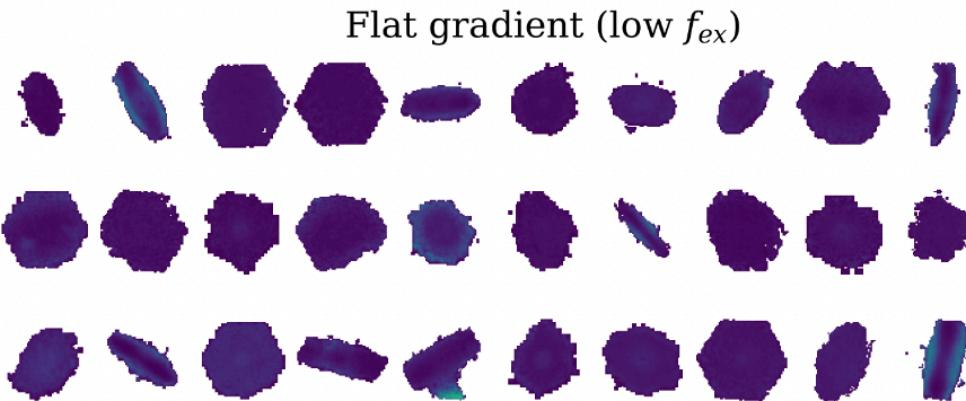
E. Angeloudi

Angeloudi, MHC+25

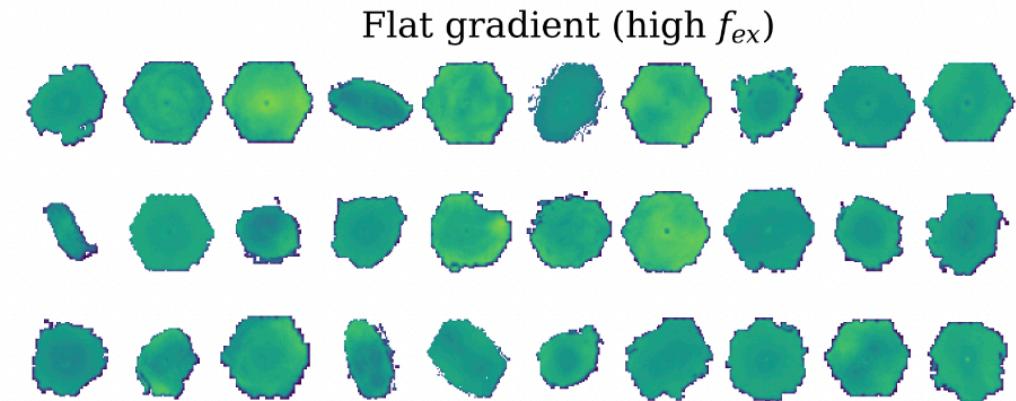
X θ 

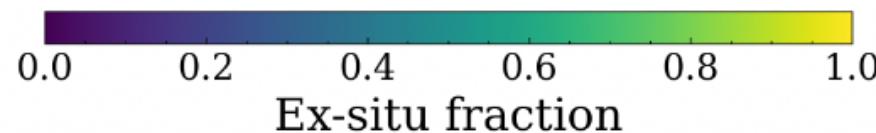


In-situ dominated



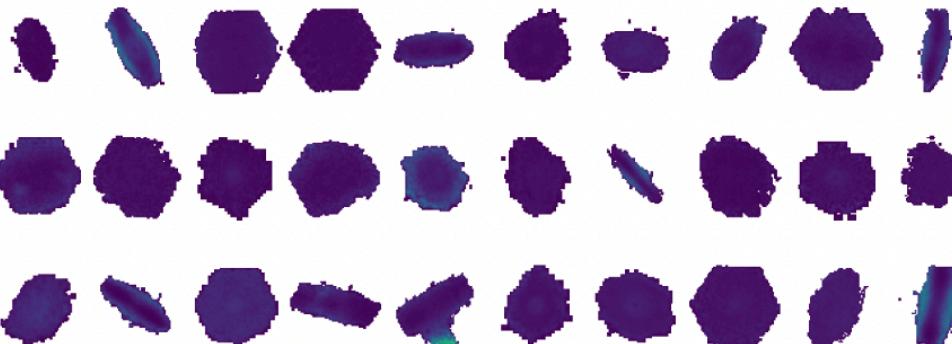
merger dominated





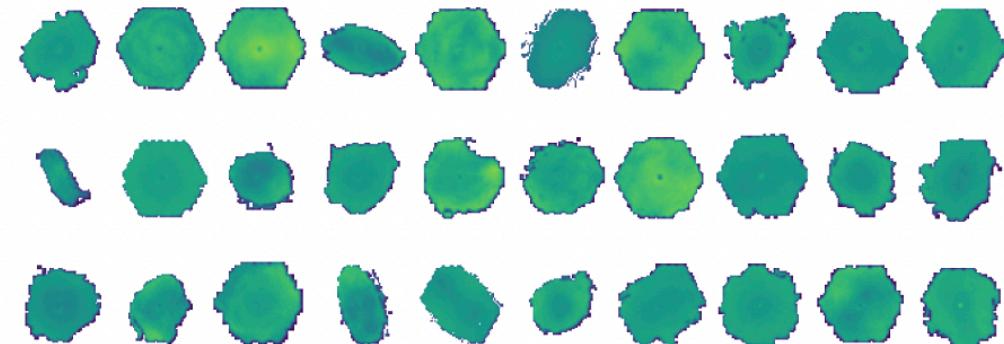
In-situ dominated

Flat gradient (low f_{ex})



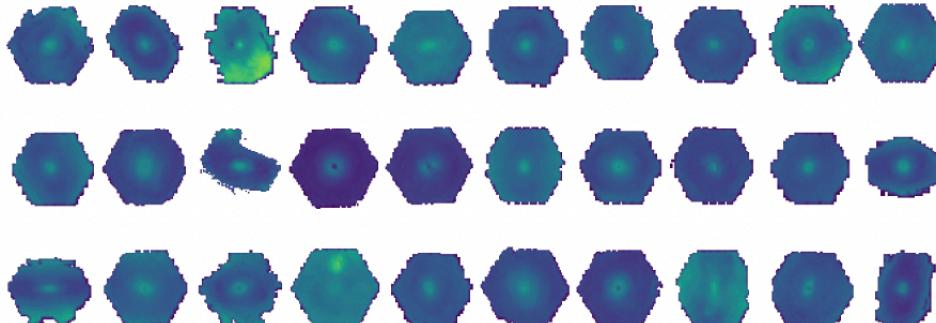
merger dominated

Flat gradient (high f_{ex})



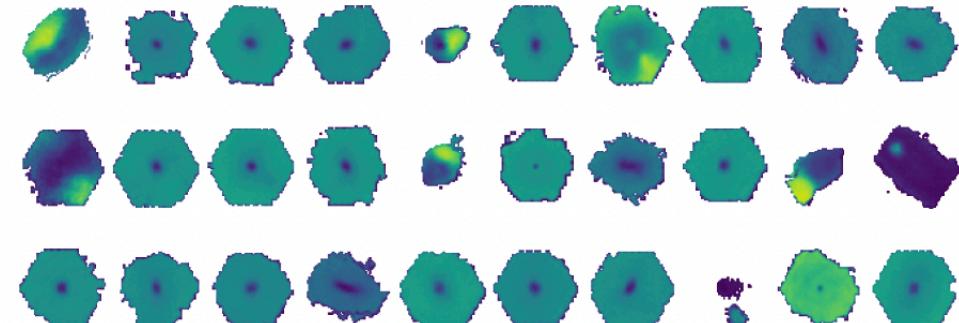
more ex-situ in center

Negative gradient

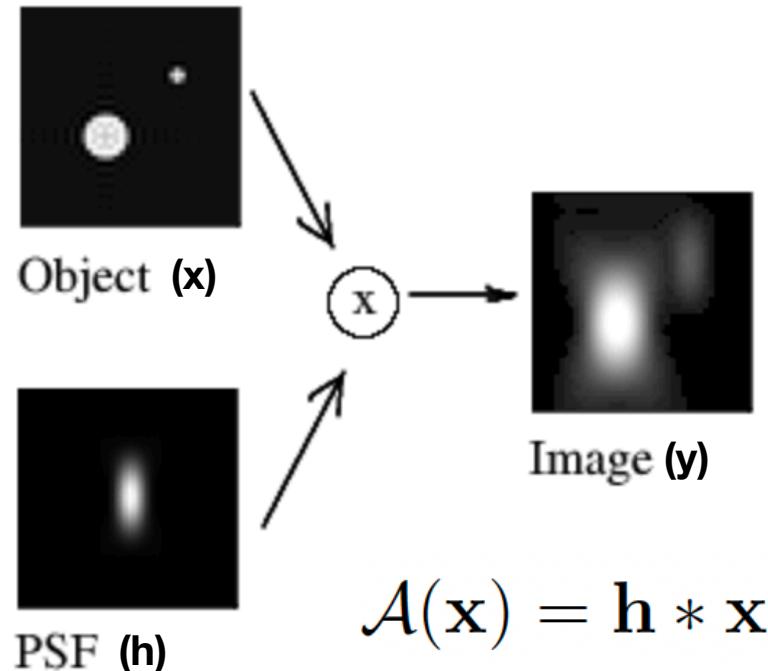


more ex-situ in outskirts

Positive gradient



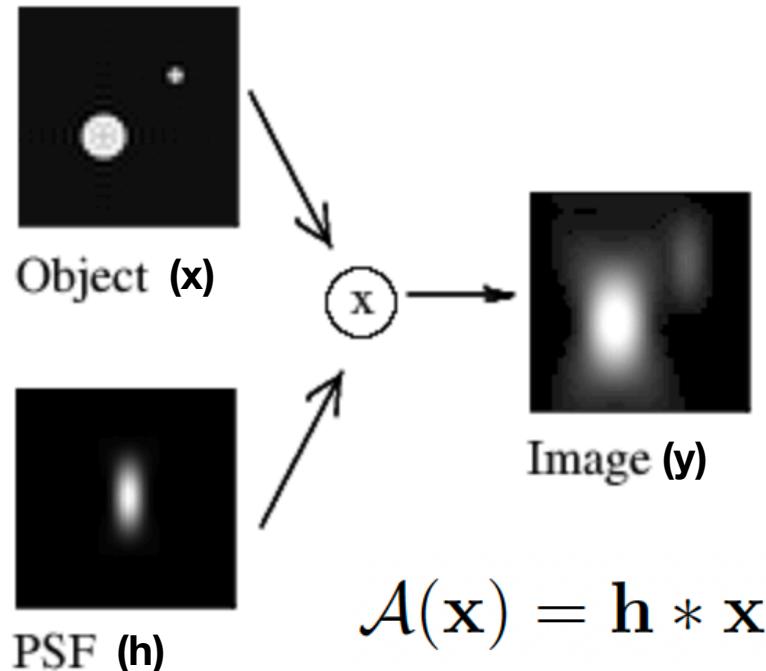
Learning the prior: Deconvolution through diffusion



$$\mathbf{y} = \mathcal{A}\mathbf{x}_0 + \sigma_y \mathbf{n}$$

$$\mathcal{A}(\mathbf{x}) = \mathbf{h} * \mathbf{x}$$

Learning the prior: Deconvolution through diffusion



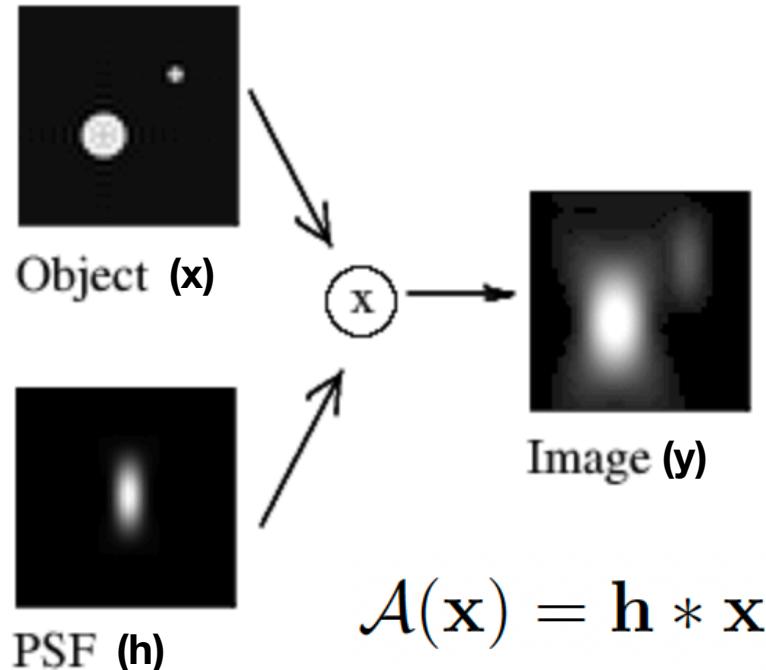
$$\mathbf{y} = \mathcal{A}\mathbf{x}_0 + \sigma_y \mathbf{n}$$

$$\log p(x|y) = \log p(y|x) + \log p(x)$$

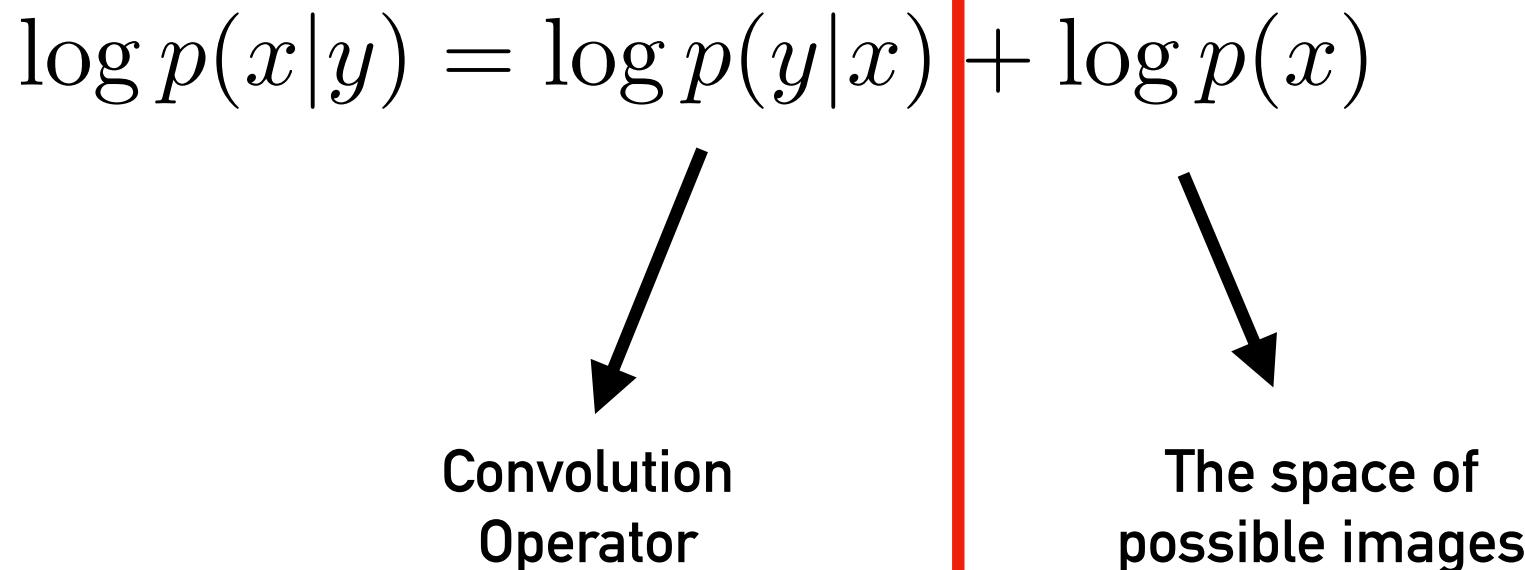
Convolution
Operator

The space of
possible images

Learning the prior: Deconvolution through diffusion



$$\mathbf{y} = \mathcal{A}\mathbf{x}_0 + \sigma_y \mathbf{n}$$



Learning the prior: Deconvolution through diffusion

$$\mathbf{y} = \mathcal{A}\mathbf{x}_0 + \sigma_y \mathbf{n}, \quad \mathcal{A}(\mathbf{x}) = \mathbf{h} * \mathbf{x}$$

$$\log p(x|y) = \log p(y|x) + \log p(x)$$

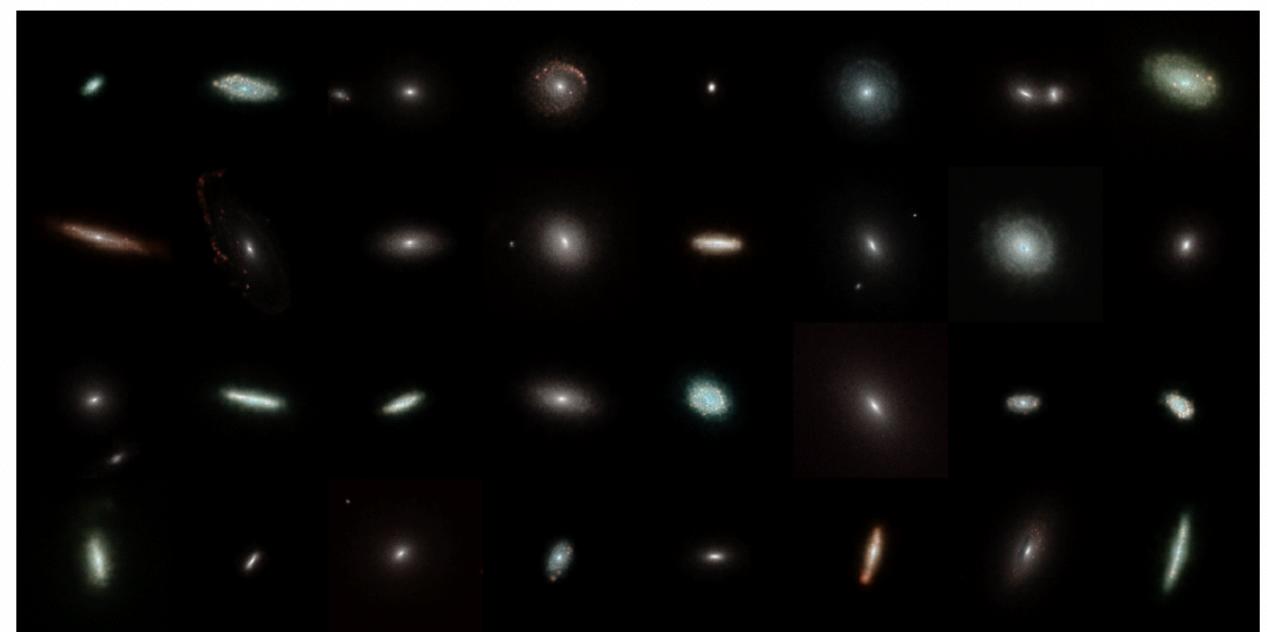
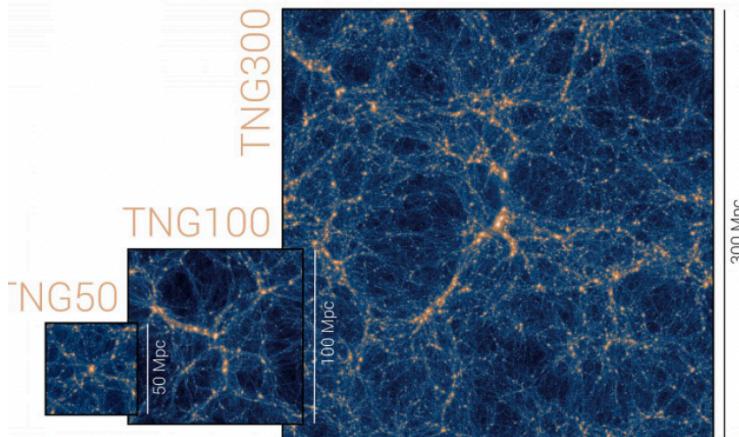
Convolution
Operator

The space of
possible images

neural network

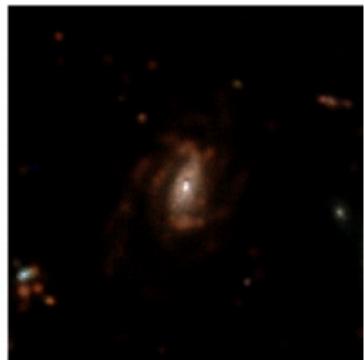
$$p(x)$$

Learned through diffusion models and cosmological simulations

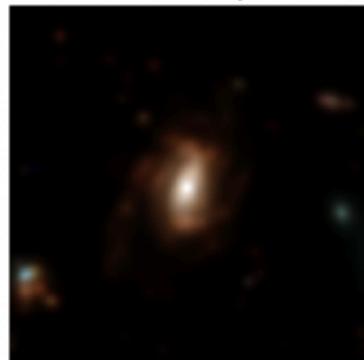


Samples from the prior

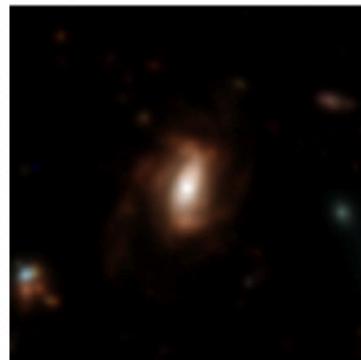
Deconvolved



Processed x_0



HSC



HST

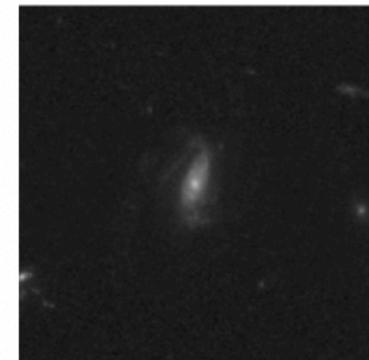
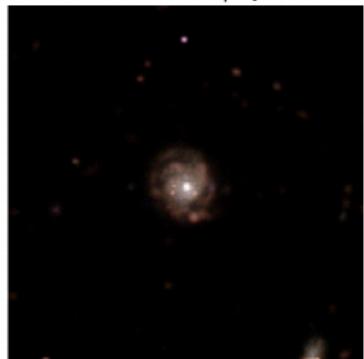


Figure 25: Object 1

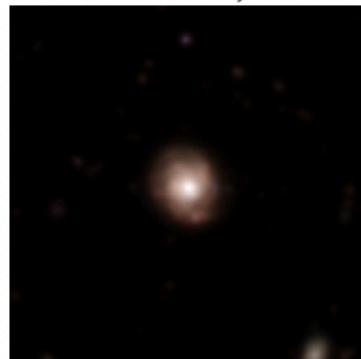
Diffusion step x_0



Processed x_0



Observation y



HST

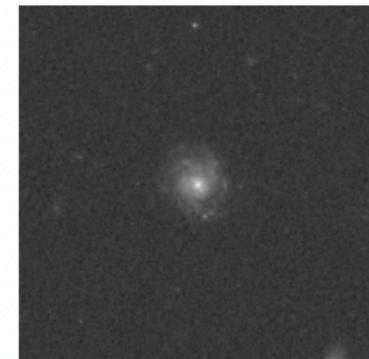
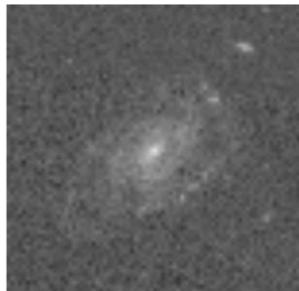


Figure 26: Object 2

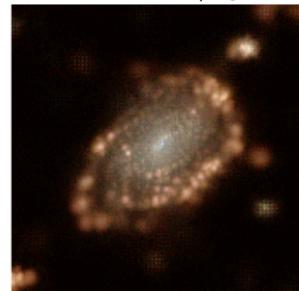
Observation y



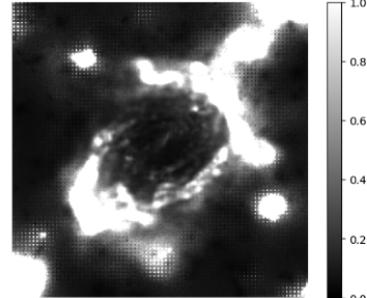
HST



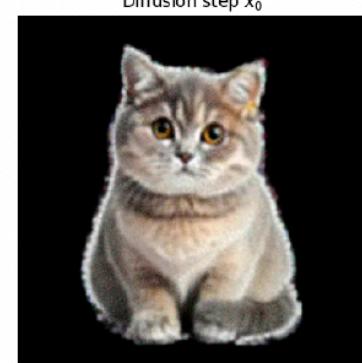
Diffusion step x_0



Variance Ratio



Deconvolved



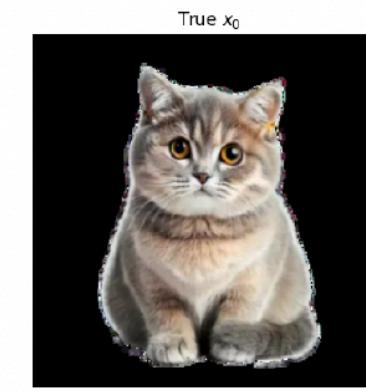
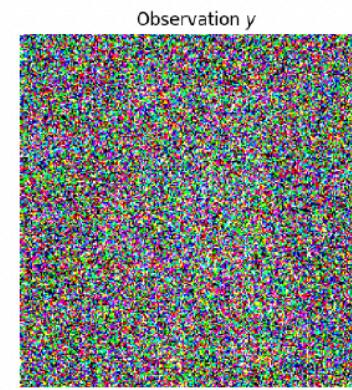
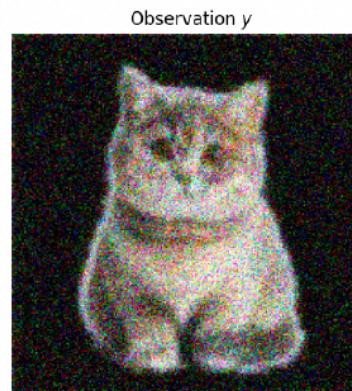
Observed

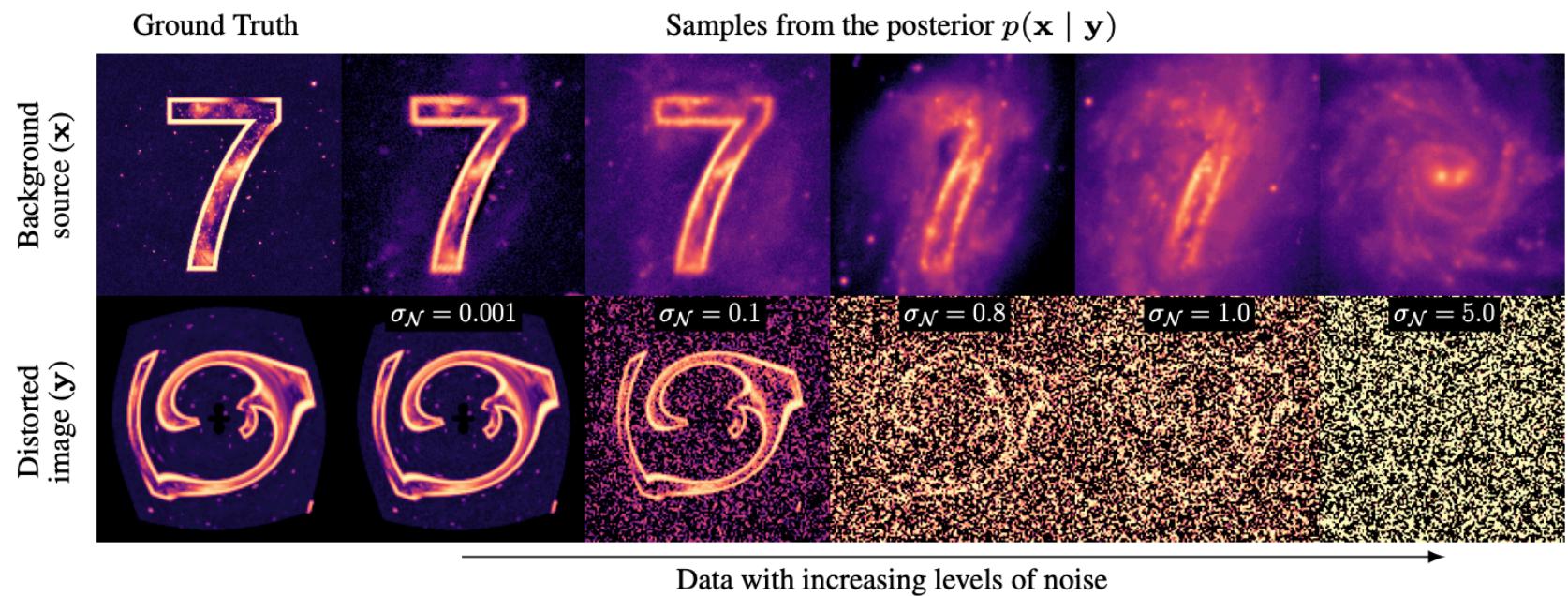
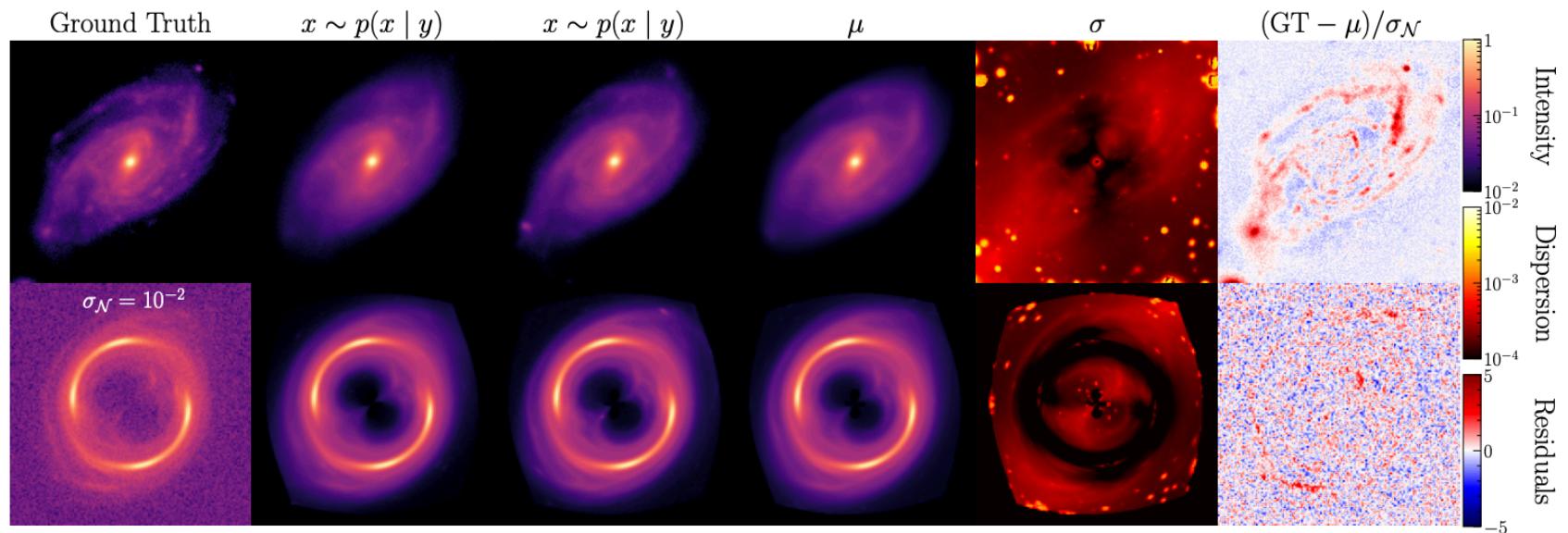


True



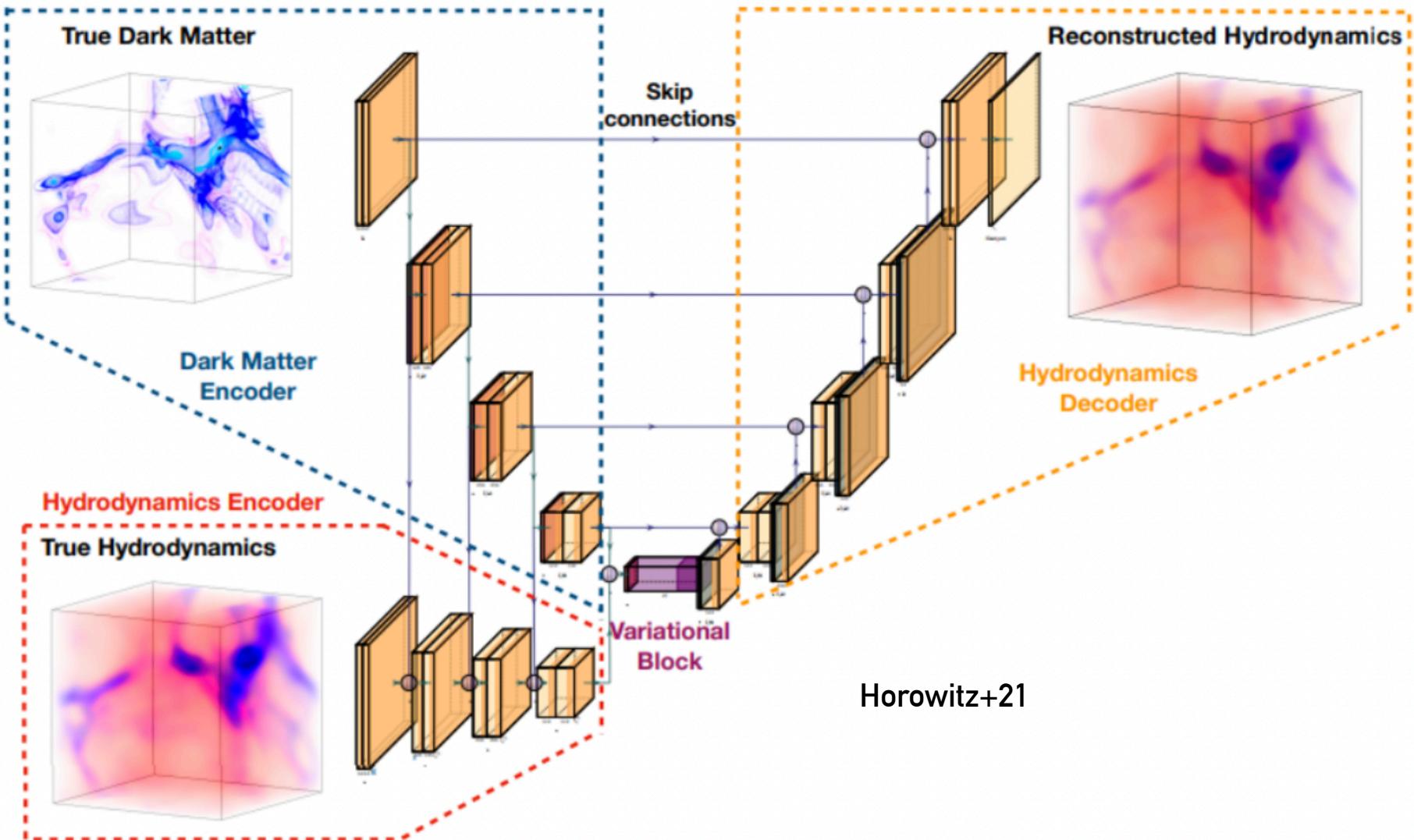
increasing
noise





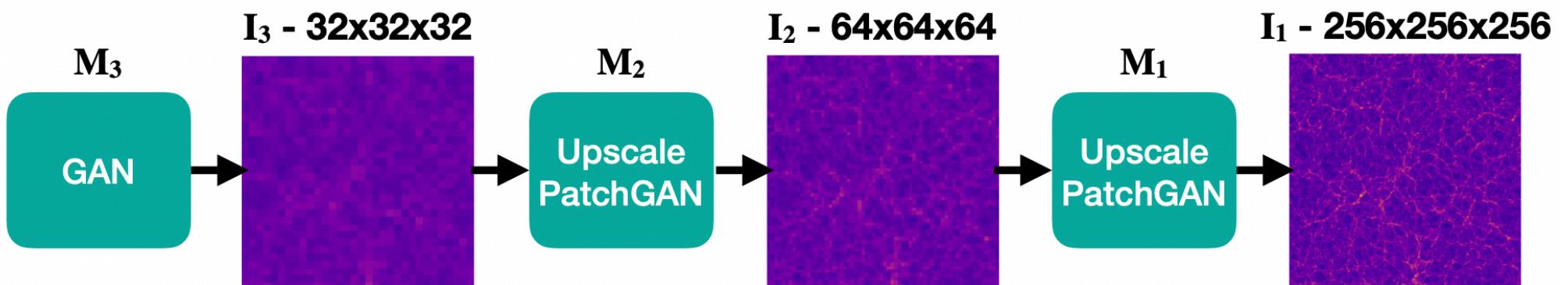
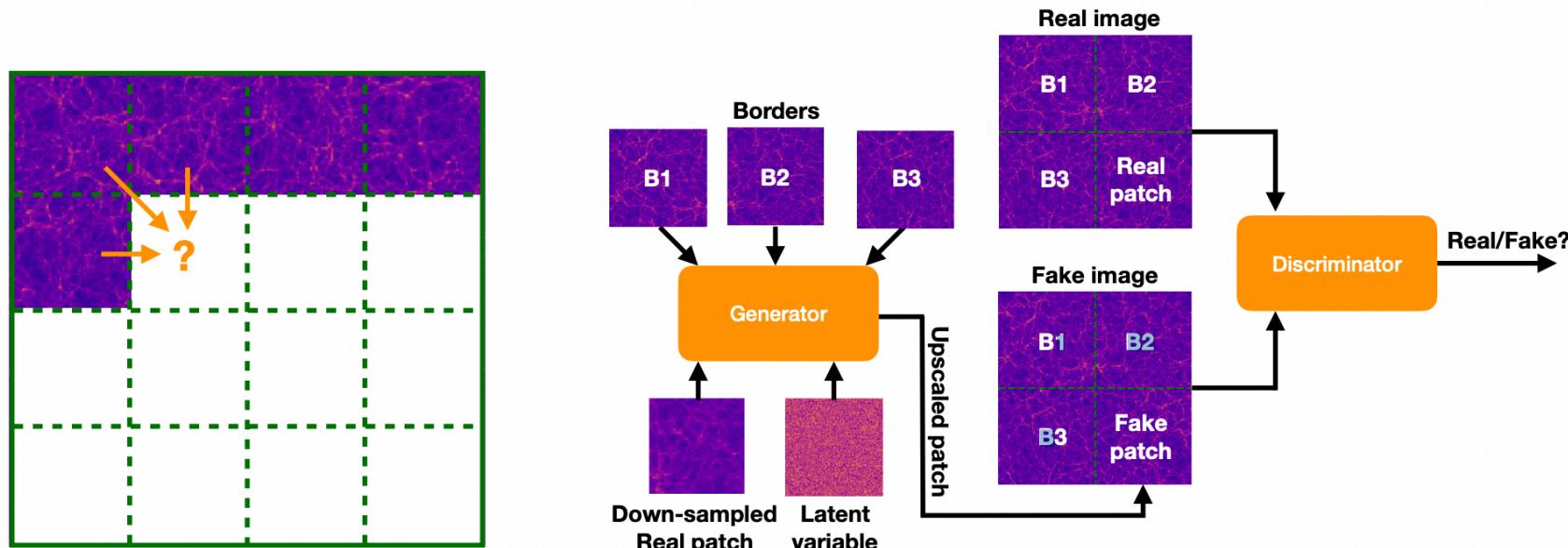
Adam+22

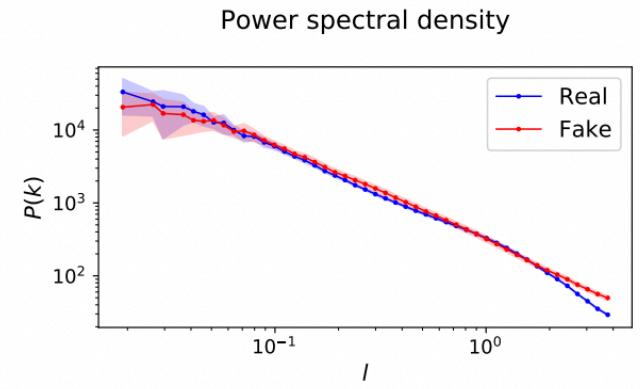
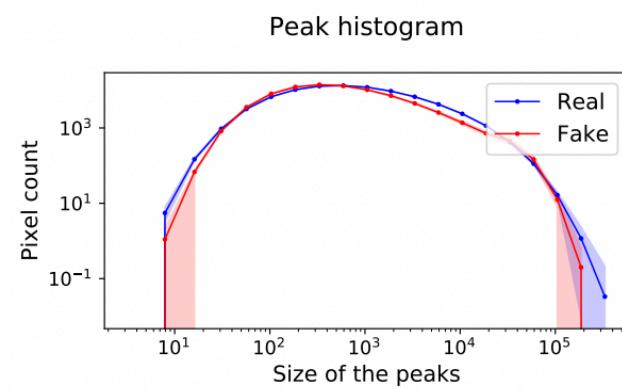
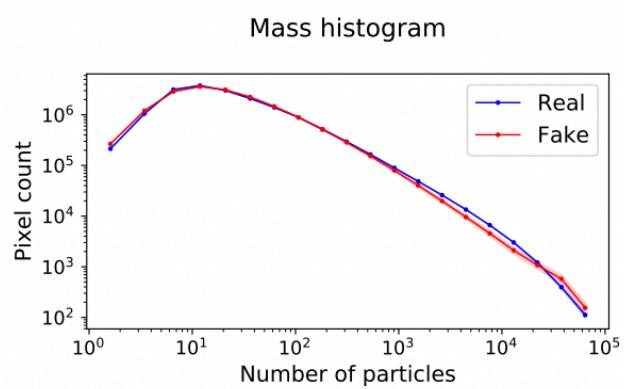
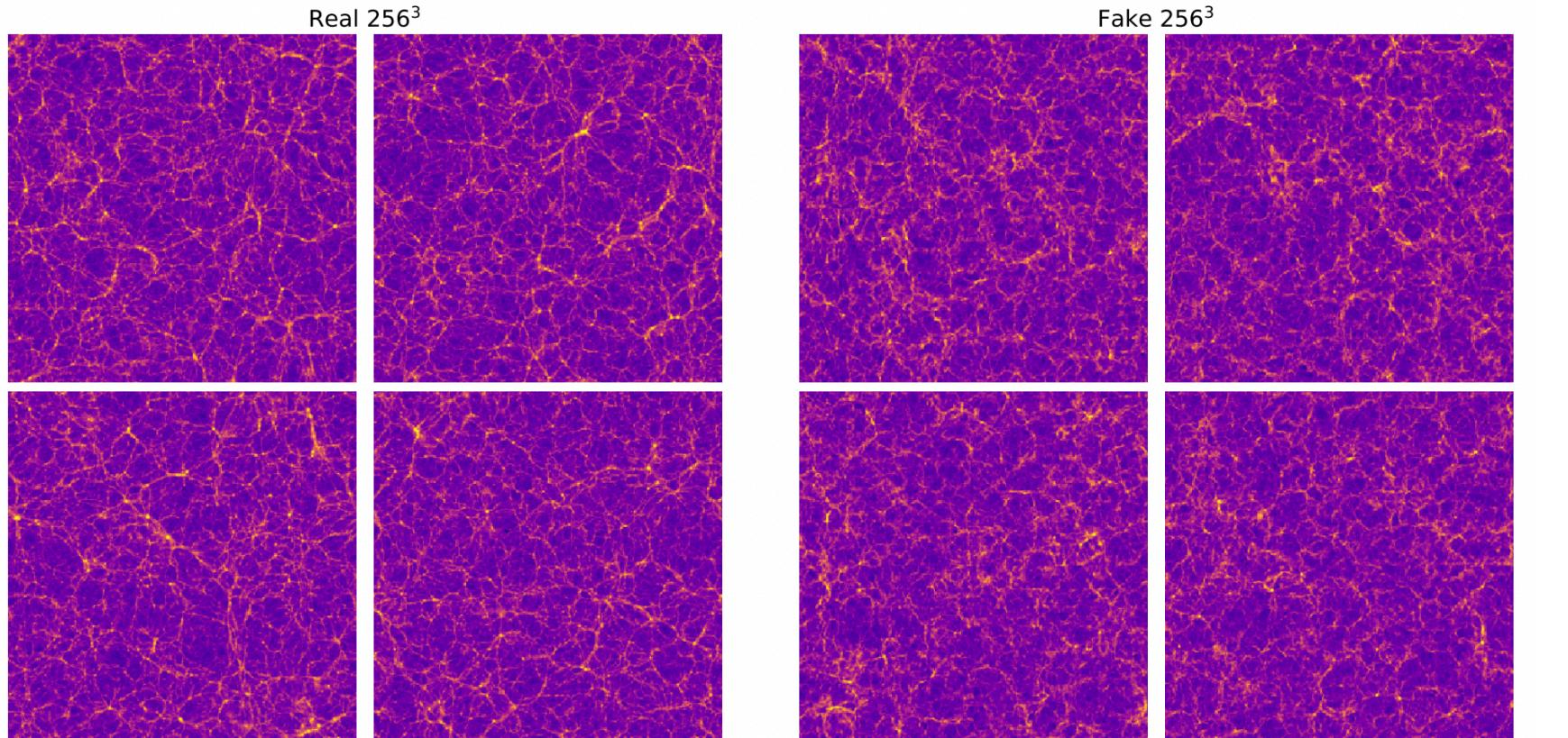
Painting Baryons



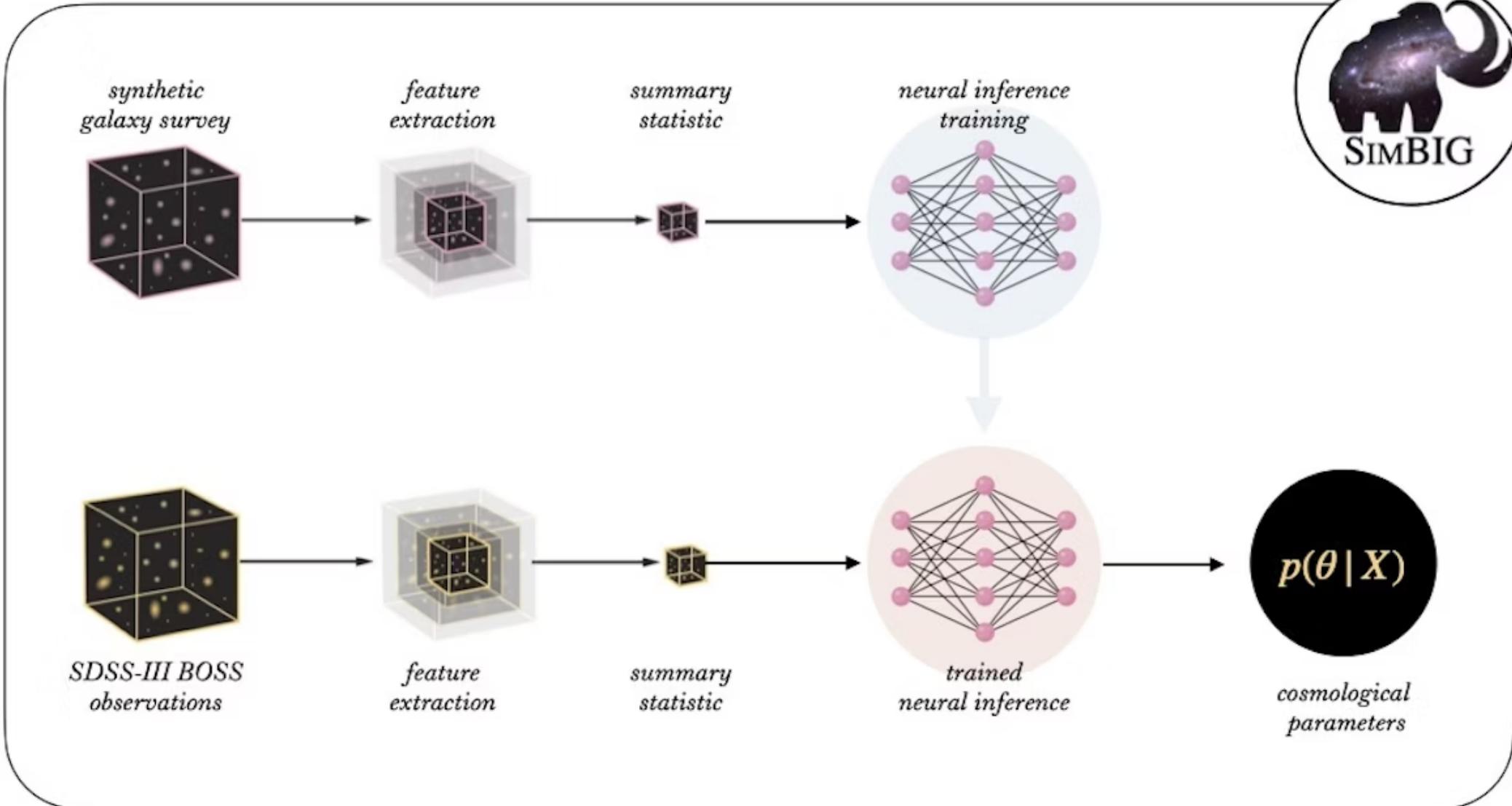
Horowitz+21

N-body emulation by Deep Generative Modelling



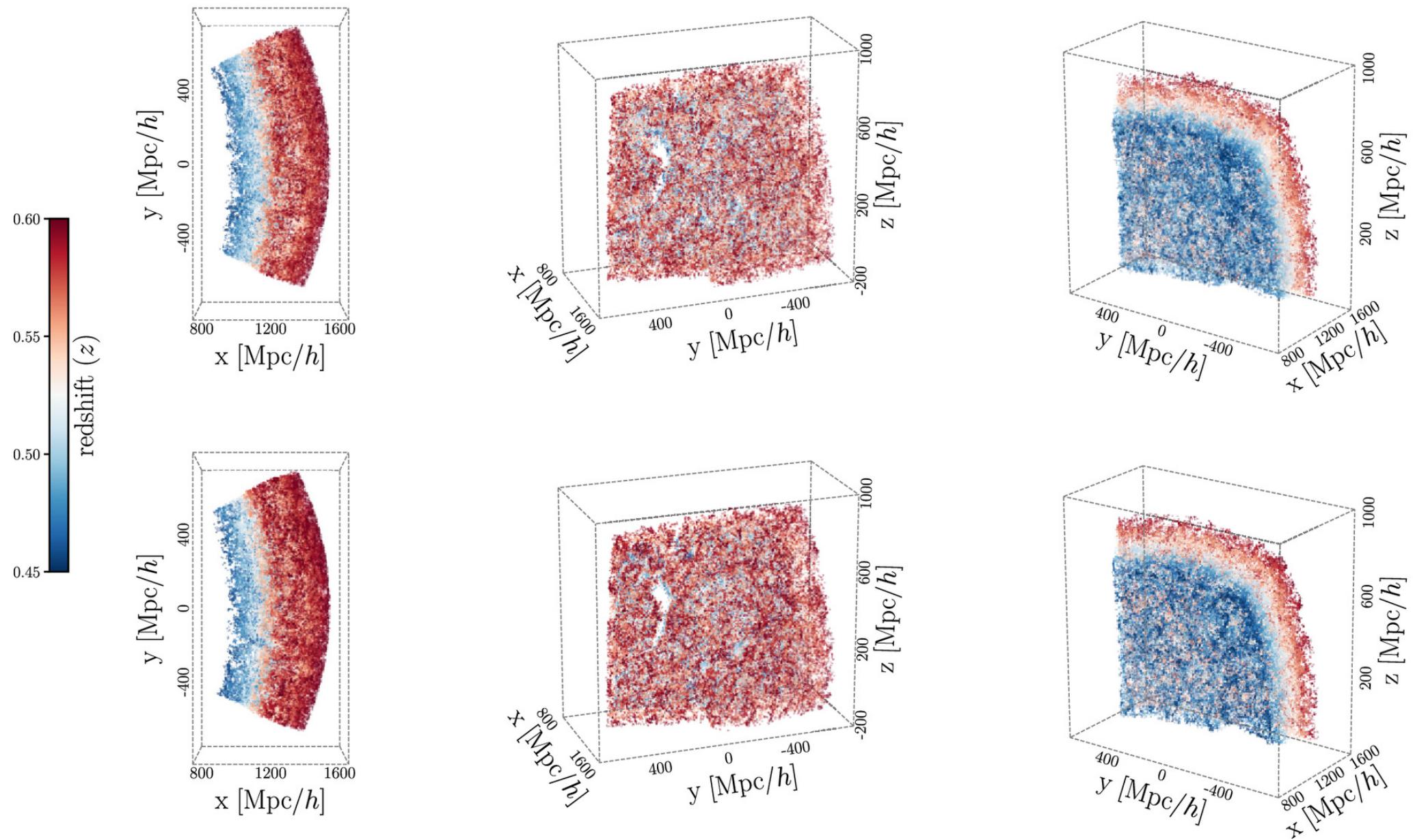


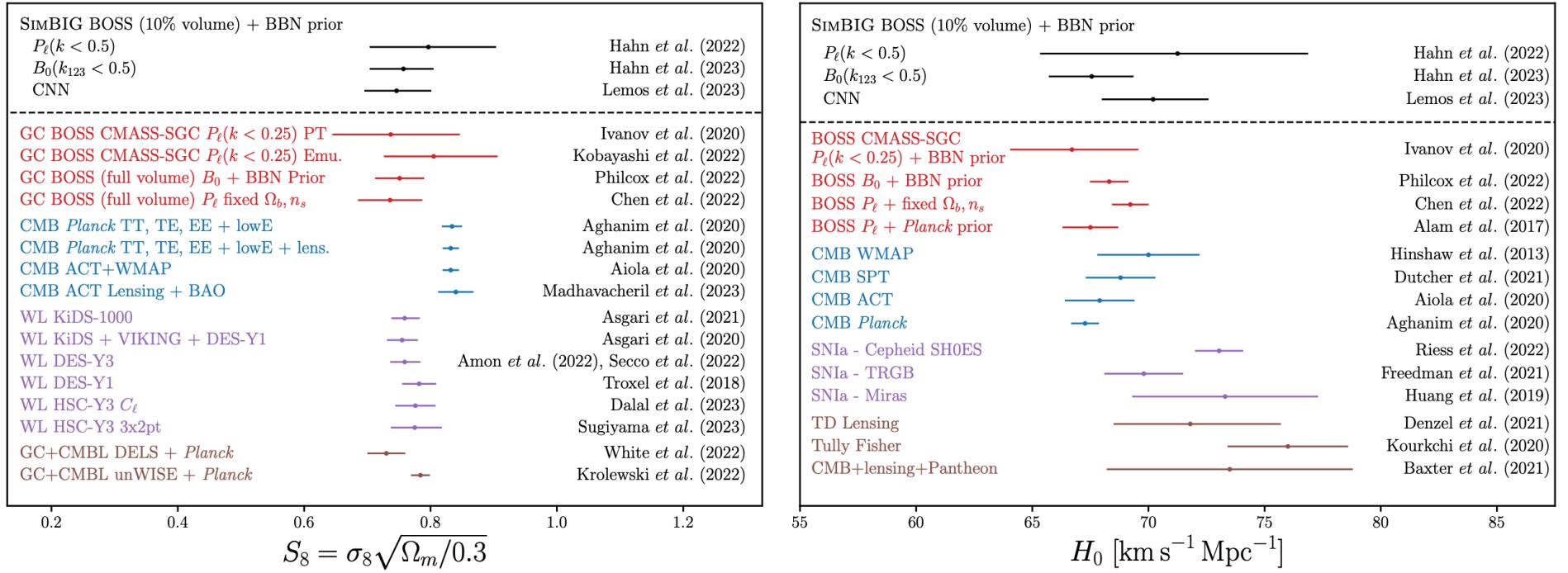
SBI based cosmological inference



“If you can simulate, you can do inference”

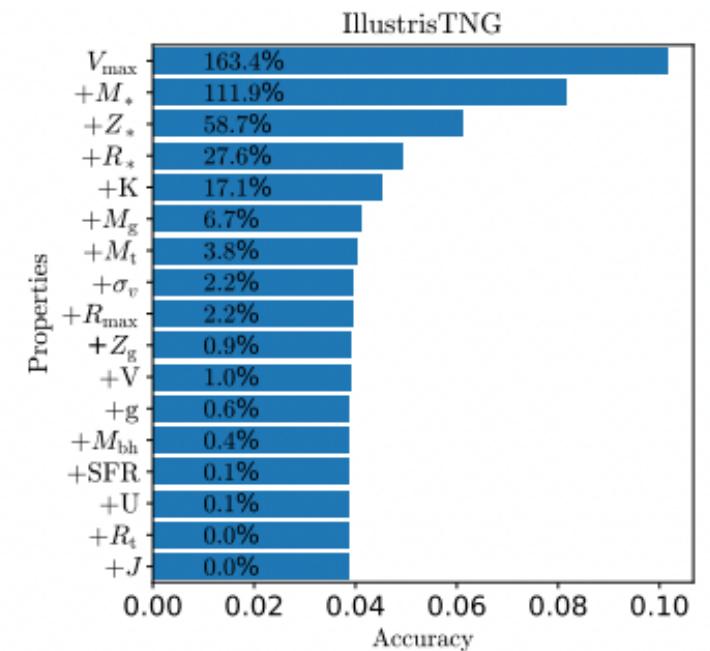
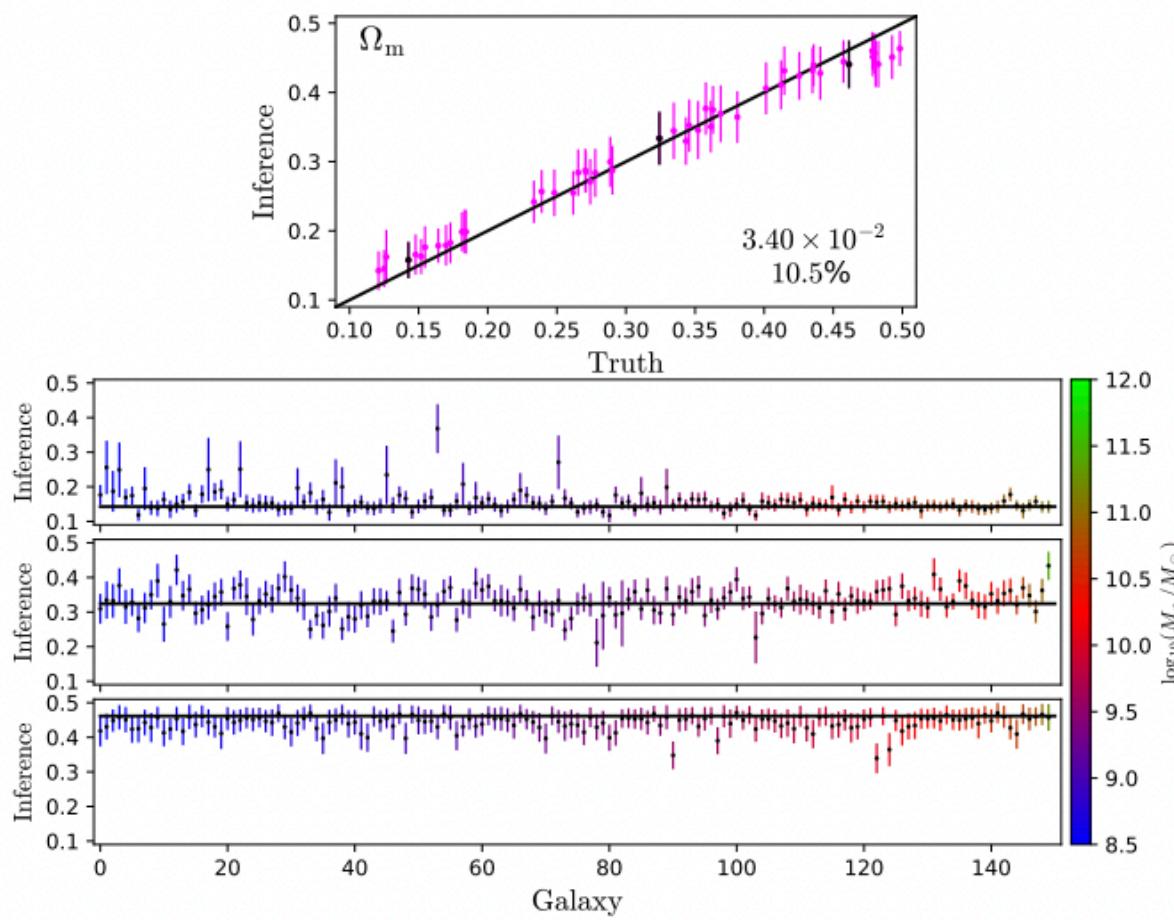
Towards Field Level Cosmological Inference





Hahn+23

Towards Field Level Cosmological Inference

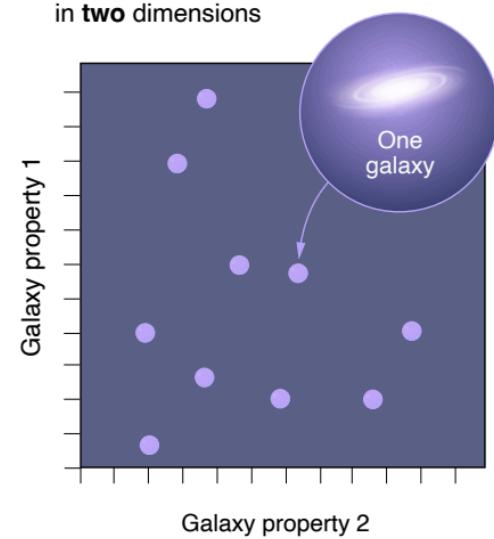


(Cosmology and Astrophysics with Machine Learning Simulations)

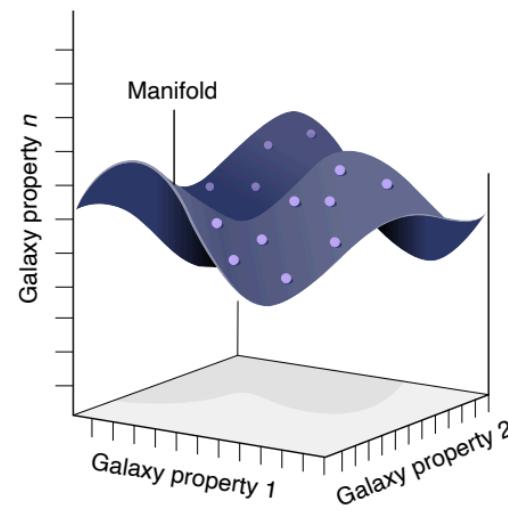
CAMELS

Villaescusa-Navarro+22

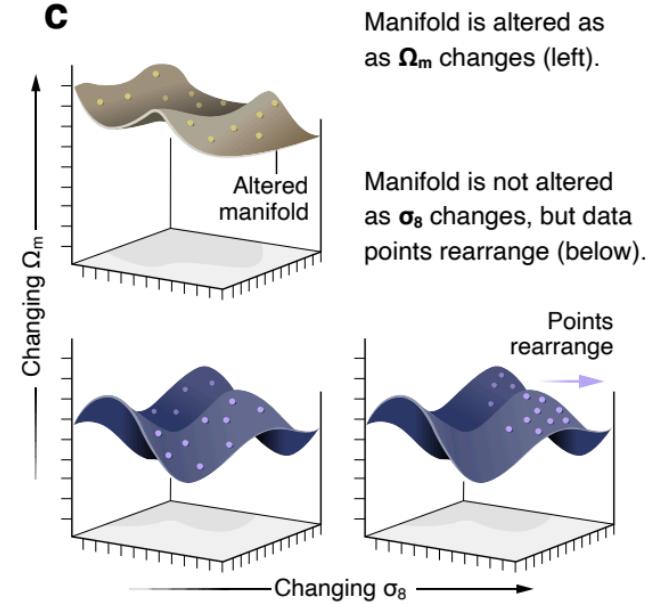
a Plotting galaxies in two dimensions



b Plotting galaxies in n dimensions



c



Lue, Genel, MHC+25