



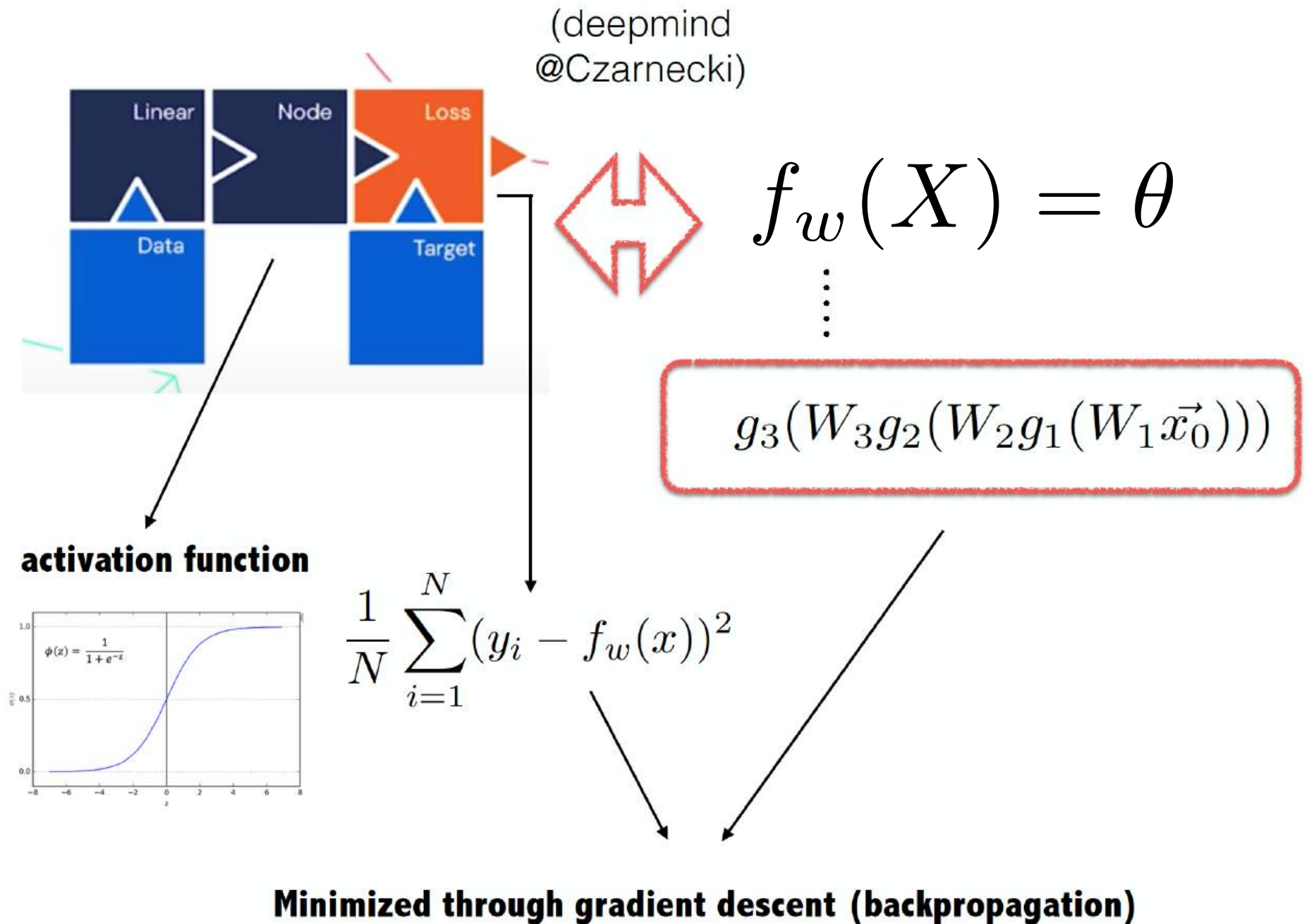
Generative Modelling, Density Estimation and Simulation Based Inference

A. Boucaud, M. Huertas-Company

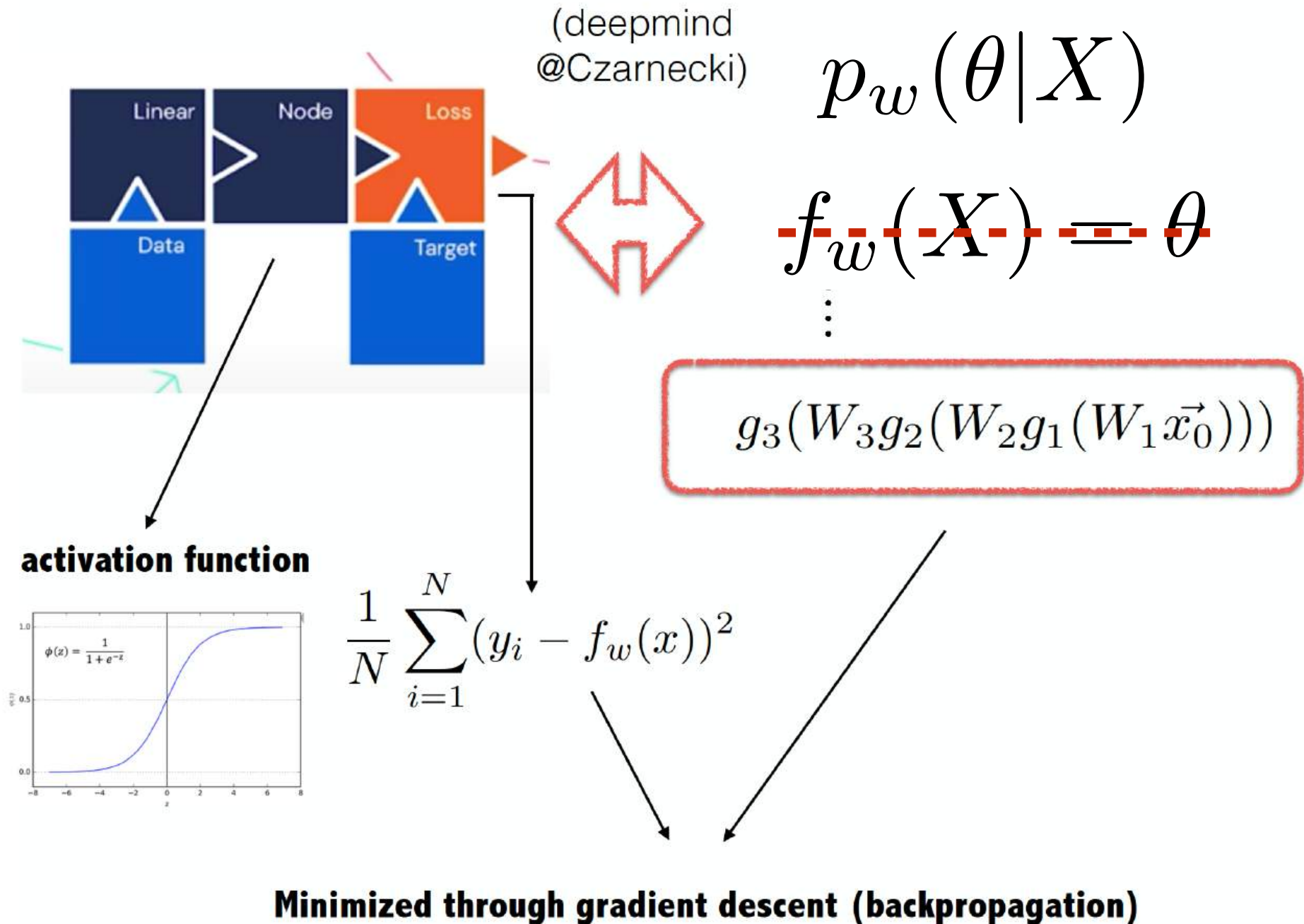


Ecole d'été Rodolphe Clédassou - Session 2024 - Cycle3

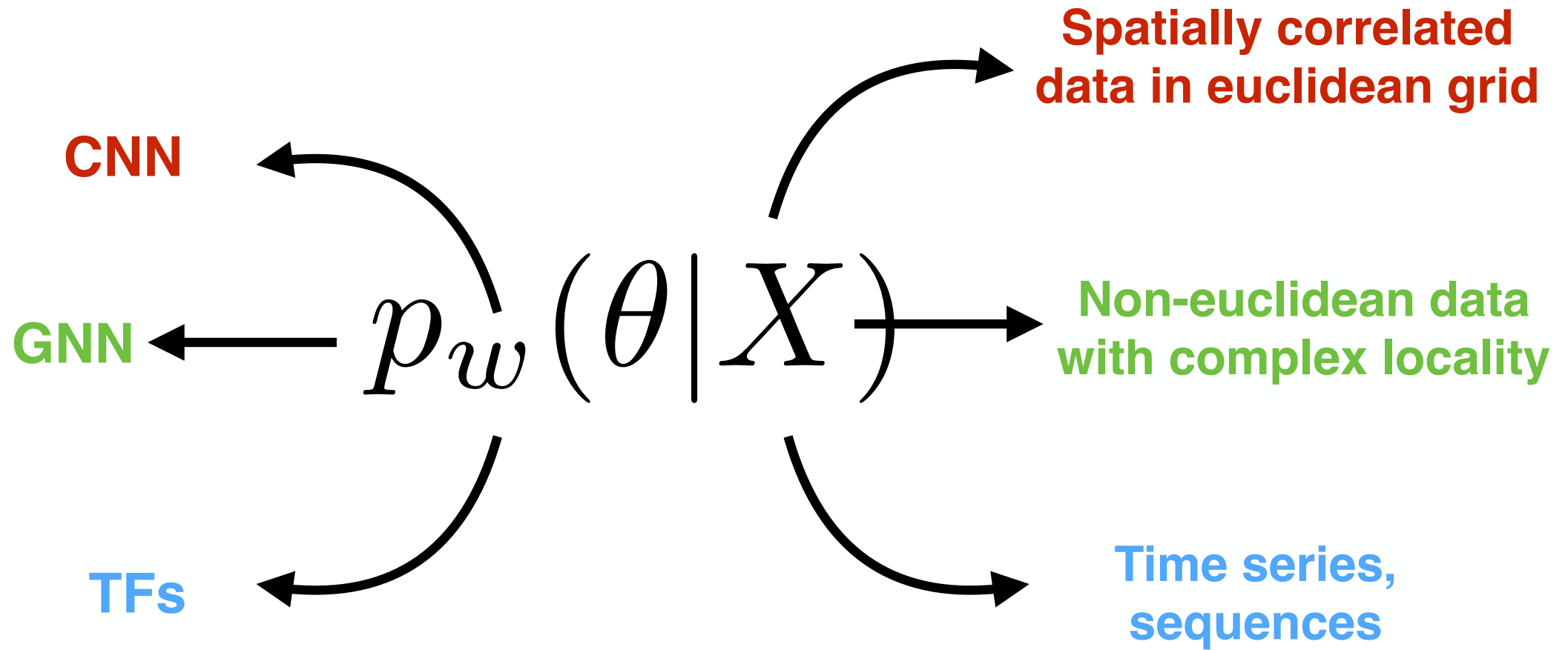
RECAP Cycle 1: NNs as universal approximators



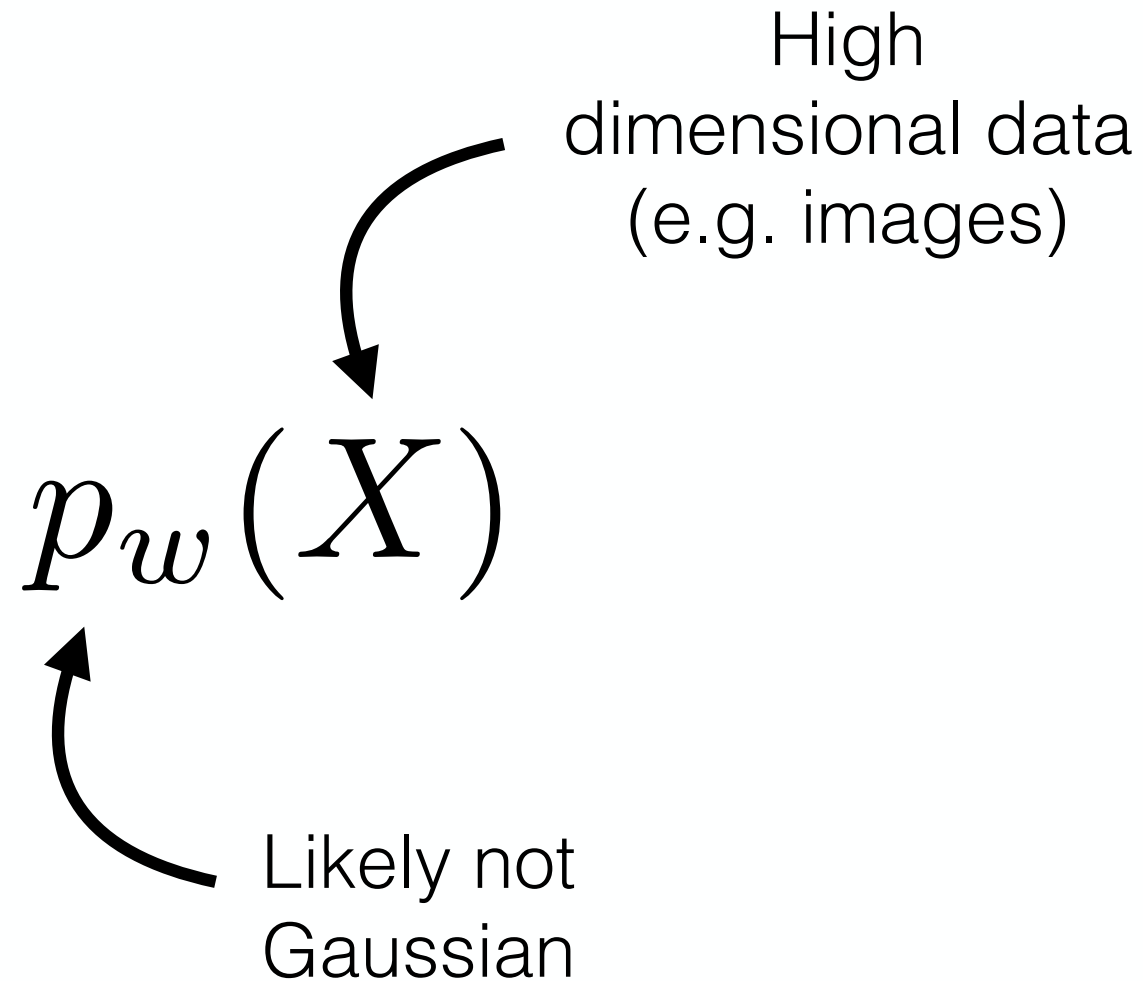
RECAP Cycle 1: NNs as universal approximators



RECAP Cycle 2: supervised deep learning



Estimate a probability density function of an arbitrarily high dimensional data



Estimate a probability density function of an arbitrarily high dimensional data

High
dimensional data
(e.g. images)

$p_w(X)$

Likely not
Gaussian

Why?

1. Sampling. how can I interpolate / extrapolate a (small, sparse) dataset to generate new data sampled from the same (unknown) distribution?

2. Likelihood evaluation. what is the probability that a new observation is drawn from the same (unknown) distribution as some reference (small, sparse) dataset?

Evaluation

High
dimensional data
(e.g. images)

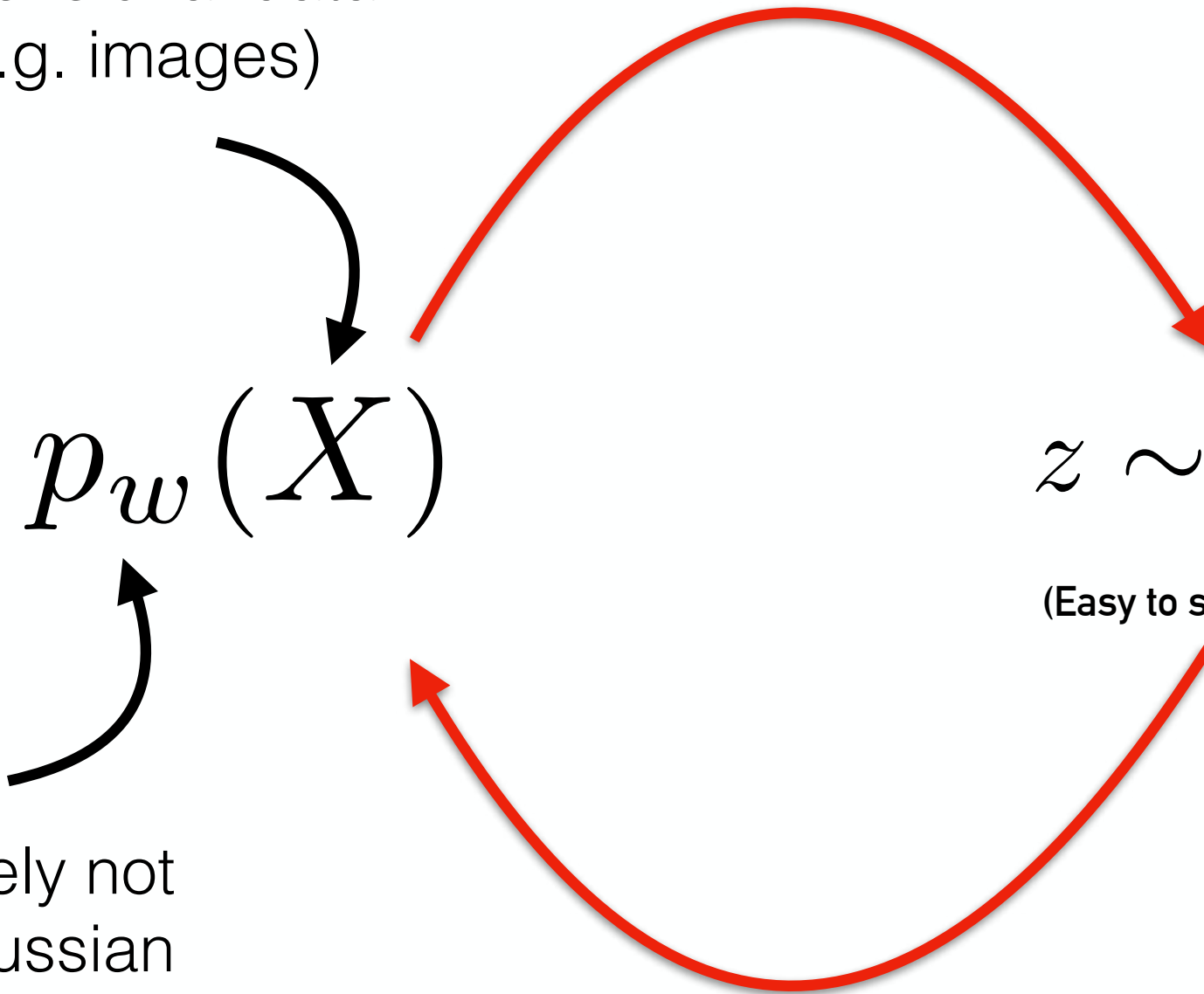
$$p_w(X)$$

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

(Easy to sample, easy to evaluate)

Likely not
Gaussian

Sampling

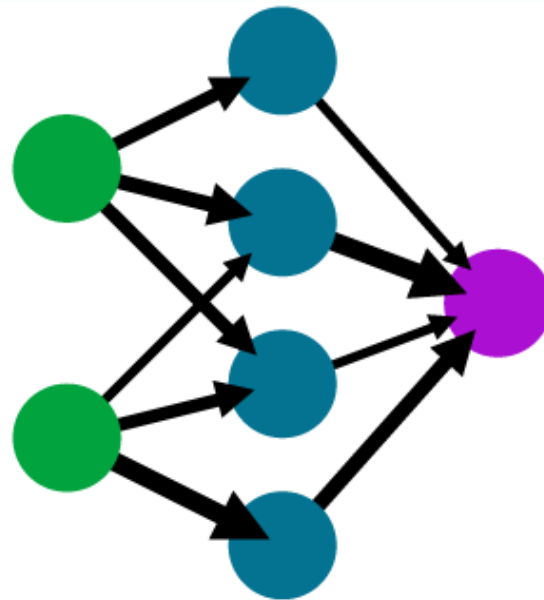


Evaluation

High dimensional data
(e.g. images)

NNs == universal
approximators

$p_w(X)$



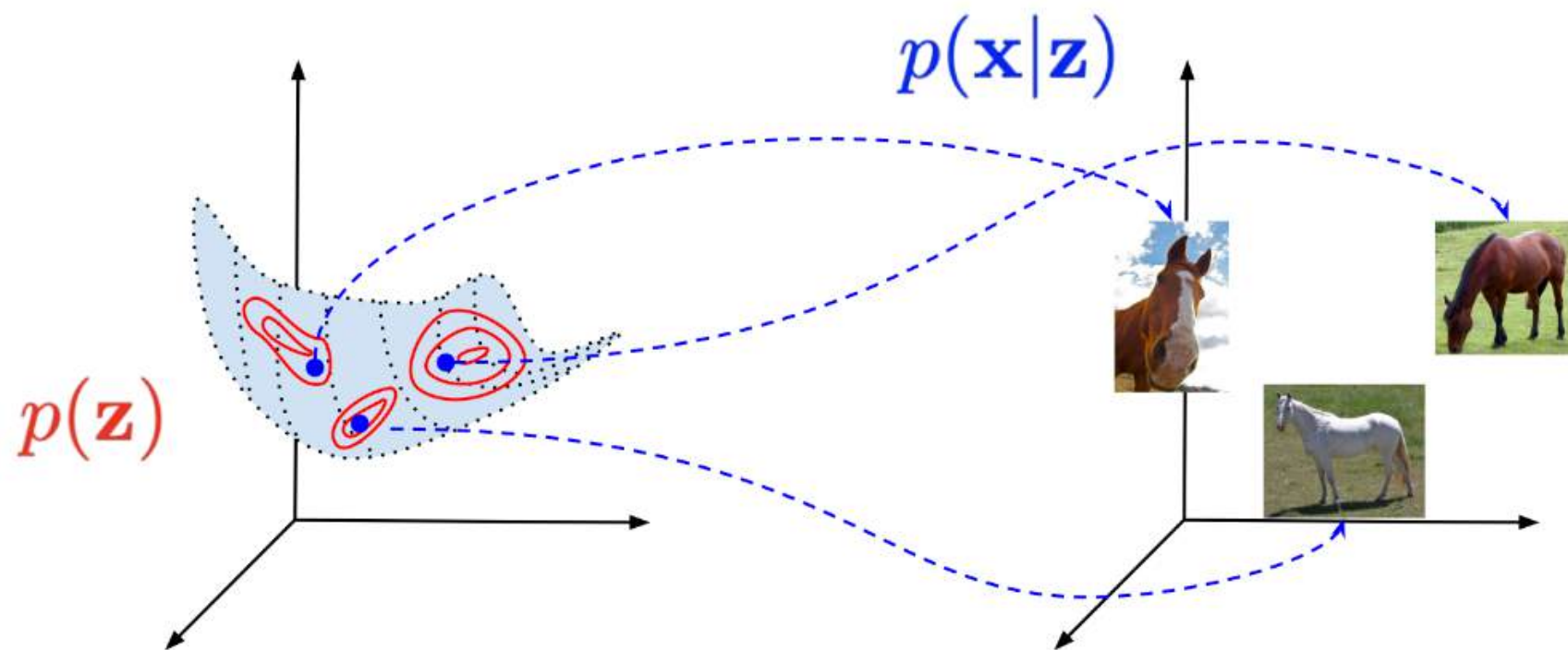
$z \sim \mathcal{N}(0, 1)$

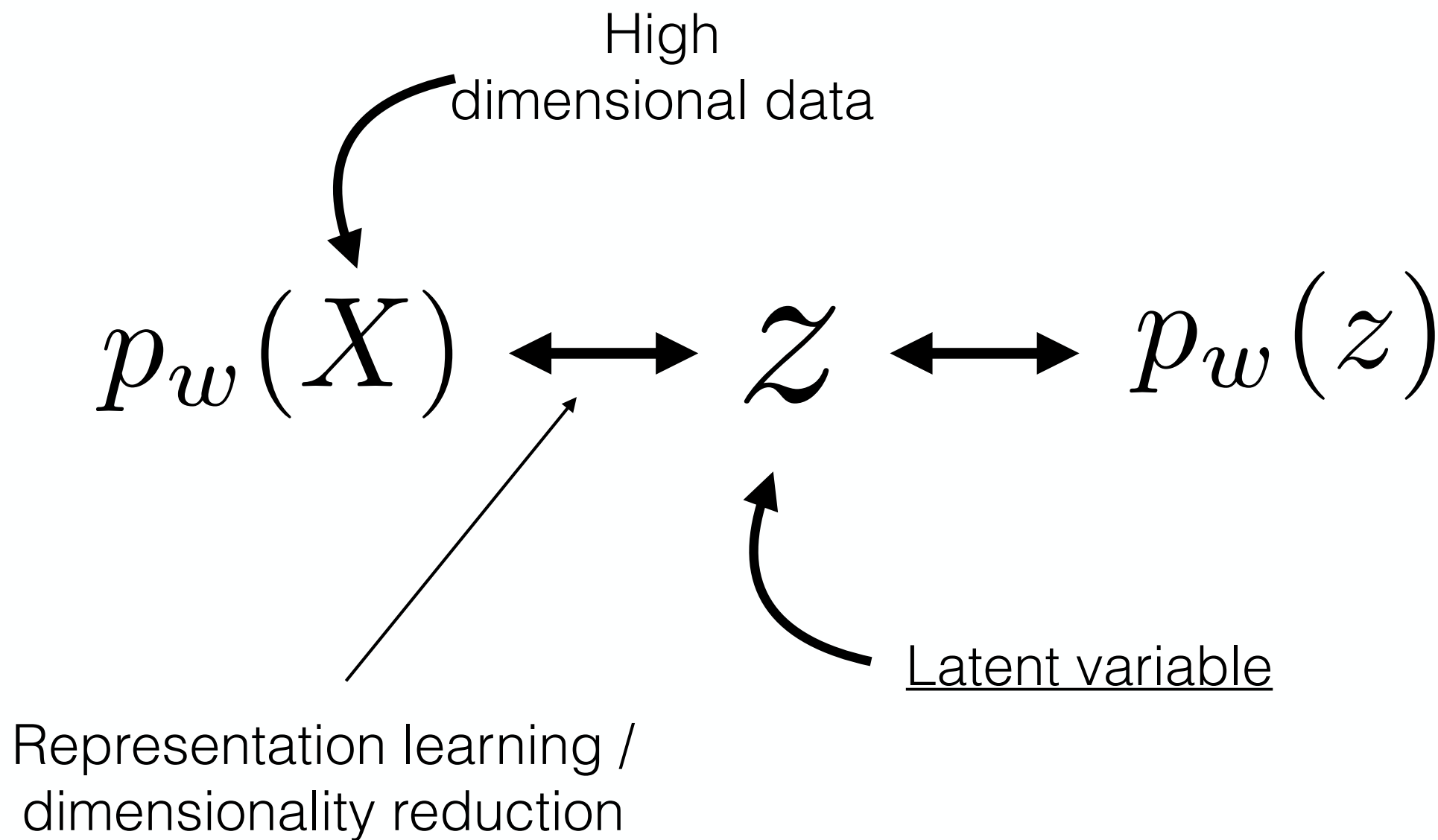
(Easy to sample, easy to evaluate)

(Optimization problem)

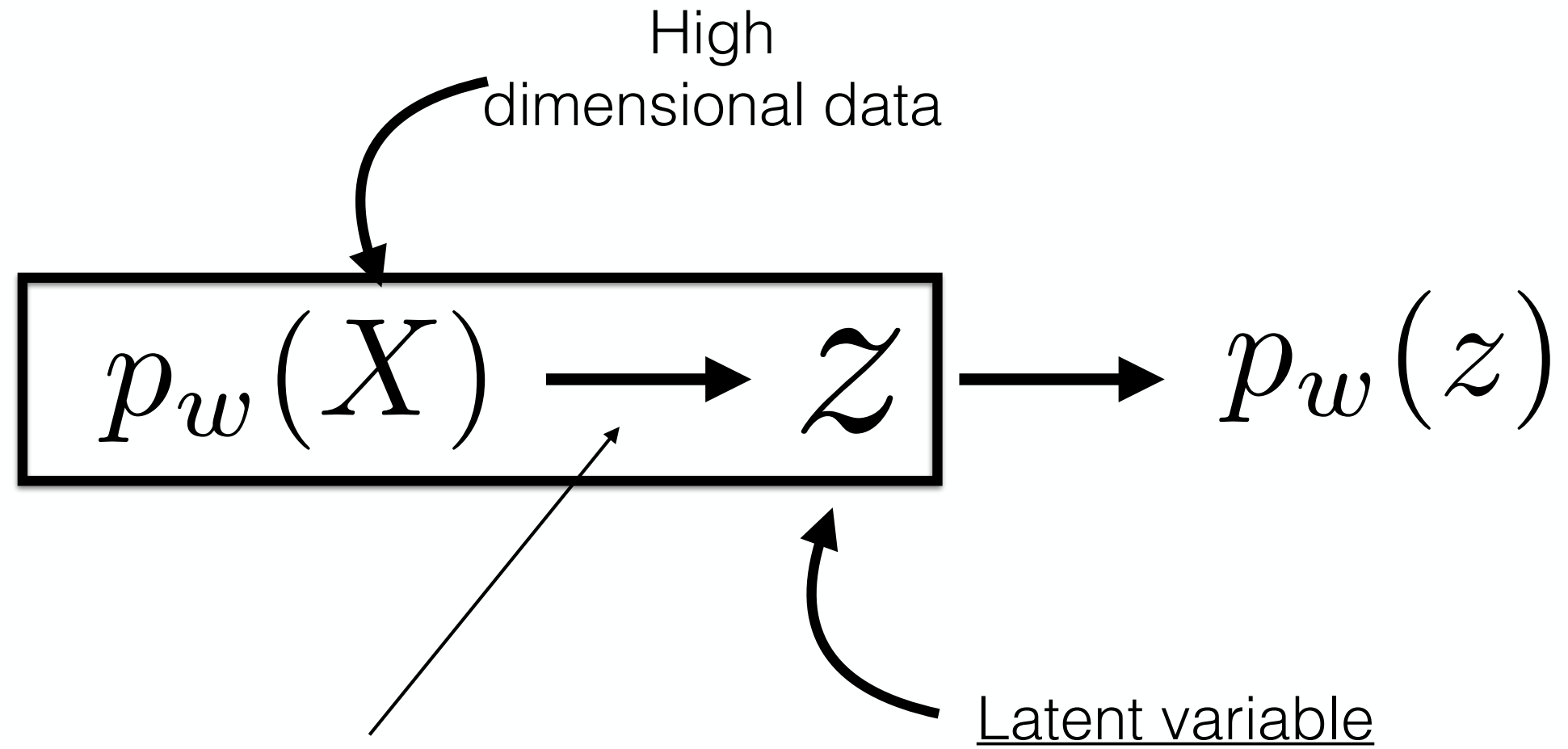
Likely not
Gaussian

Sampling





1. Representation learning

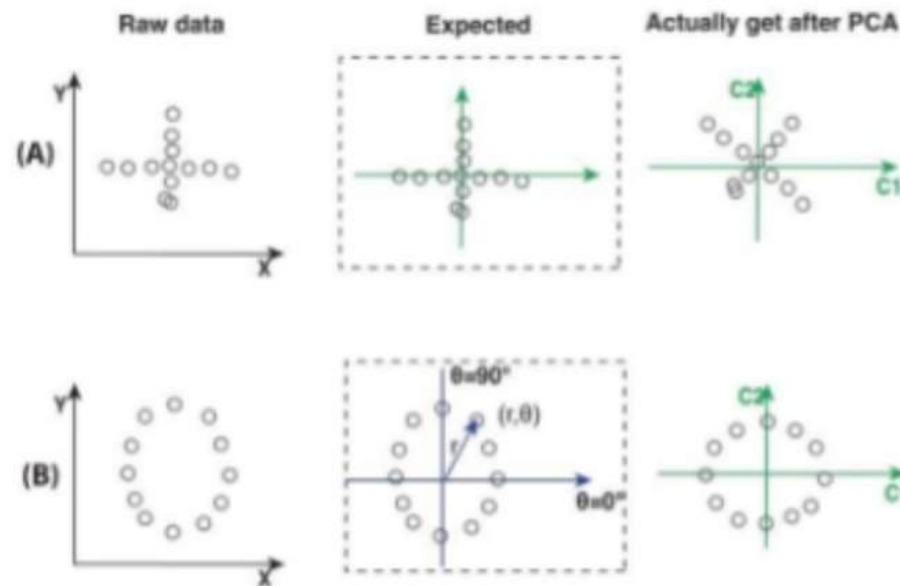


Representation learning /
dimensionality reduction

LIMITATIONS OF PCA

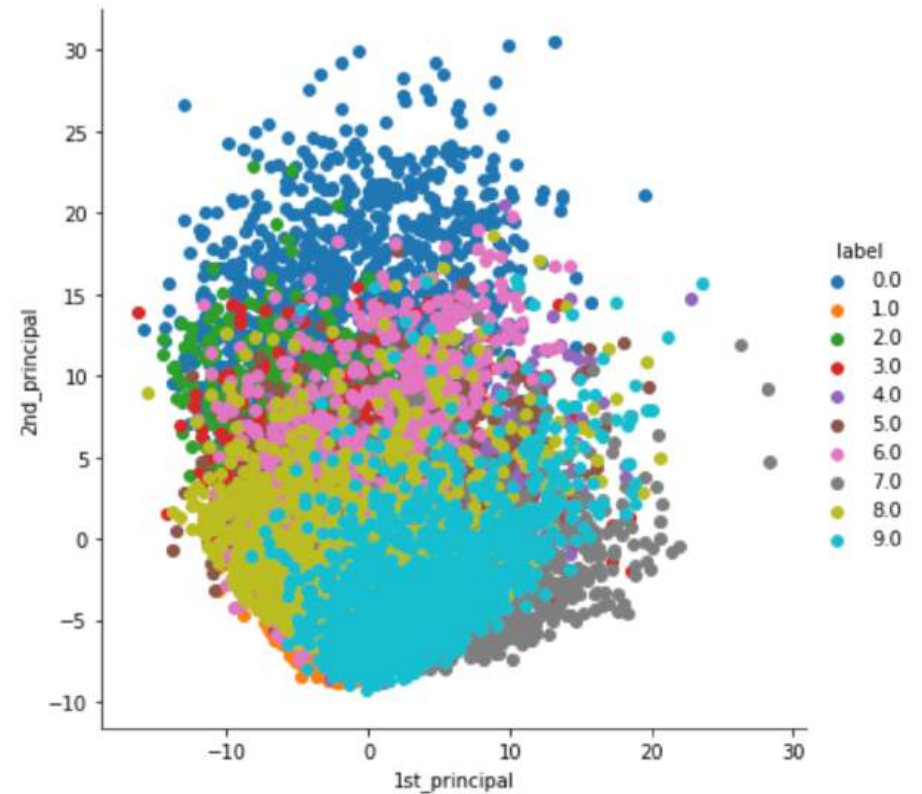
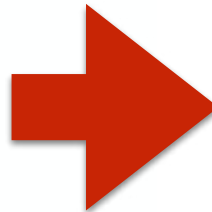
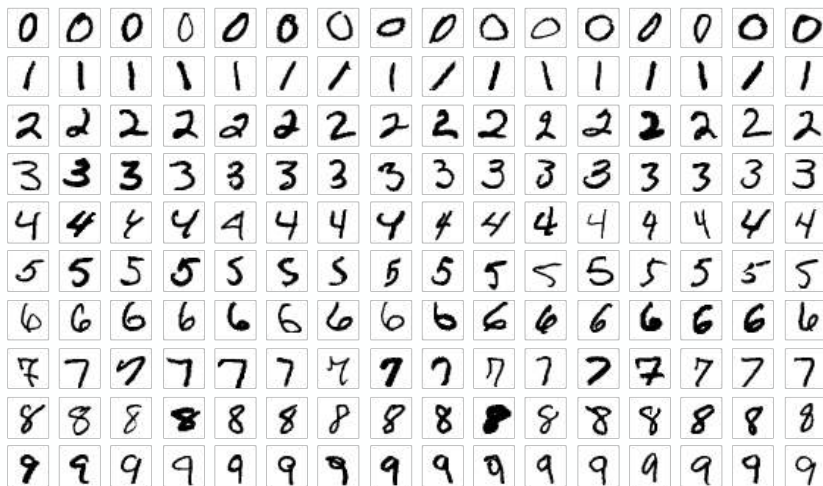
PCA APPLY LINEAR TRANSFORMATIONS

SINCE WE USE THE COVARIANCE MATRIX, IT ASSUMES
THAT THE DATA FOLLOWS A **MULTIDIMENSIONAL
NORMAL DISTRIBUTION**



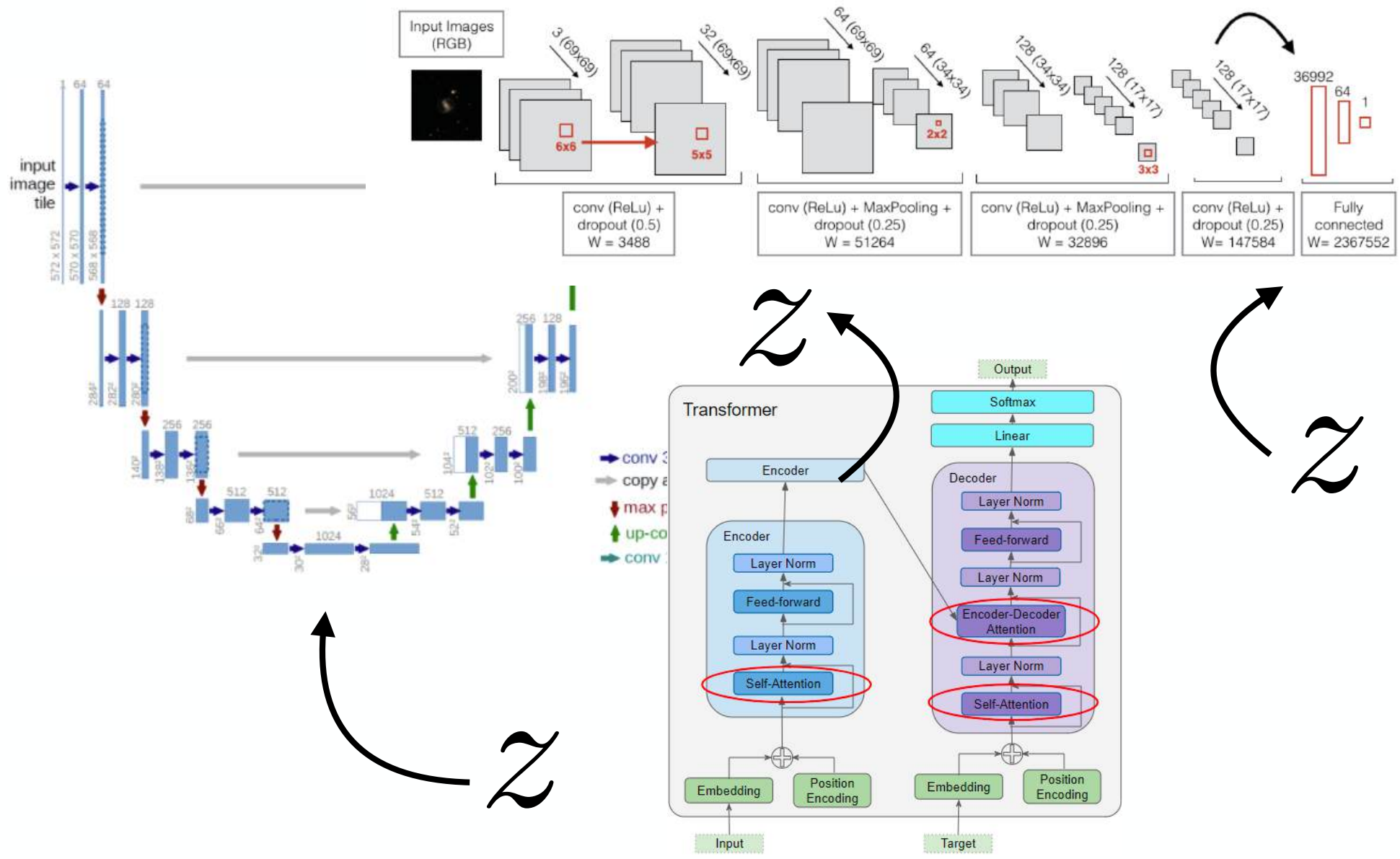
LIMITATIONS OF PCA

AND DATA IS NOT ALWAYS GAUSSIAN....

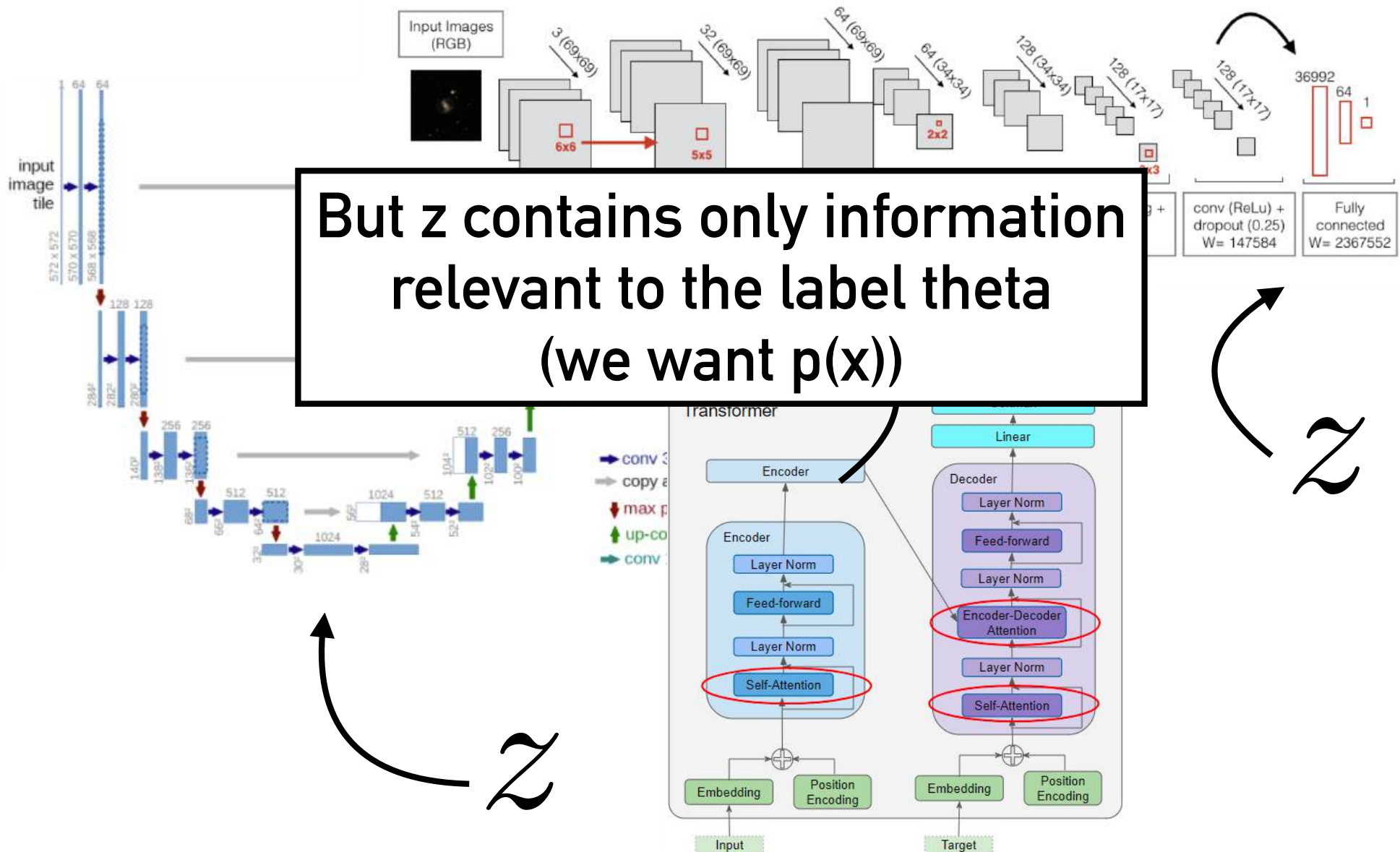


CAN WE GENERALIZE THAT?

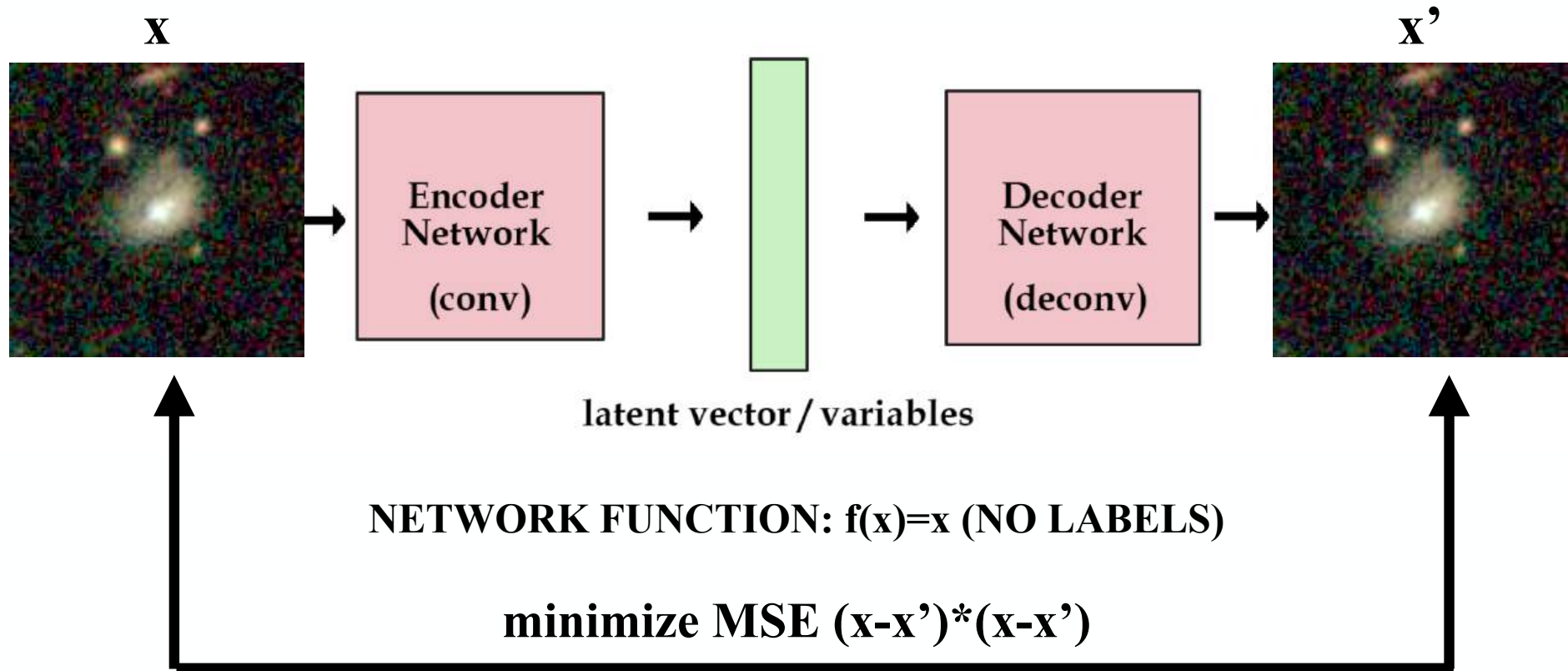
Any supervised NN obtains a non linear representation “z”



Any supervised NN obtains a non linear representation “z”

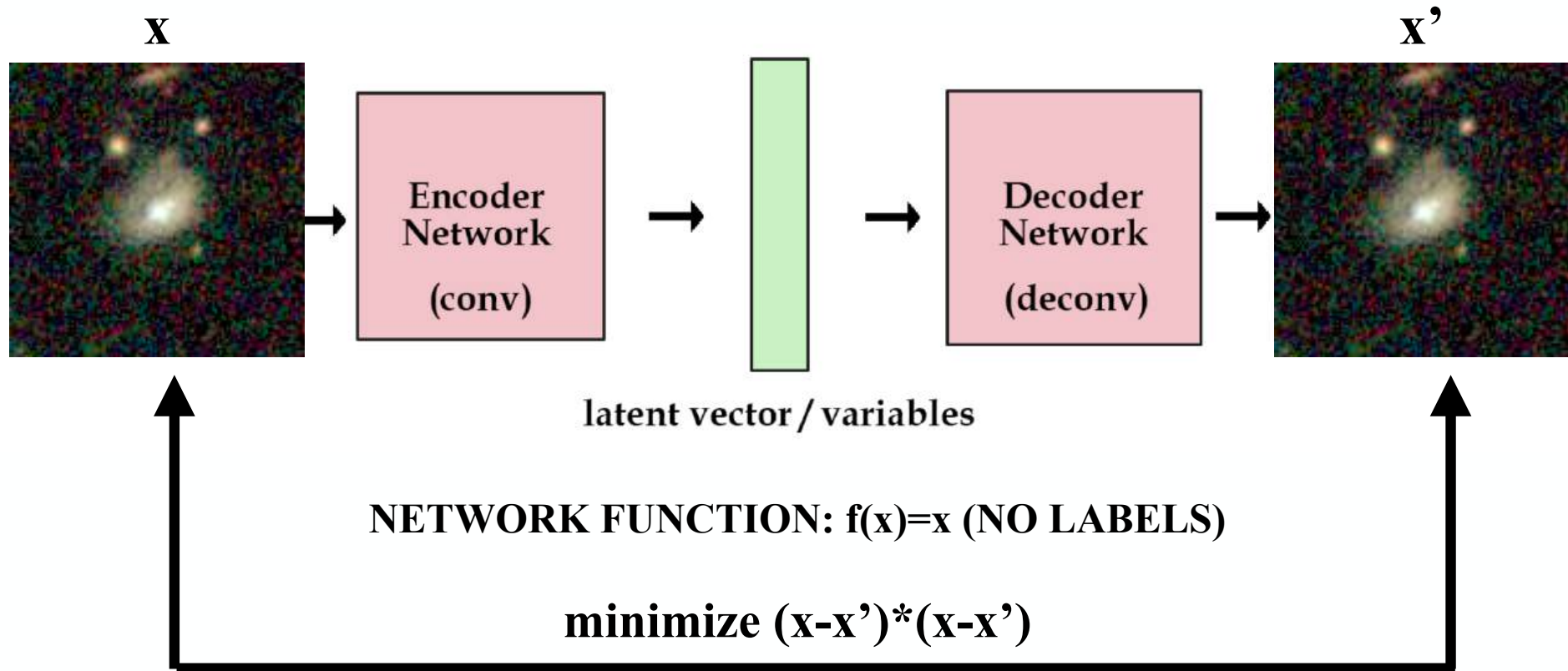


(CONVOLUTIONAL) AUTO-ENCODER



AN AUTO-ENCODER IS SIMPLY ANY NETWORK WITH IDENTICAL INPUT AND OUTPUT

CONVOLUTIONAL AUTO-ENCODER

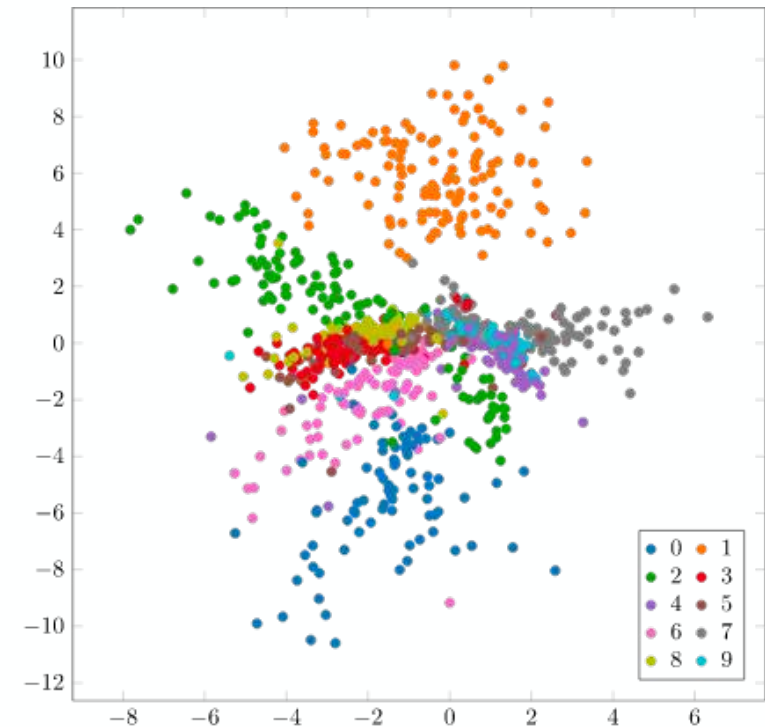
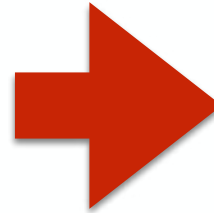
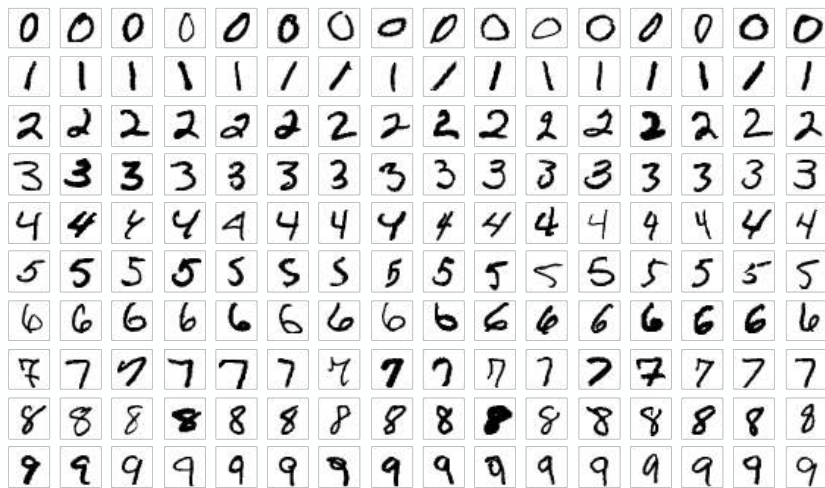


BY REDUCING THE DIMENSIONALITY IN THE LATENT SPACE
WE FORCE THE NETWORK TO LEARN A REPRESENTATION
OF THE INPUT DATA IN A LOWER DIMENSIONALITY SPACE

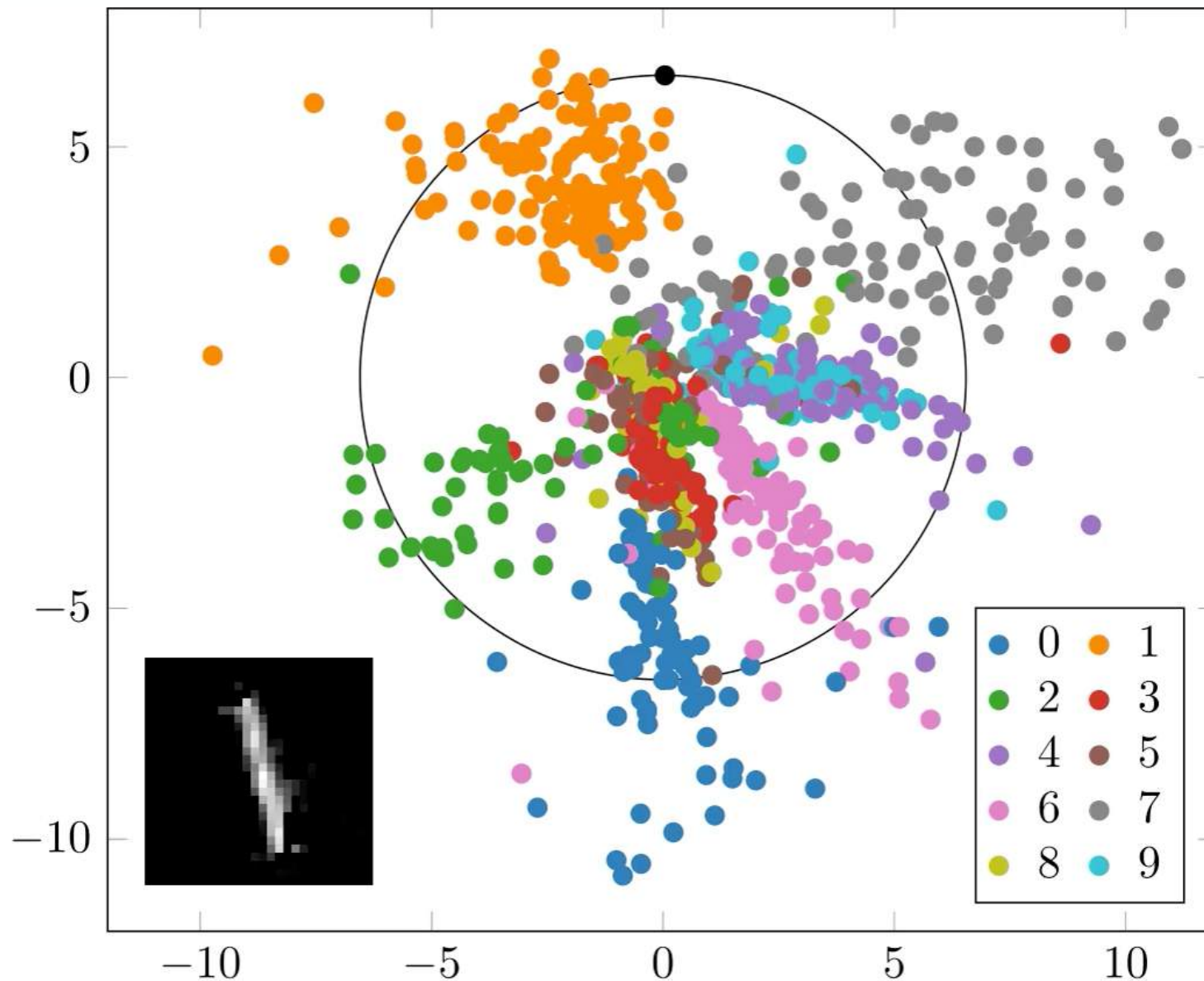
* NO NEED TO BE CONVOLUTIONAL - ANY NEURAL NETWORK WITH A BOTTLENECK WILL DO THE JOB

* **QUESTION:** WHAT WOULD HAPPEN IF WE SET AN AUTOENCODER WITH NO ACTIVATION FUNCTIONS?

AUTOENCODER REPRESENTATION OF MNIST



AUTOENCODER REPRESENTATION OF MNIST



SELF-SUPERVISED LEARNING

Another alternative is to invent a target to obtain z

$$p_w(\hat{\theta} | X)$$

(Still using a discriminative model, with a general target, and use the latent z as a representation of the data)

SELF-SUPERVISED LEARNING

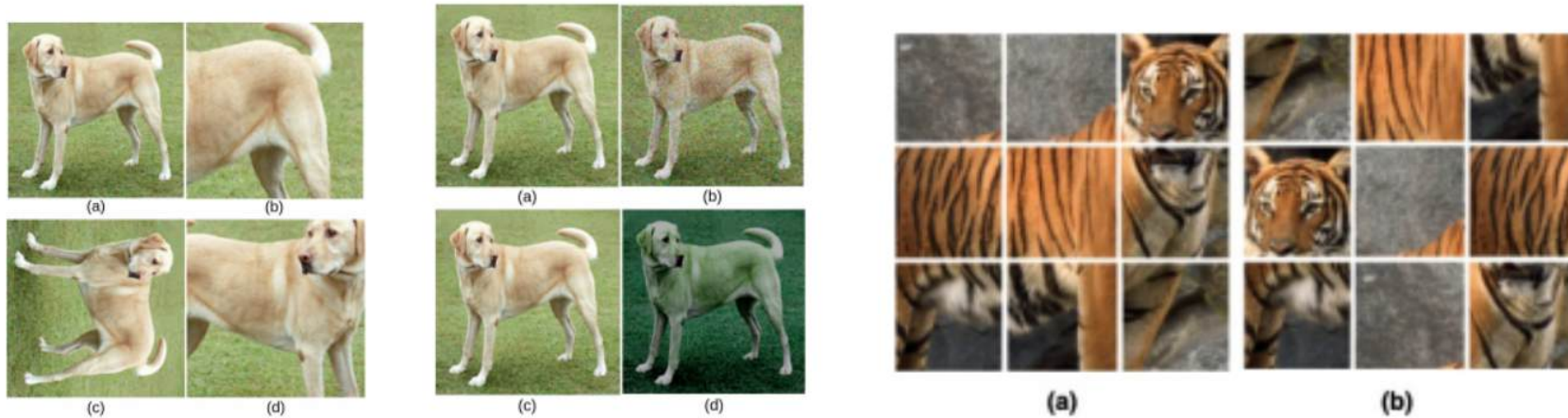
Another alternative is to invent a target to obtain z
(Pretext tasks)

$$p_w(\hat{\theta} | X)$$

 Rotation, zoom etc..

(Still using a discriminative model, with a general target, and use the latent z as a representation of the data)

e.g. contrastive learning uses the similarity between augmented version as a pretext task



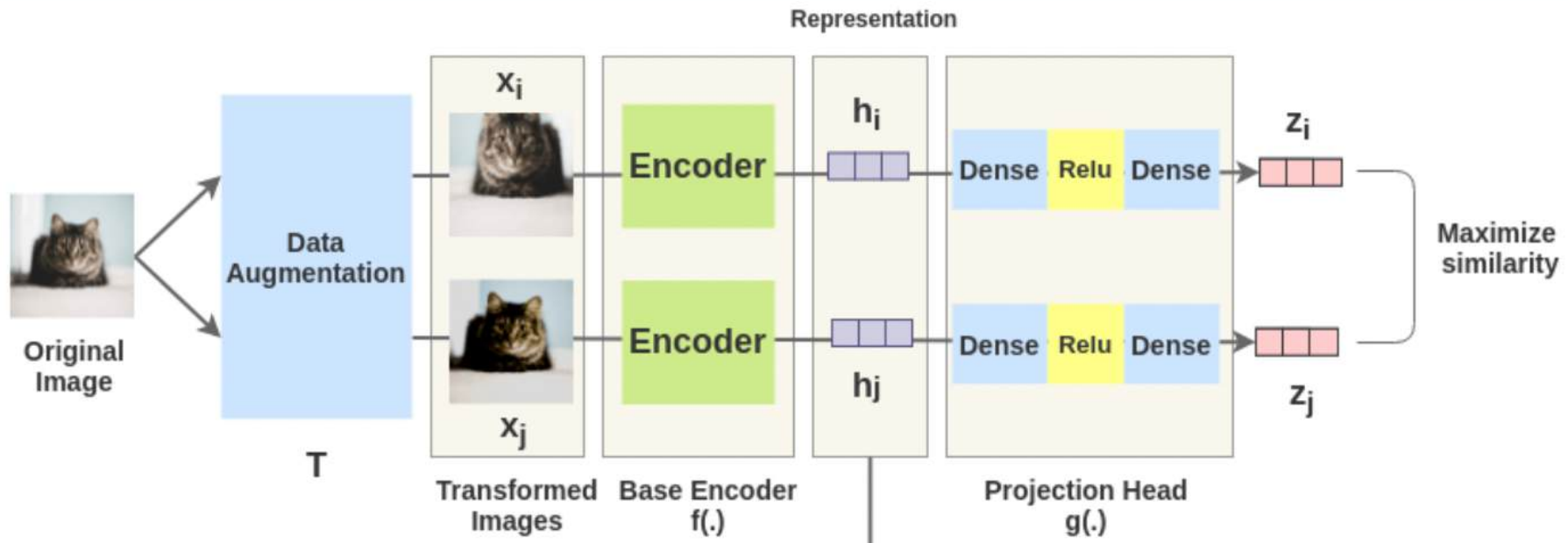
Color Augmentation

Image Rotation

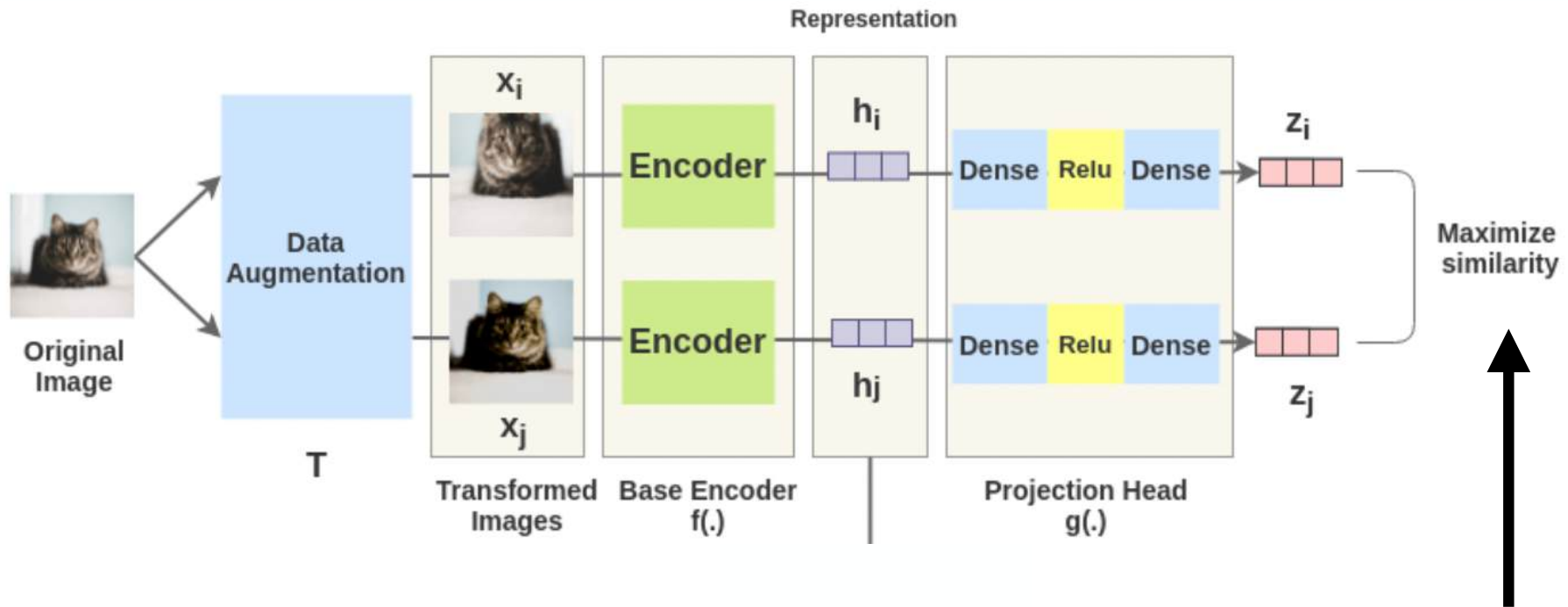
Image Cropping

Any geometrical transformation ...

The augmented versions of the images are passed through siamese networks and projected into a latent variable z



The augmented versions of the images are passed through siamese networks and projected into a latent variable z



The key is
the loss
function

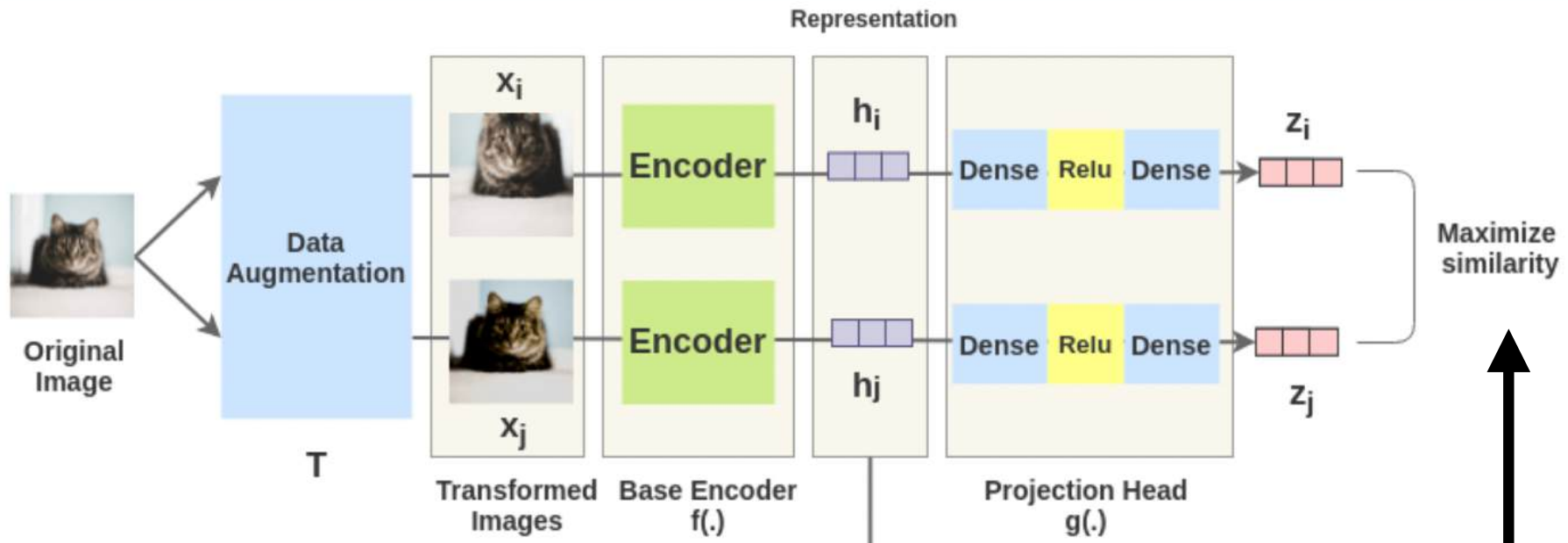
The contrastive loss:

Similarly between
two representations of positive pairs

$$l_{i,j} = -\log \frac{\exp(\langle z_i, z_j \rangle / h)}{\sum_{k=1, k \neq i}^{2N} \exp(\langle z_i, z_k \rangle / h)},$$

Sum of all similarities between
negative pairs

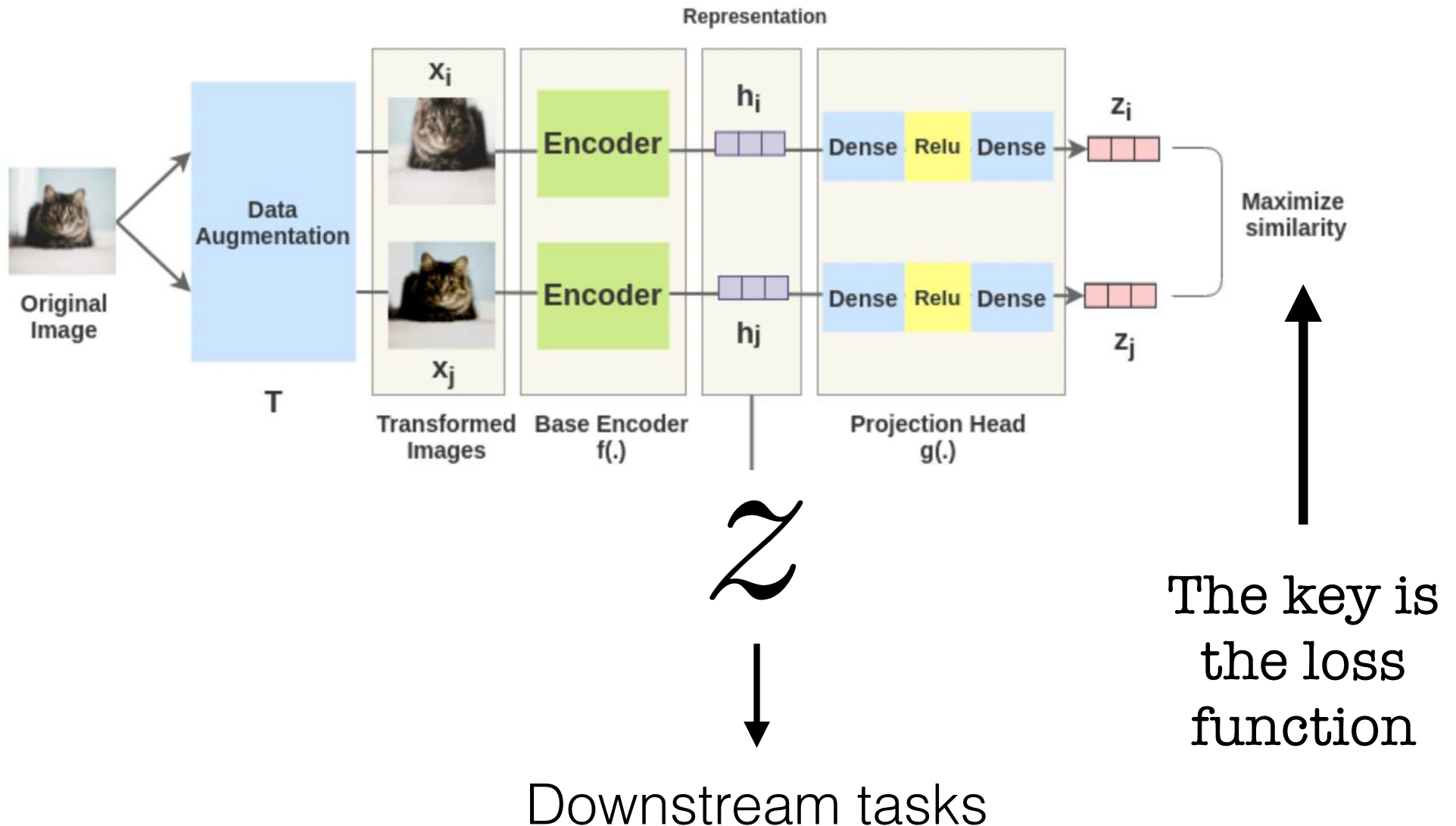
The augmented versions of the images are passed through siamese networks and projected into a latent variable z



z

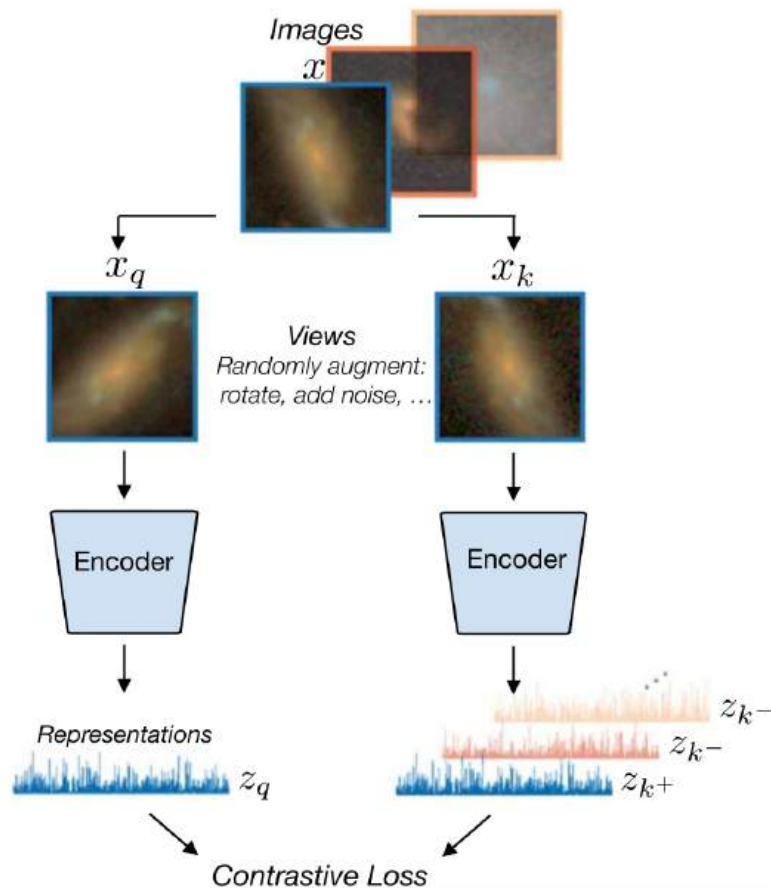
The key is
the loss
function

The augmented versions of the images are passed through siamese networks and projected into a latent variable z



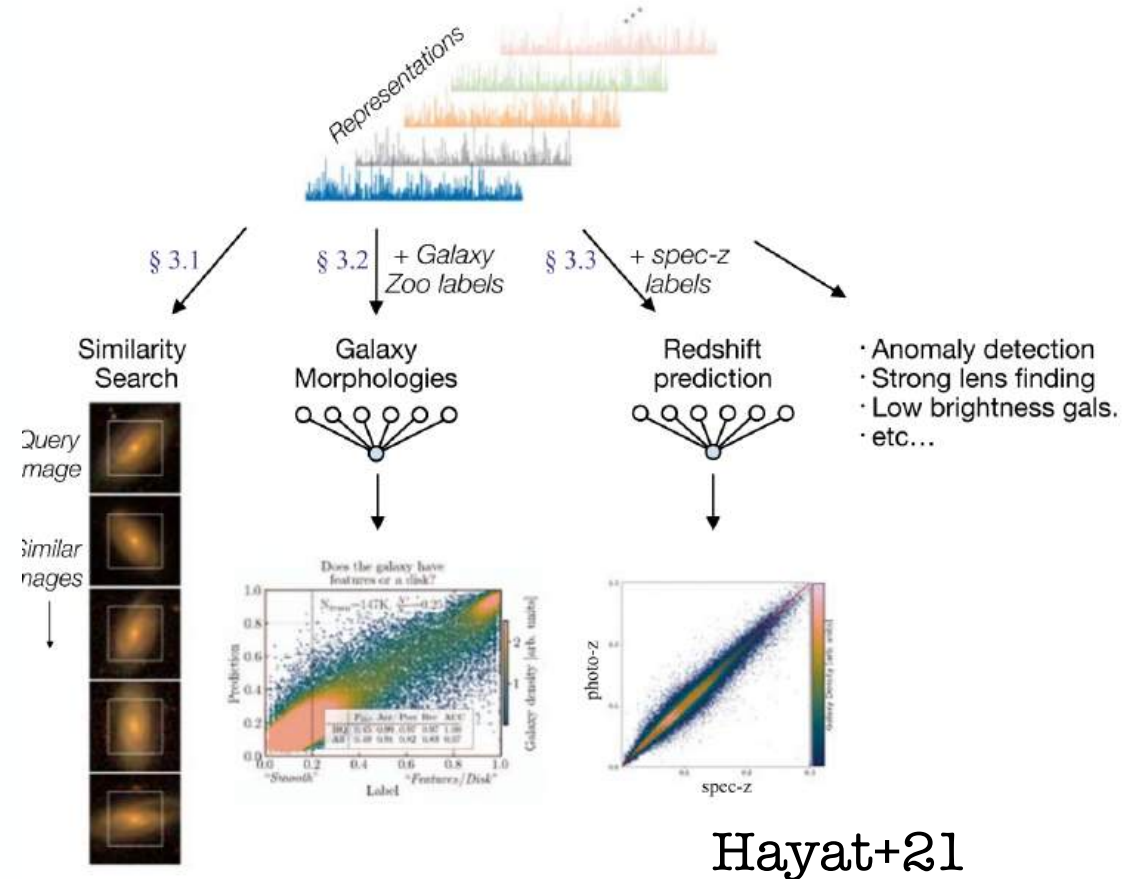
1. Self-supervised contrastive representation learning

Learn representations in an unsupervised manner



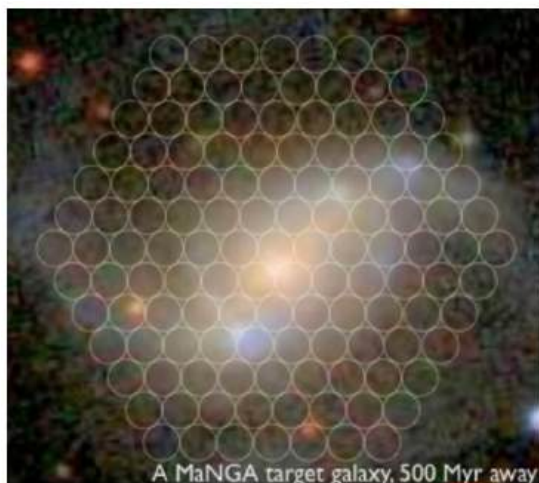
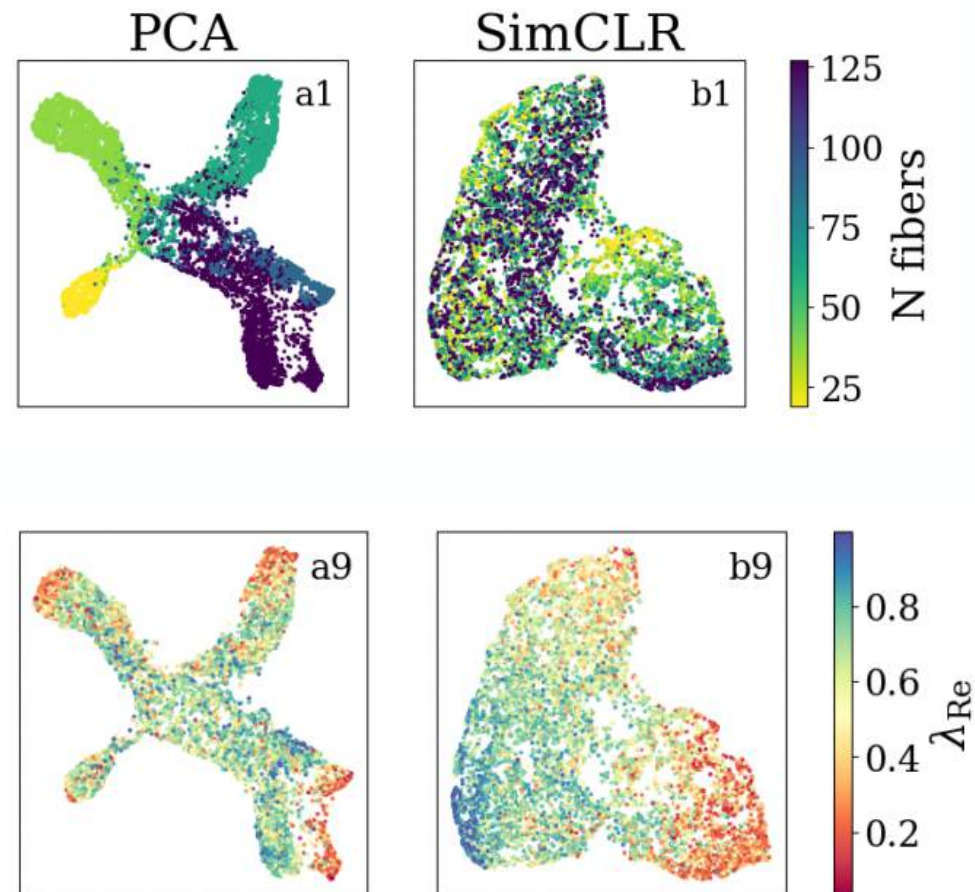
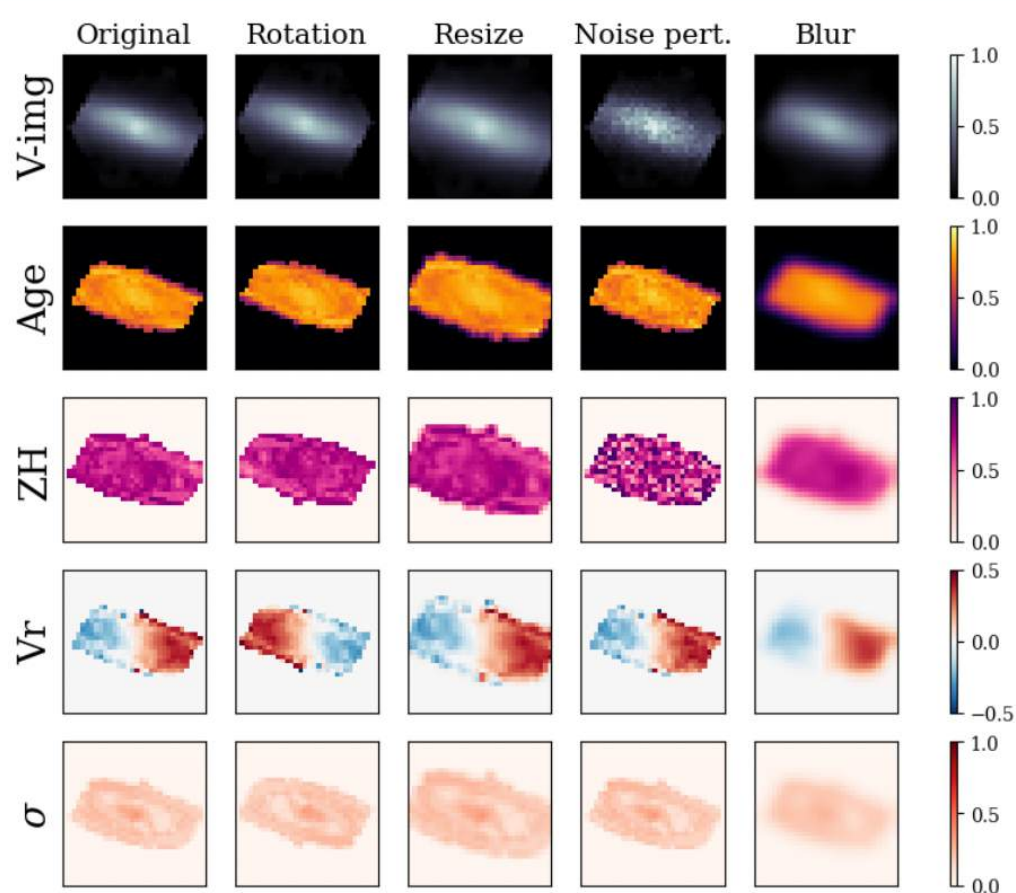
2. Downstream tasks

Use representations for a variety of applications



It turns out that the representations learned are very general and can be used for a variety of “downstream tasks”

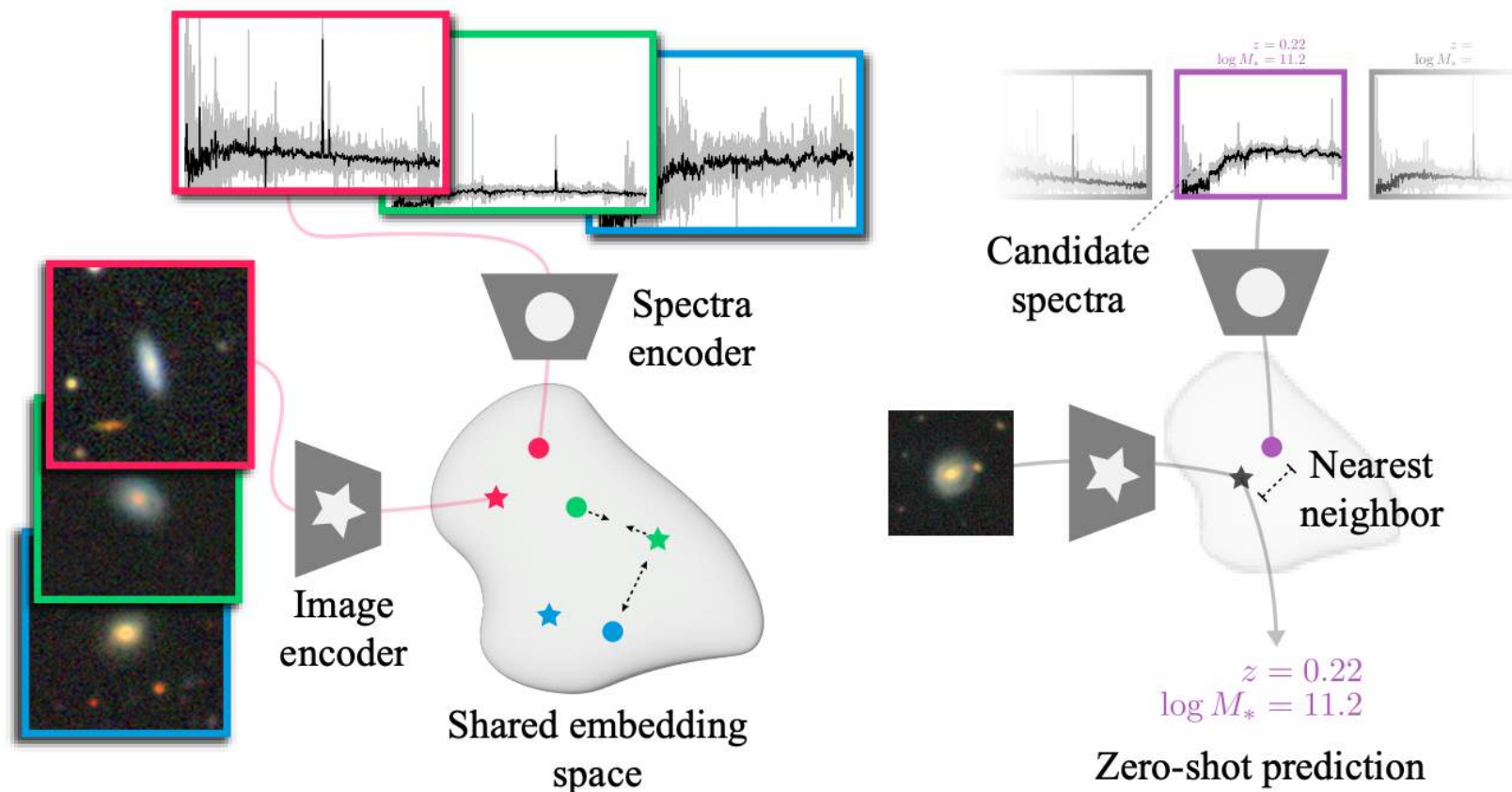
(Foundation Models) = Model not trained for a specific task, which can be easily fine tuned for a variety of downstream tasks



Contrastive learning representation of Manga galaxies

Sarmiento+21

Allows for “straightforward” multimodal representation learning

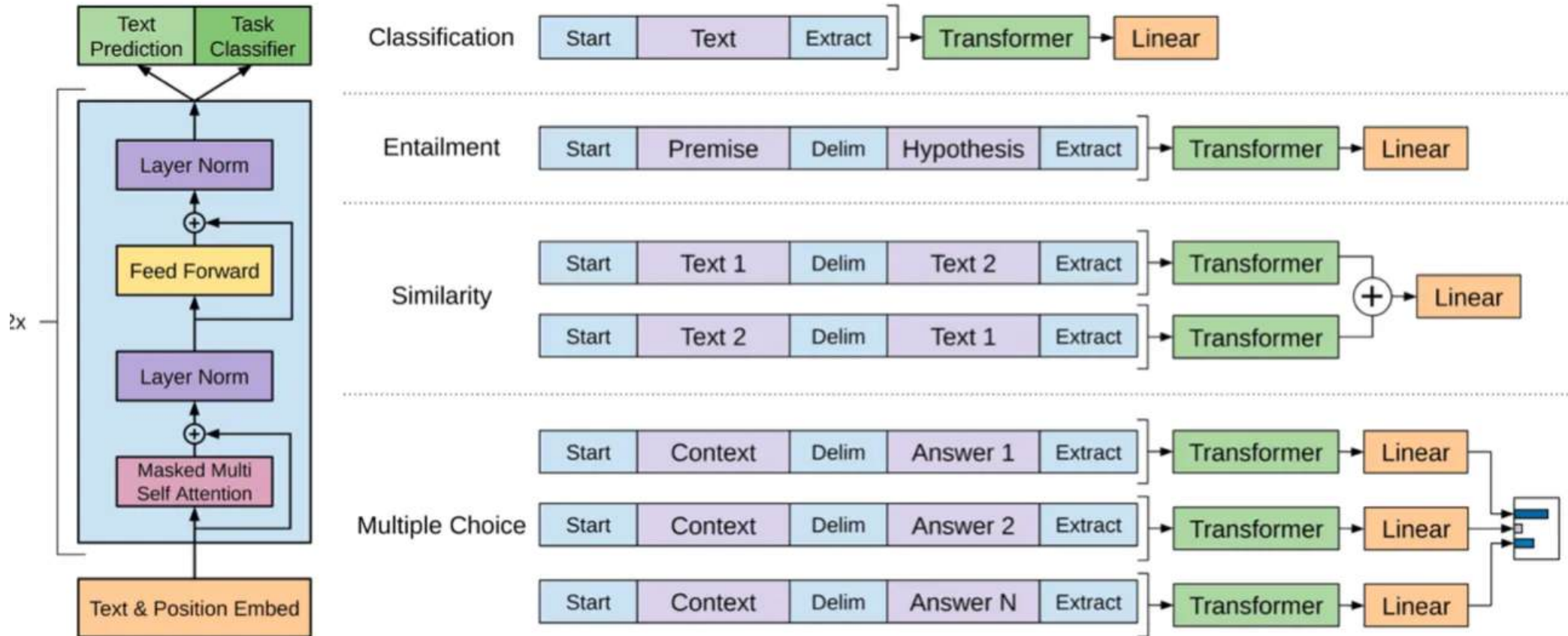


Parker+24

$$l_{i,j} = -\log \frac{\exp(\langle z_i, z_j \rangle / h)}{\sum_{k=1, k \neq i}^{2N} \exp(\langle z_i, z_k \rangle / h)},$$

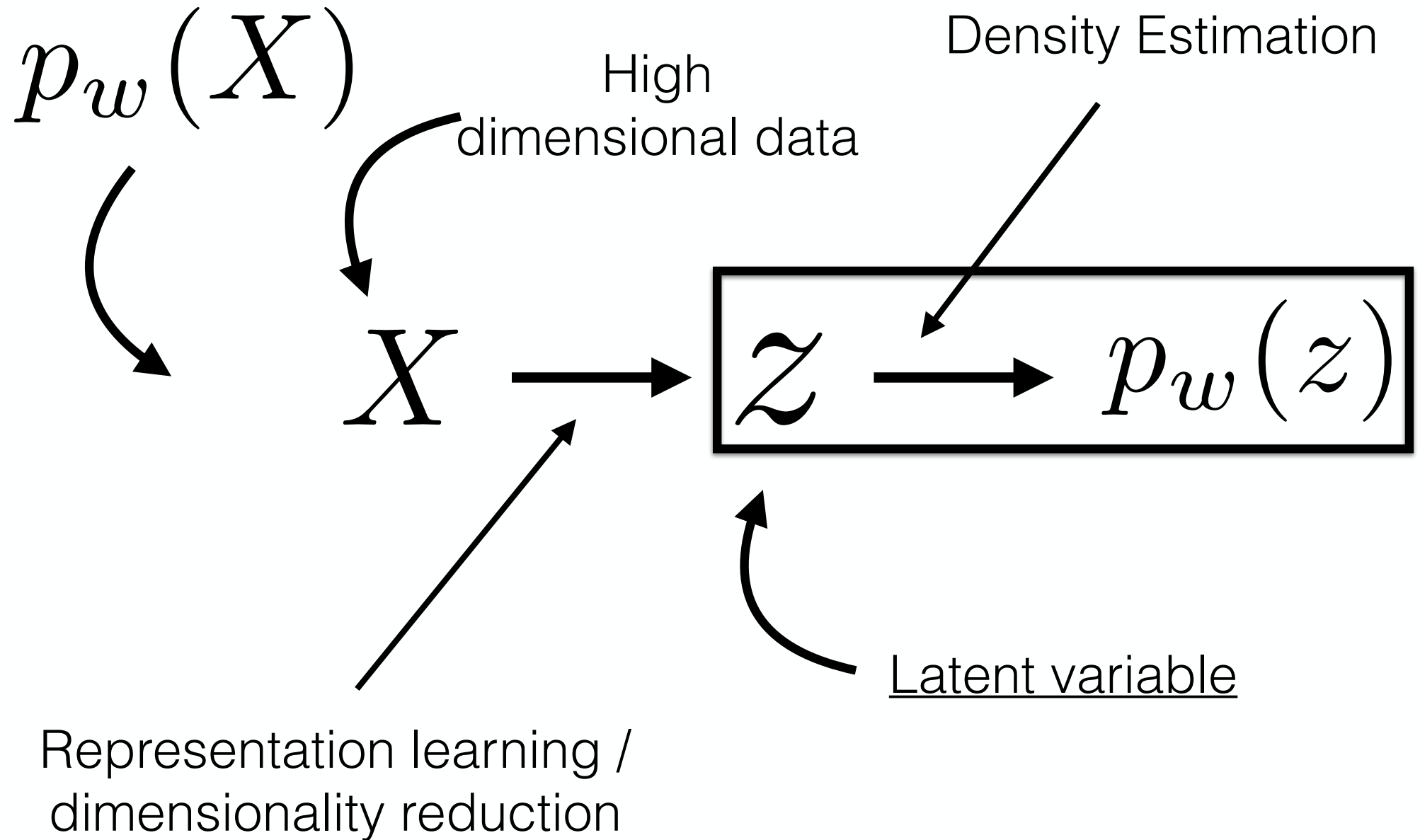
The loss function lives in the latent space

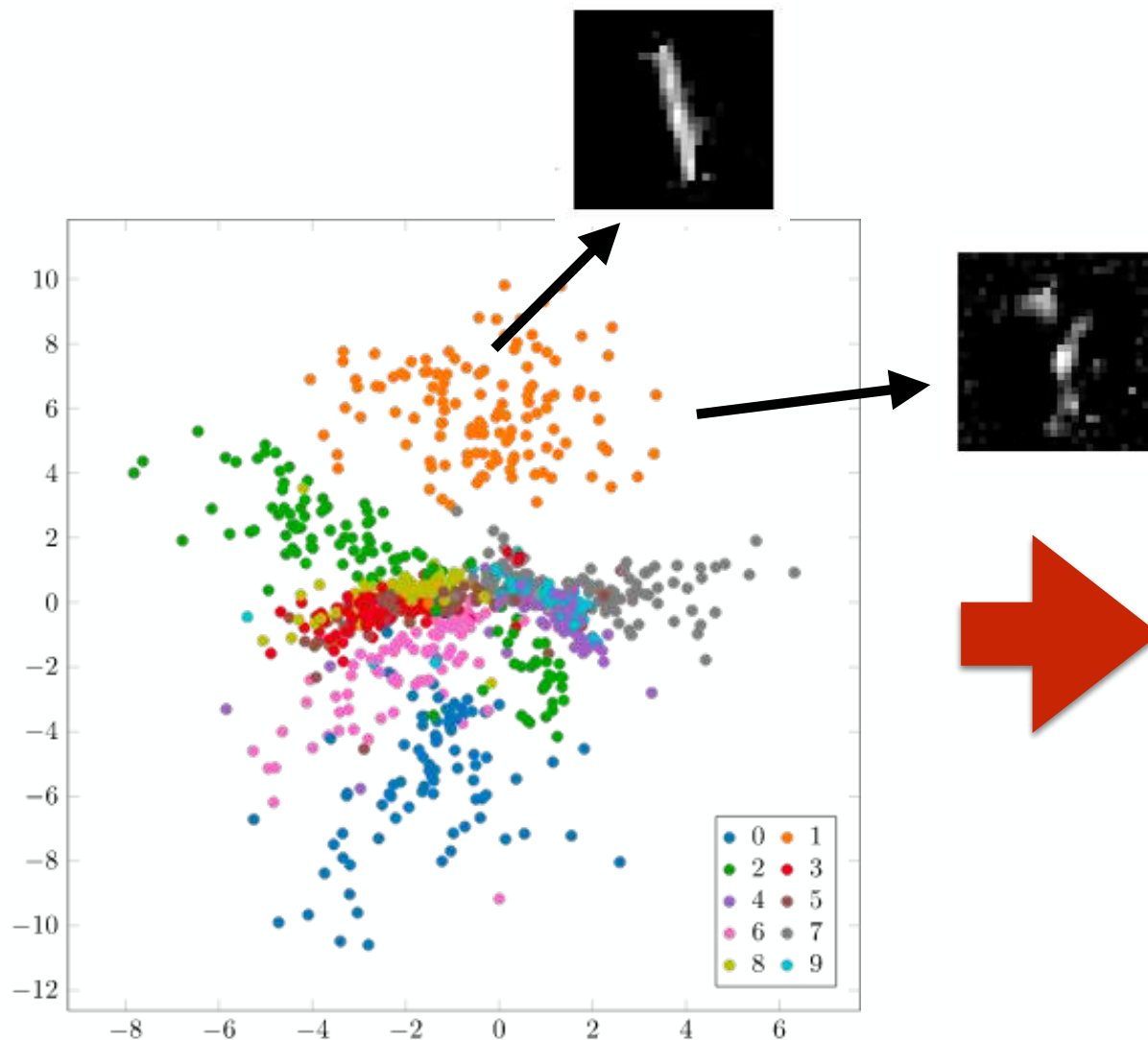
LLMs are Foundation Models



2. Sampling

RECAP5b:





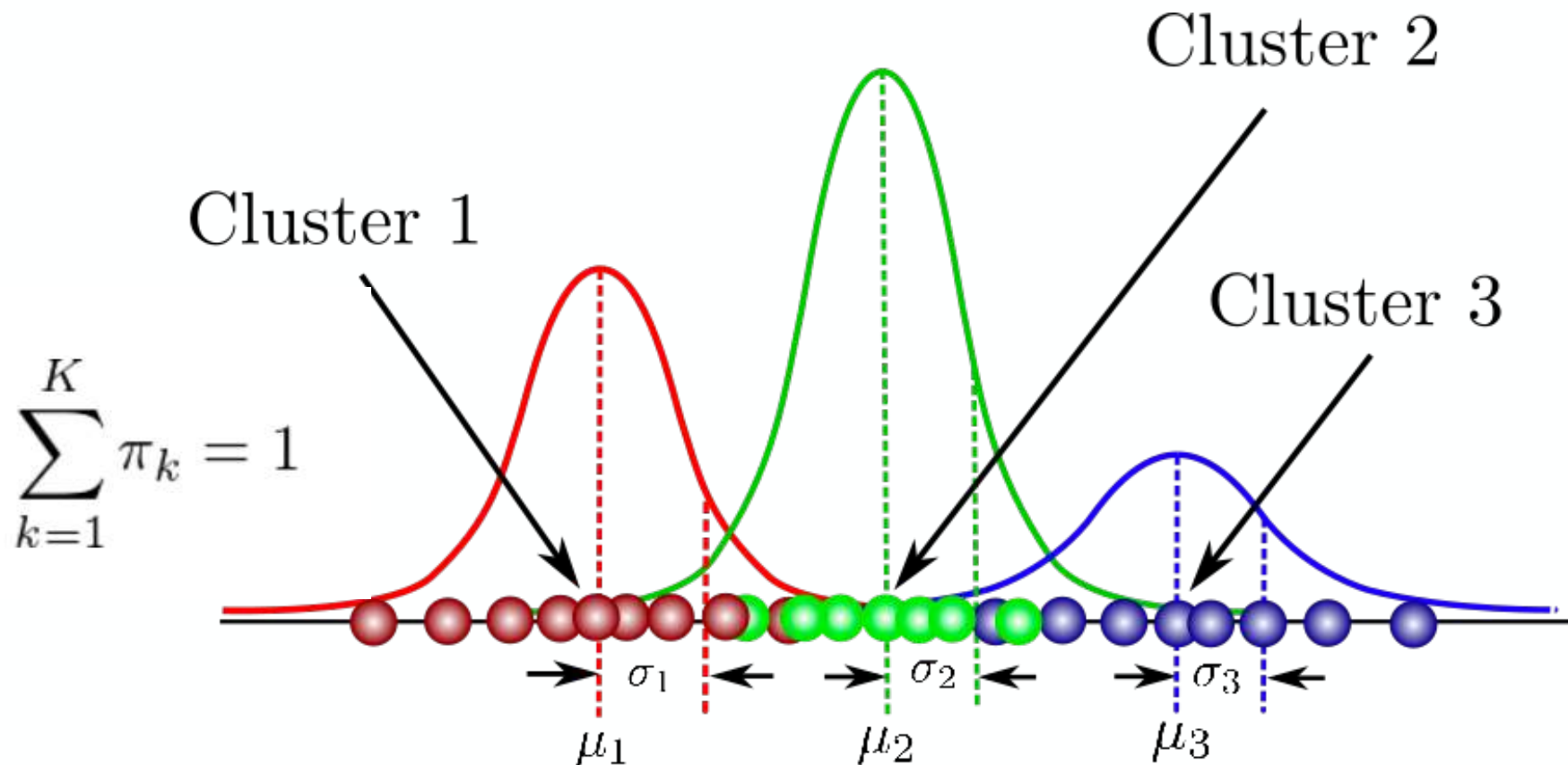
z space (latent space)

HOW CAN I
ESTIMATE $P(X)$?

$(P(z|x))$

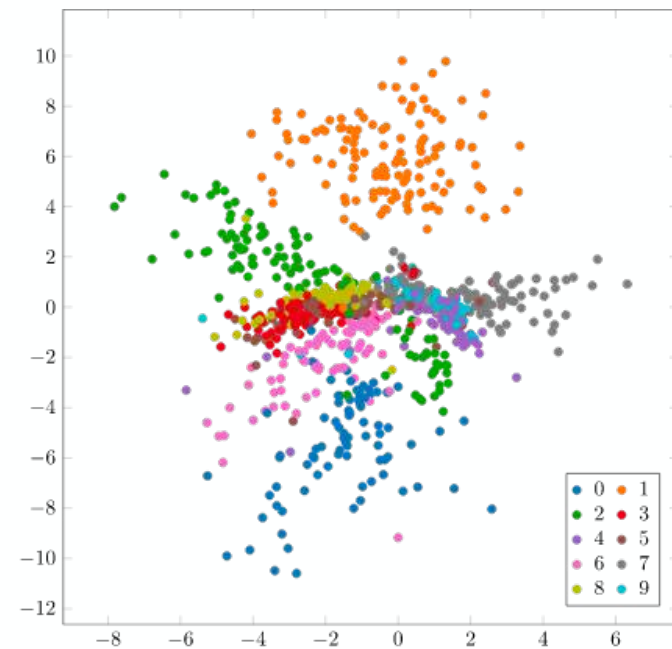
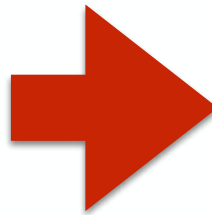
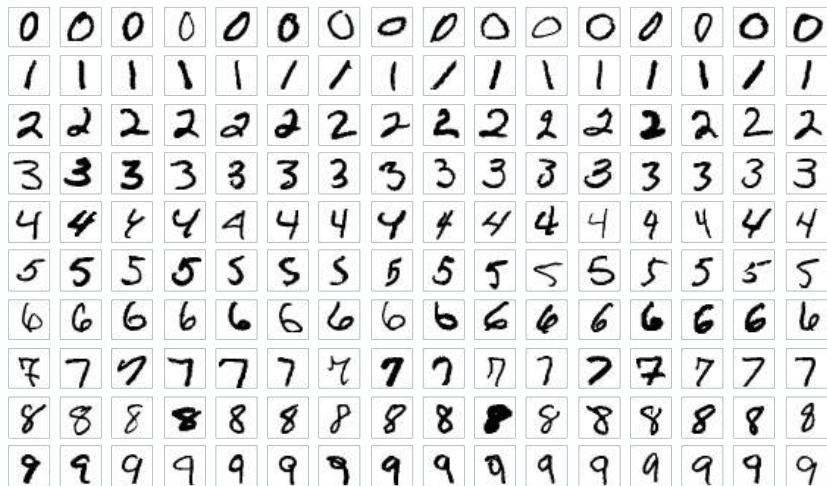
When you do not know, assume it is close to Gaussian...

GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION

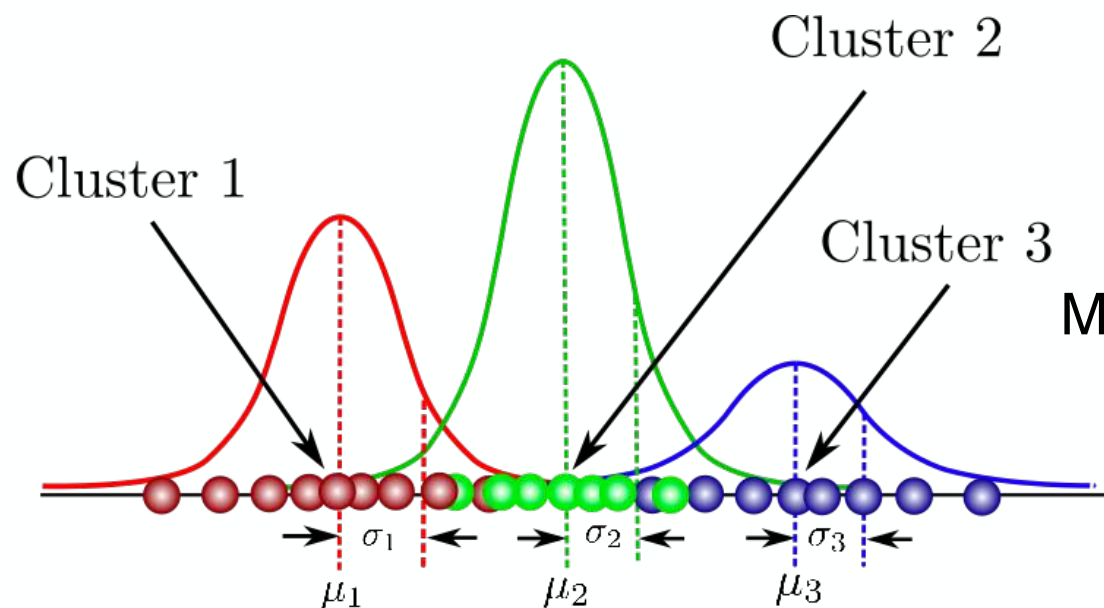


means, sigmas and scale factors of each gaussian are free parameters

DATA



REPRESENTATION
(AUTOENCODER - z)

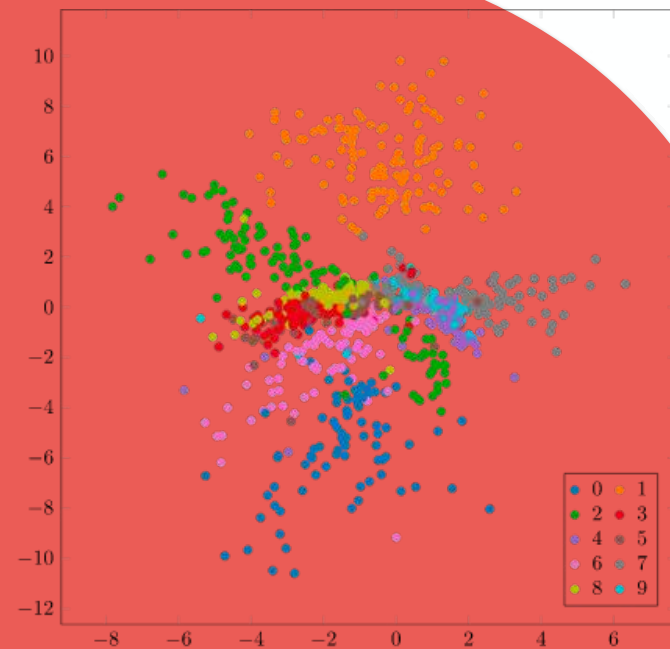


MODELLING OF $P(z|x)$ WITH GMMs

GMMs

- ++ Both sampling and evaluating are straightforward with GMMs
- – However, performance scales poorly with increasing dimensionality
- – Depends on initial choices...
- And nothing guarantees that the latent space is well behaved as a mixture of gaussians

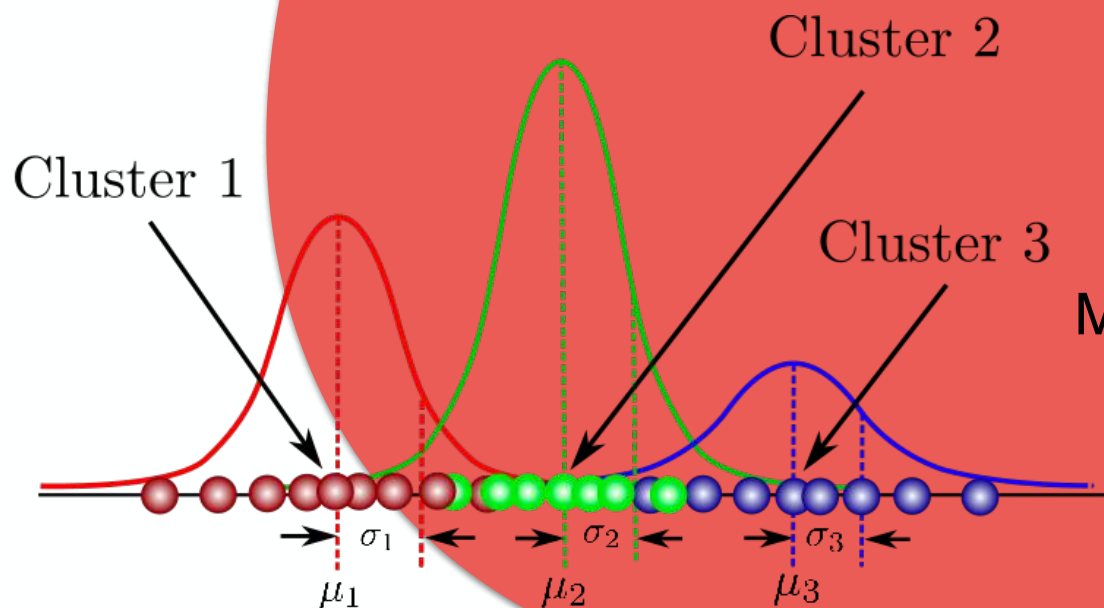
DATA



CAN WE COMBINE
THESE 2 STEPS?



REPRESENTATION
(AUTOENCODER)



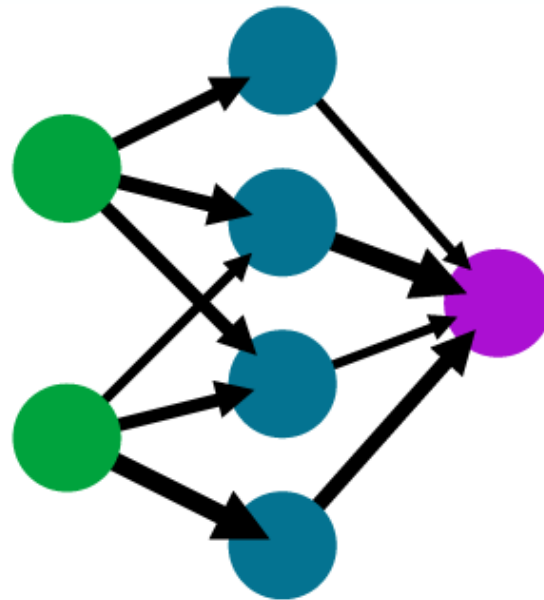
MODELING OF $P(z|x)$ WITH GMMs

Evaluation

High dimensional data
(e.g. images)

NNs == universal
approximators

$p_w(X)$



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

(Easy to sample, easy to evaluate)

(Optimization problem)

Likely not
Gaussian

Sampling

*** Not addressing score based models (diffusion) in this lecture**

Much of the recent progress in unsupervised deep learning has been to invent network architectures that are capable of solving either or both of these related problems directly, without resorting to any auxiliary methods

VAE

(VARIATIONAL AUTOENCODER)

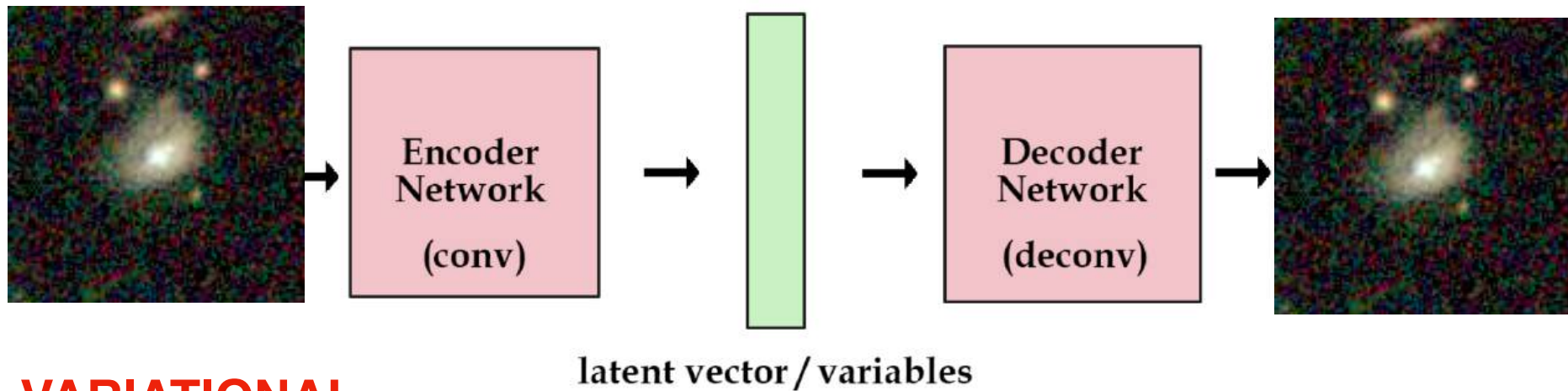
GAN

(GENERATIVE ADVERSARIAL NETWORK)

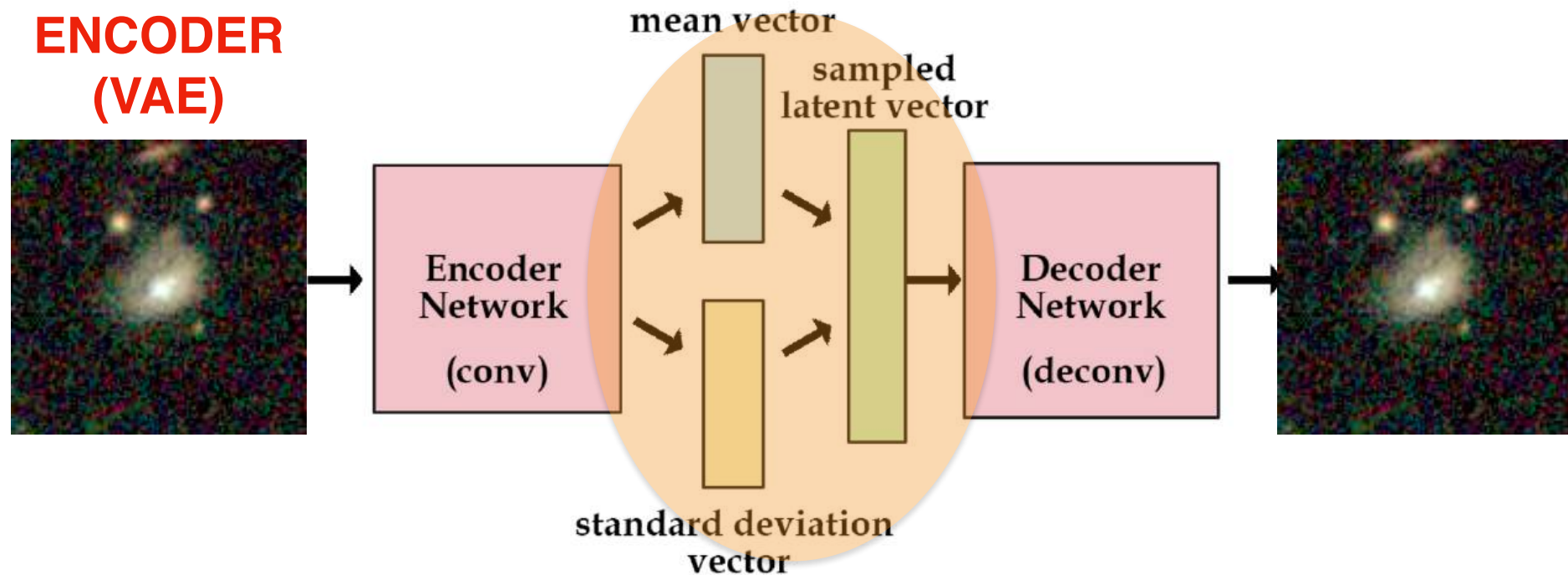
NF-ARF

(NORMALIZING FLOW, AUTOREGRESSIVE FLOW)

AUTO-ENCODER

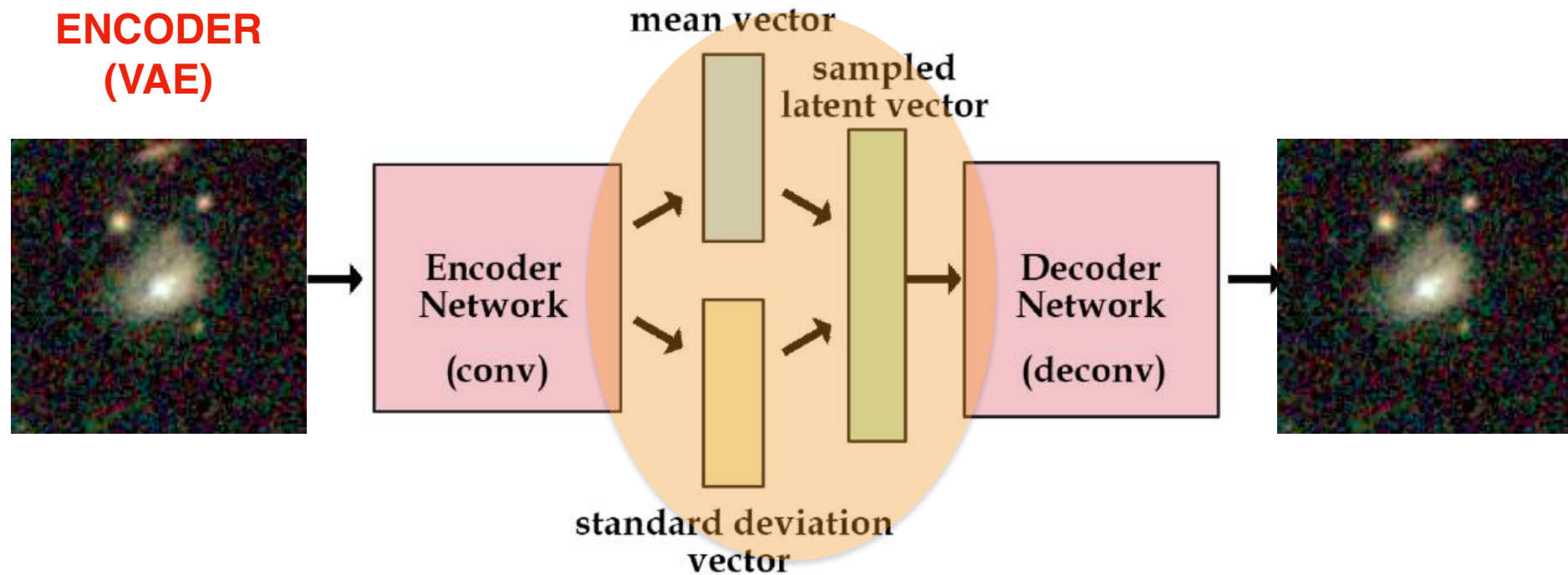


VARIATIONAL AUTO- ENCODER (VAE)



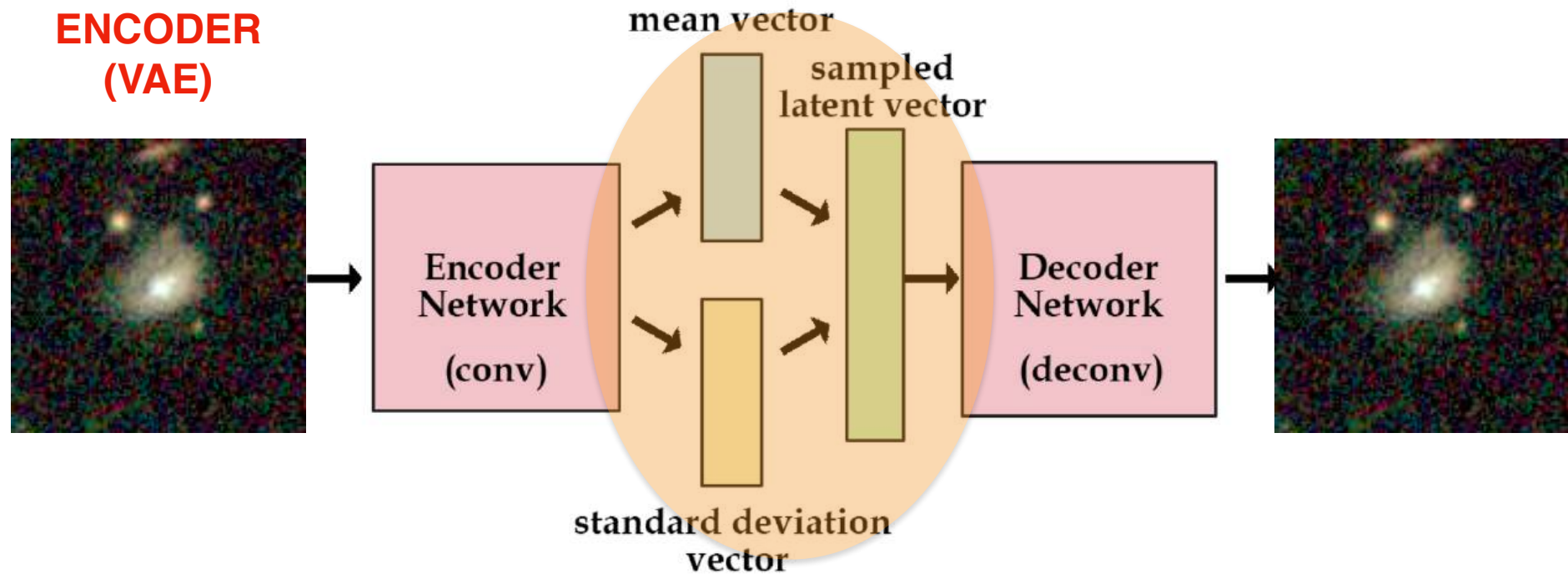
LET'S MODEL THE LATENT SPACE WITH A MIXTURE OF
GAUSSIANS

VARIATIONAL AUTO- ENCODER (VAE)



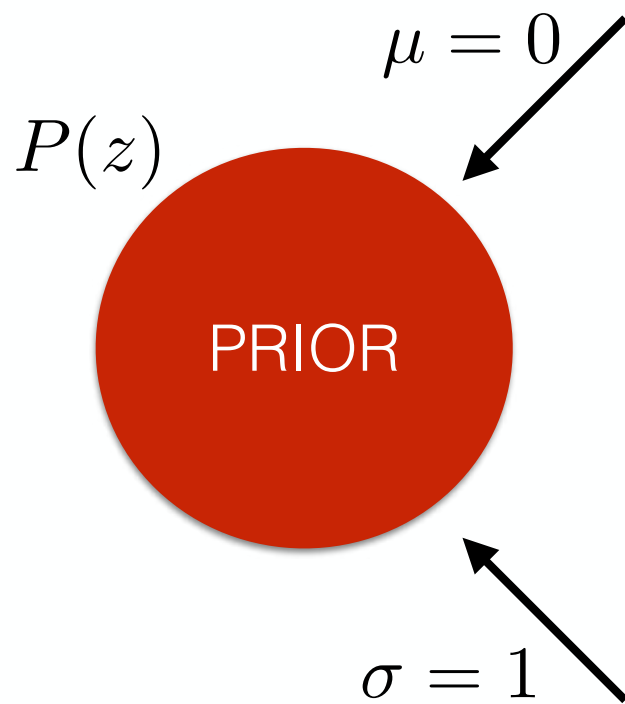
HOWEVER, NOTHING GUARANTEES US THAT $P(z)$ CAN BE ACCURATELY MODELLED BY A MIXTURE OF GAUSSIANS....

VARIATIONAL AUTO- ENCODER (VAE)

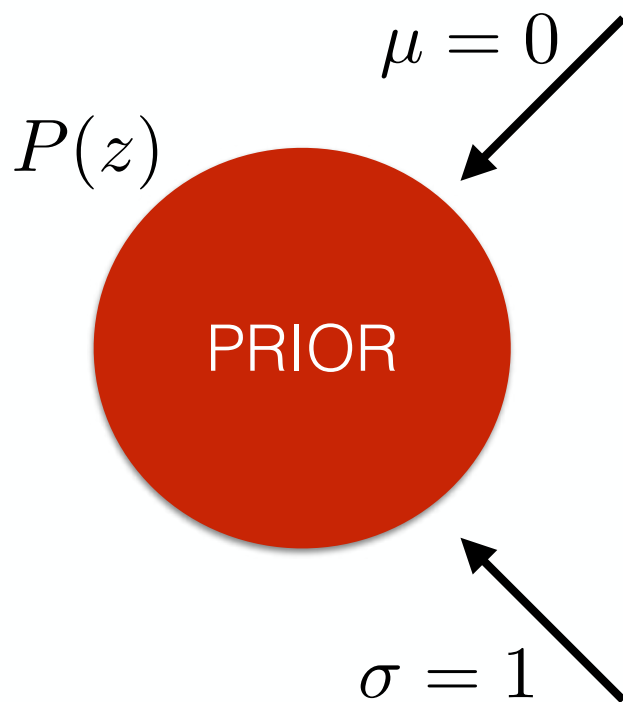


HOWEVER, NOTHING GUARANTEES US THAT THE LATENT SPACE CAN BE MODELLED BY A MIXTURE OF GAUSSIANS....

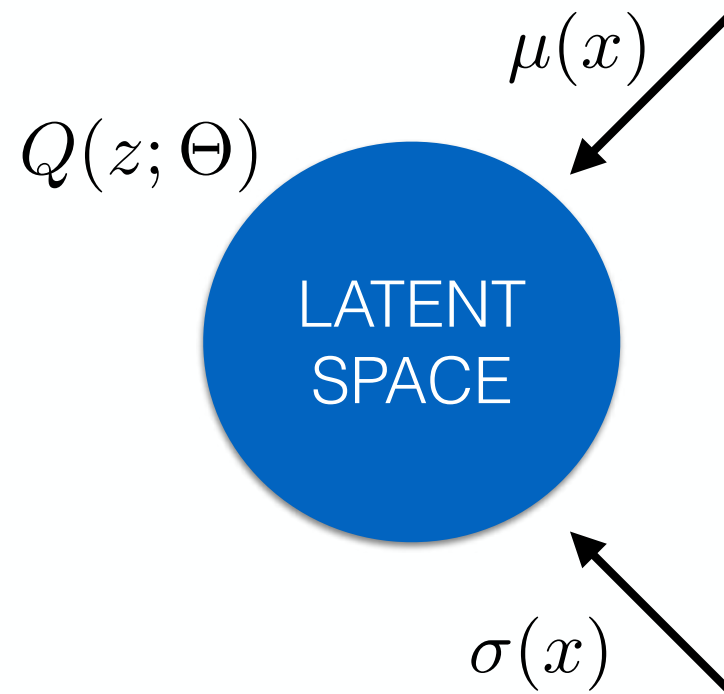
... LET'S FORCE IT TO BE GAUSSIAN LIKE!



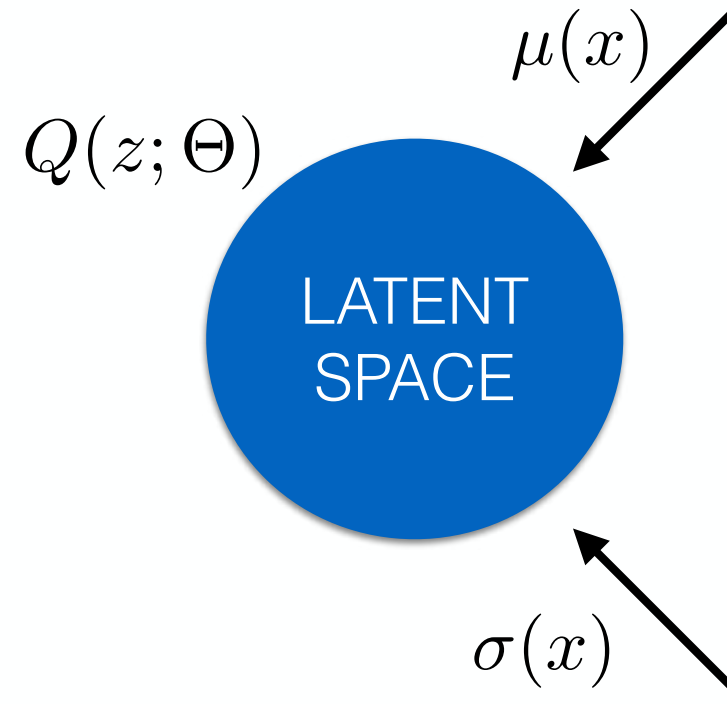
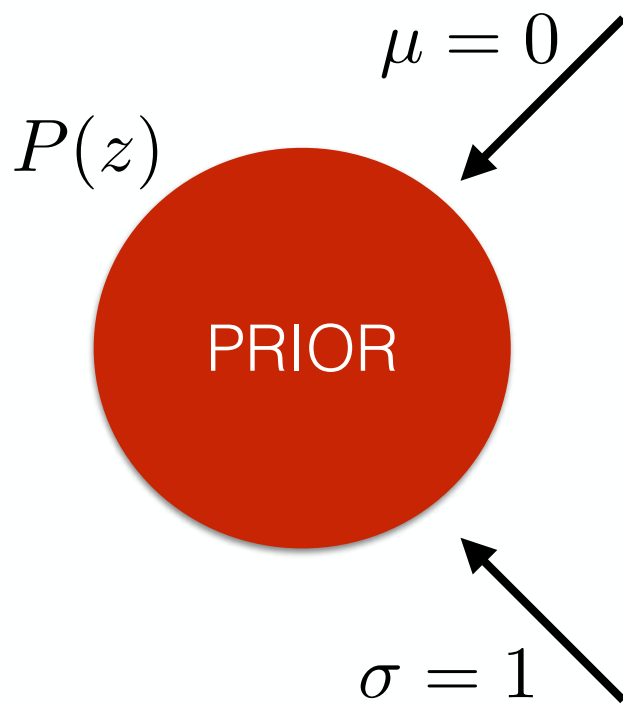
WE ASSUME A SIMPLE PRIOR



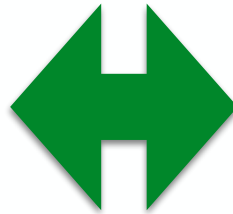
WE ASSUME A SIMPLE PRIOR



LATENT SPACE MODELING

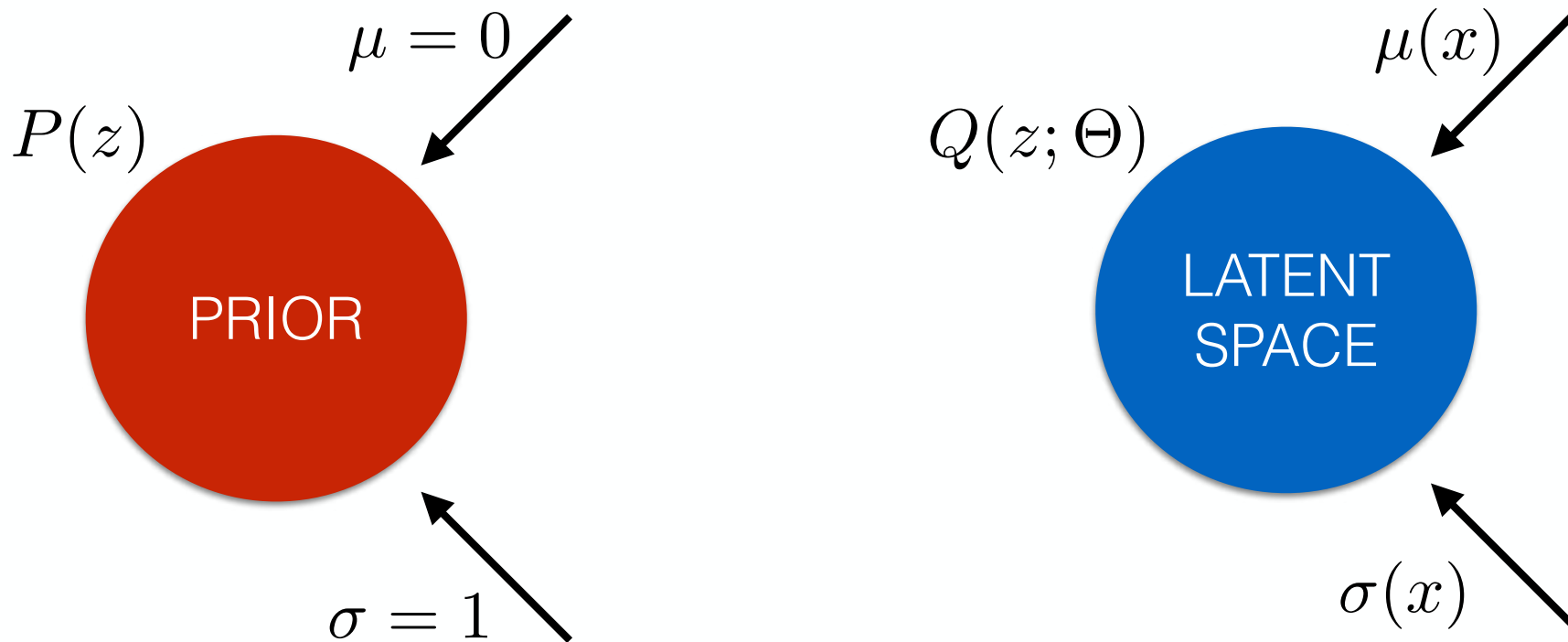


WE ASSUME A SIMPLE PRIOR

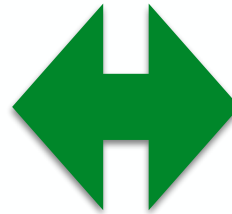


LATENT SPACE MODELLING

WE WANT Q TO BE CLOSE TO THE PRIOR...



WE ASSUME A SIMPLE PRIOR



LATENT SPACE MODELLING

WE WANT Q TO BE CLOSE TO THE PRIOR...

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \log \left(\frac{P(x)}{Q(x)} \right). \quad D_{\text{KL}}(P \parallel Q) = \int_{\mathcal{X}} \log \left(\frac{dP}{dQ} \right) \frac{dP}{dQ} dQ,$$

WE MINIMIZE THE K-L DIVERGENCE BETWEEN P AND Q

WHAT WOULD BE THEN THE LOSS FUNCTION OF A VAE?

The key insight of VAE is that we are actually performing variational inference here, which then tells us what the loss function should be...

$$-\text{ELBO} = \langle \log P(\mathbf{x} \mid \mathbf{z}) \rangle_{\mathbf{z} \sim Q} + \text{KL} (Q(\mathbf{z}; \Theta) \parallel P(\mathbf{z})) ,$$

~MSE

~Gaussian

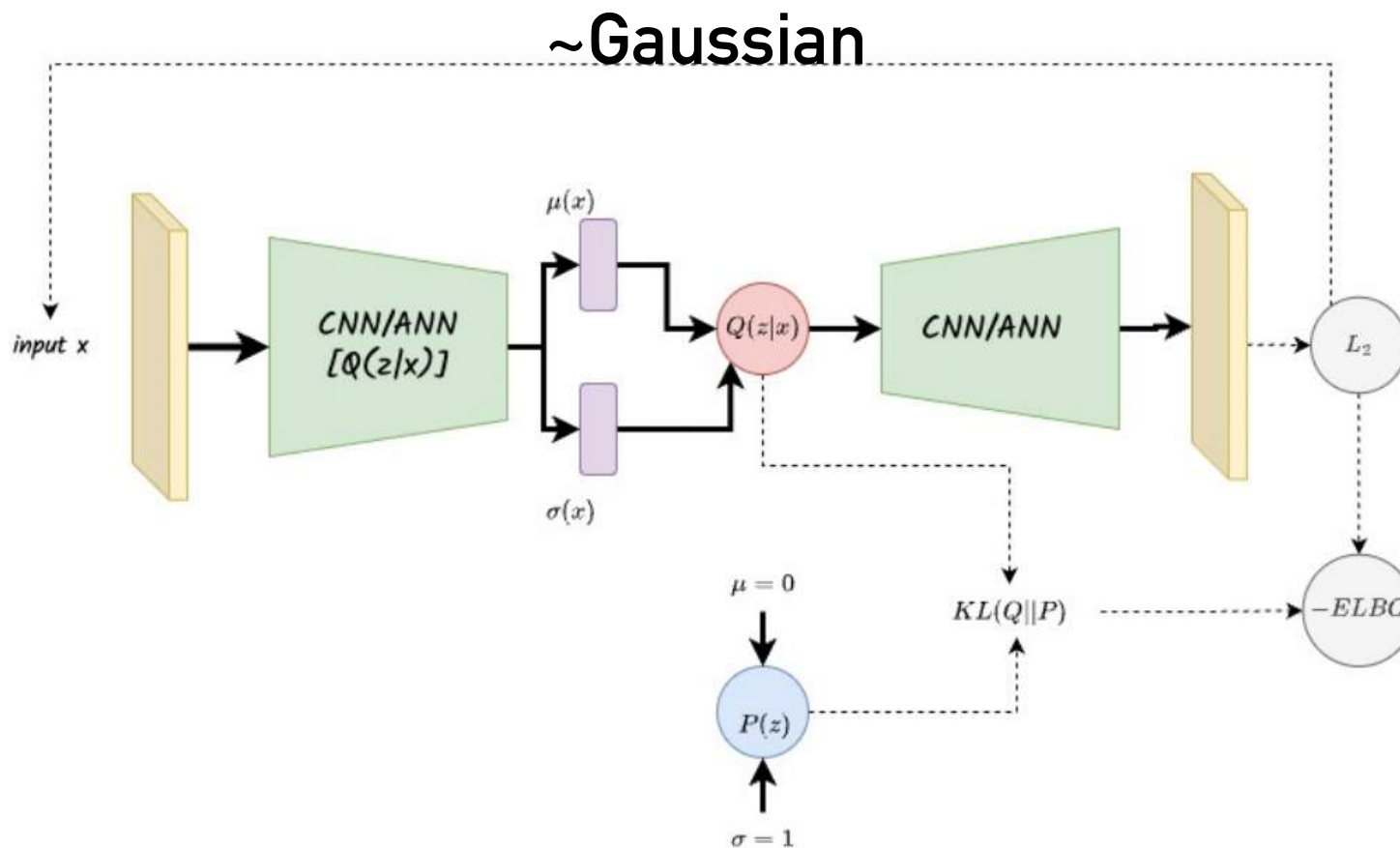
REGULARIZATION TERM

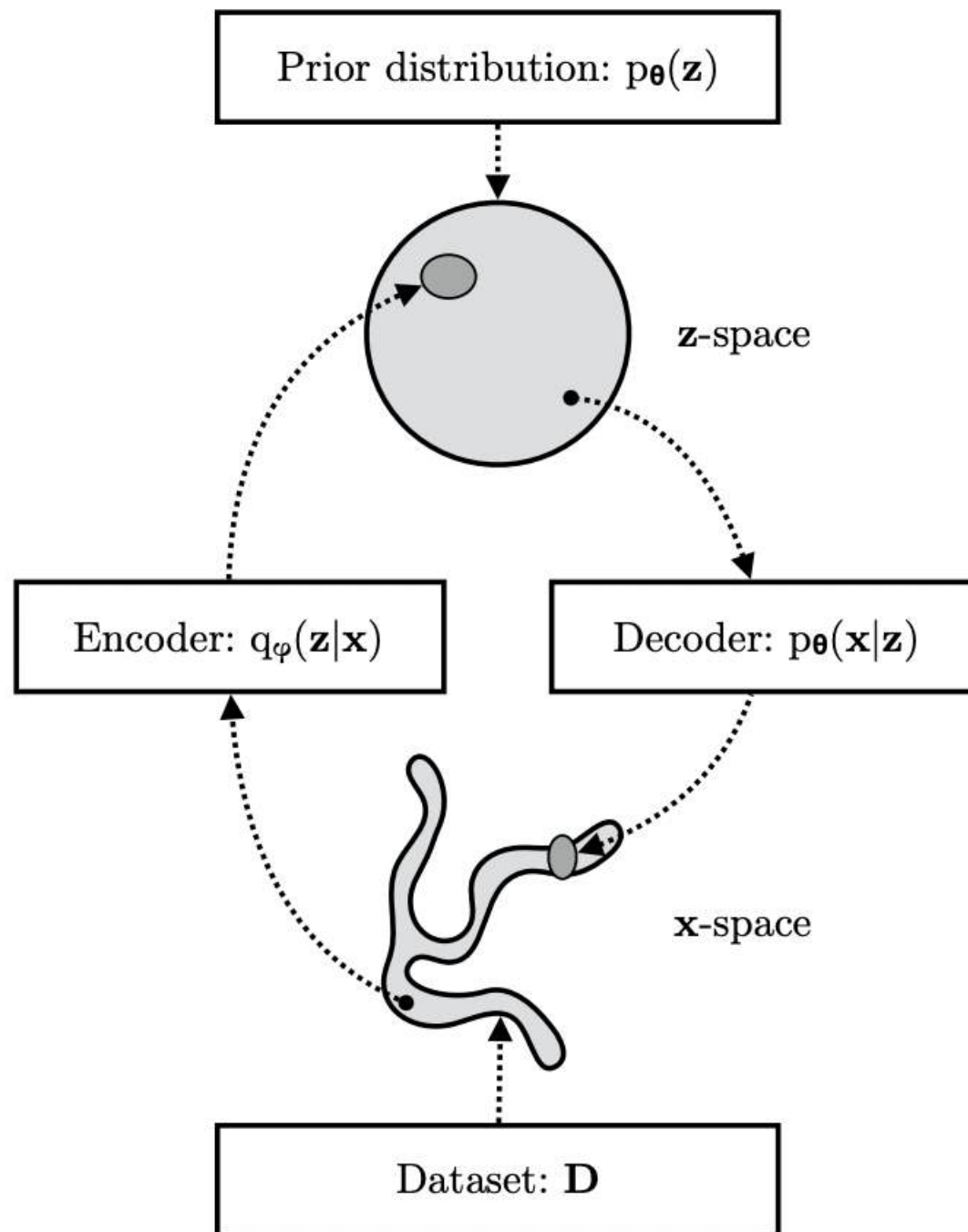
WHAT WOULD BE THEN THE LOSS FUNCTION OF A VAE?

The key insight of VAE is that we are actually performing variational inference here, which then tells us what the loss function should be...

$$-\text{ELBO} = \langle \log P(\mathbf{x} \mid \mathbf{z}) \rangle_{\mathbf{z} \sim Q} + \text{KL} (Q(\mathbf{z}; \Theta) \parallel P(\mathbf{z})) ,$$

L2 LOSS REGULARIZATION TERM





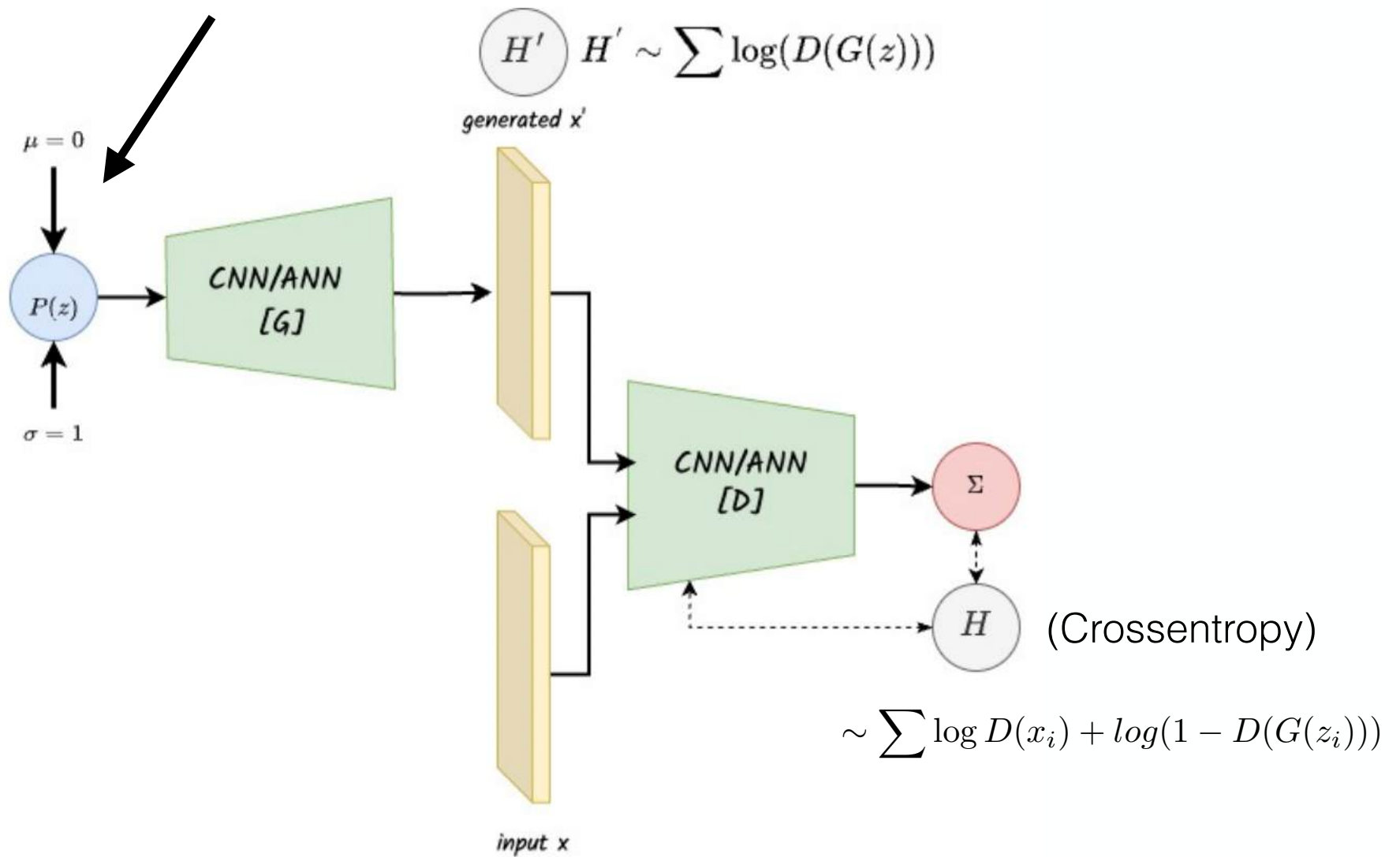
[Kingma and Welling, 2019]

Generative Adversarial Networks

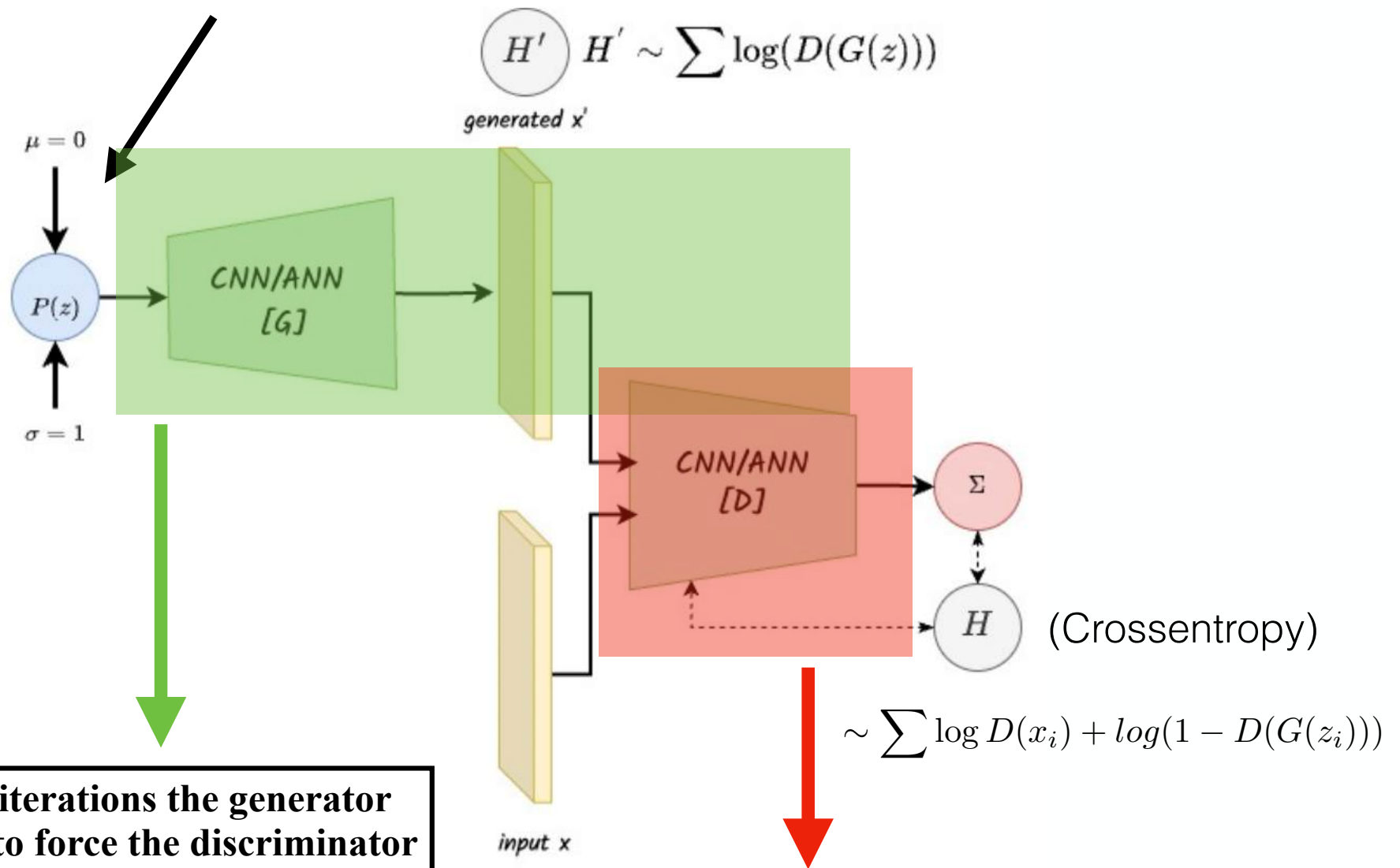
As for VAEs,
the goal of generative Adversarial Networks (GANs) is to estimate
 $p(z|x)$

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

The latent variable is
here



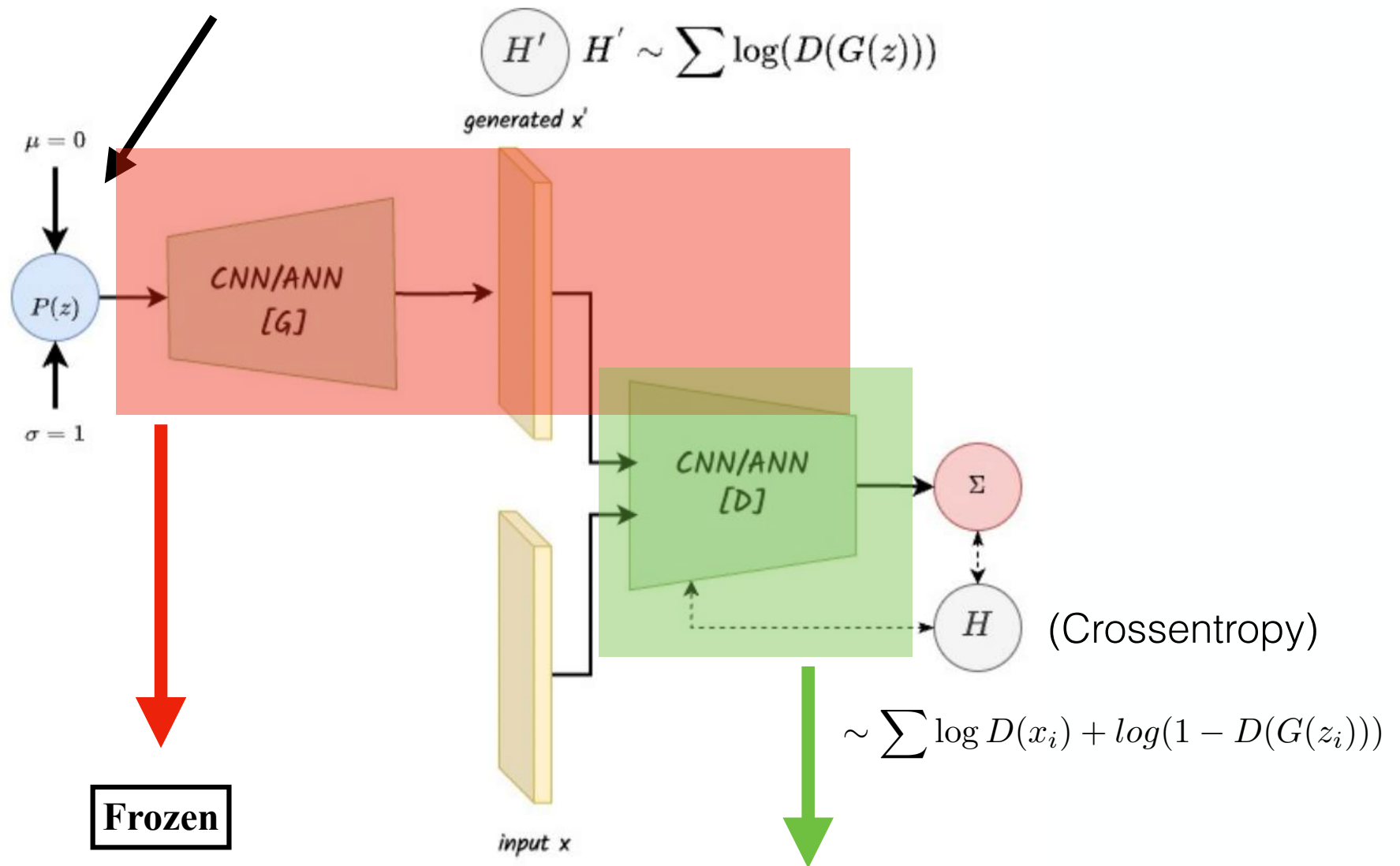
The latent variable is
here



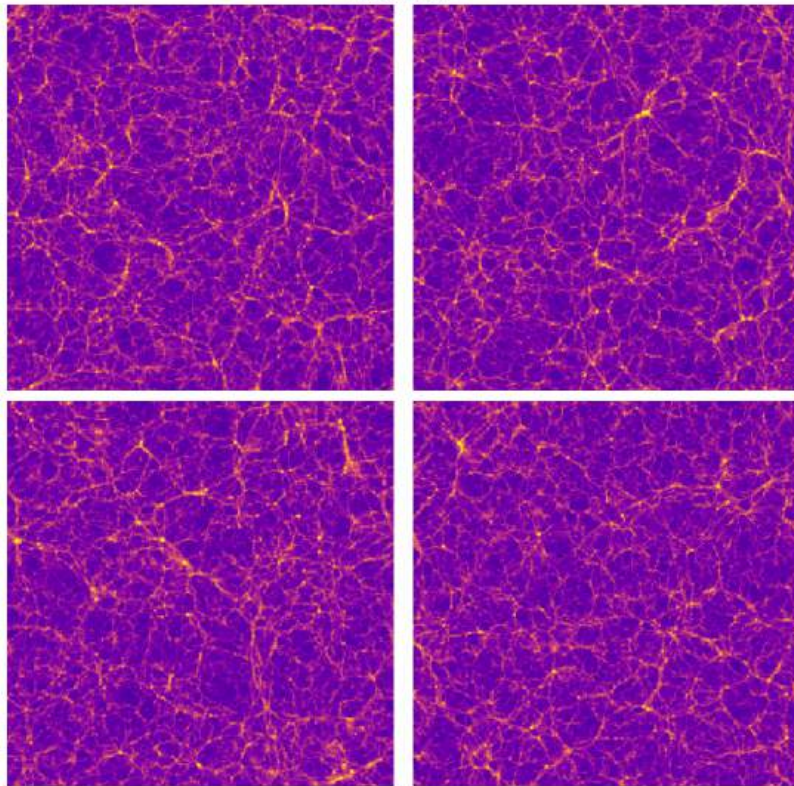
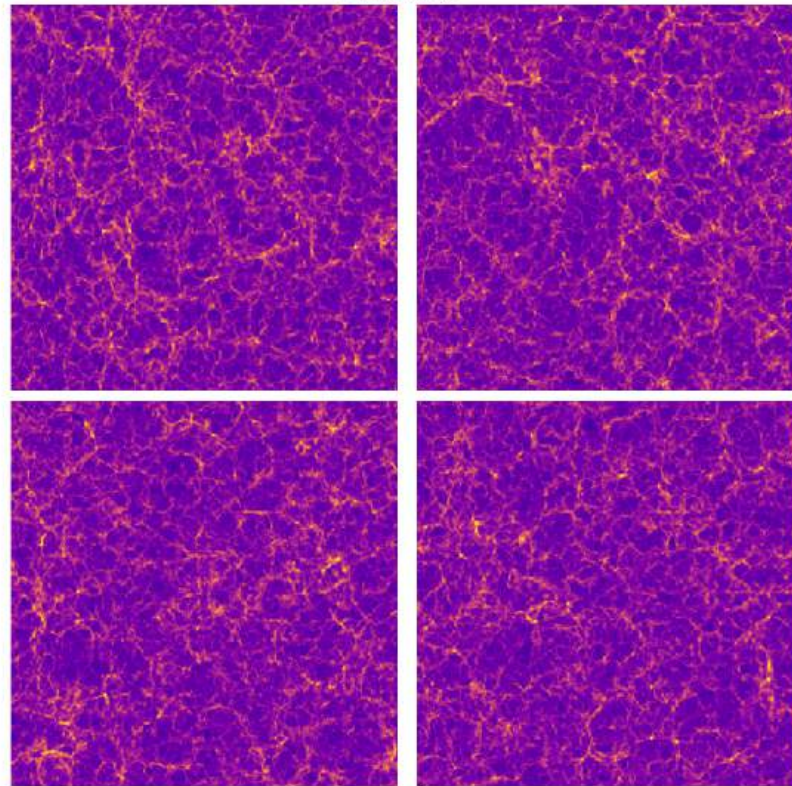
Every N iterations the generator
is trained to force the discriminator
to classify as real

Frozen

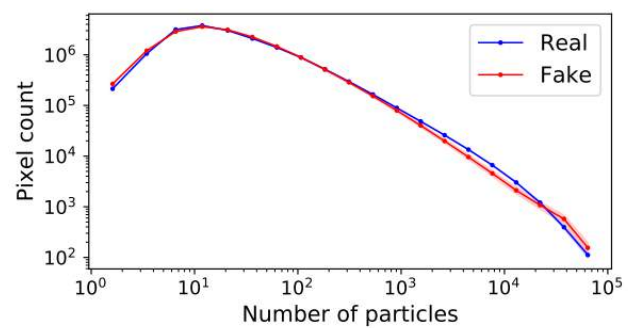
The latent variable is
here



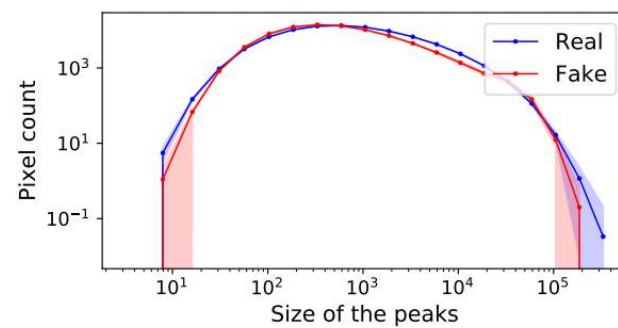
Every N iterations the discriminator
is trained to force to distinguish between
real and fake

Real 256^3 Fake 256^3 

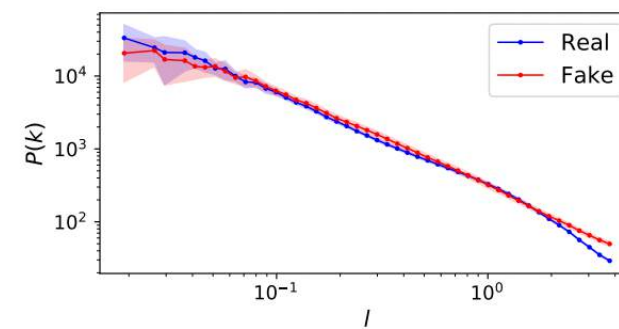
Mass histogram

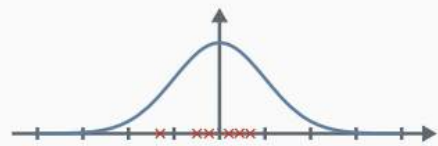


Peak histogram

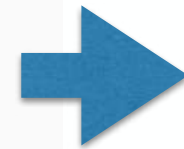
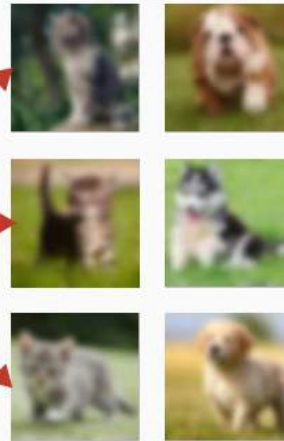


Power spectral density

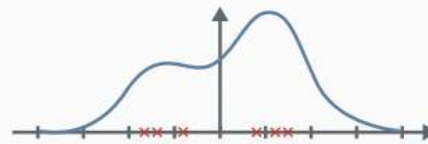




(Sampling)



easy with
VAEs and GANs



~ VAEs
difficult with
GANs

(Density estimation)

3. Density Estimation for likelihood evaluation

Normalizing Flows

Based on change of variables:

$$p_X(x) = p_Z(f(x)) | \det J_{f(x)} |$$

unknown pdf

tractable pdf
(typically Gaussian)

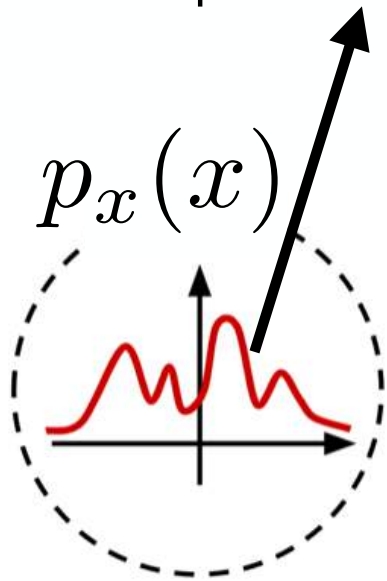
Determinant of the Jacobian
of $f(x)$

Invertible differentiable
function

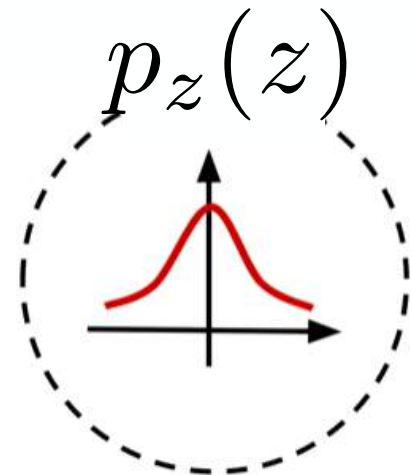
$$p_X(x) = p_Z(f(x)) | \det J_{f(x)} |$$

Can represent $p_X(x)$ in terms of $f(x)$ and $p_Z(z)$
 (notice there is no dimensionality reduction as opposed to VAEs)

a sample x_i



(unknown)

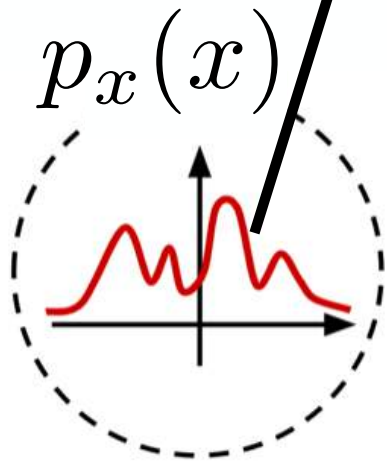


(known and easy)

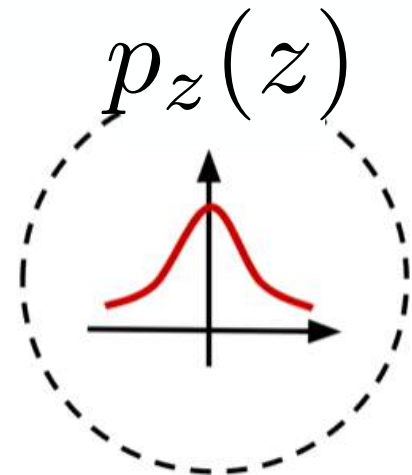
$$p_X(x) = p_Z(f(x)) |det J_{f(x)}|$$

Can represent $p_X(x)$ in terms of $f(x)$ and $p_Z(x)$
 (notice there is no dimensionality reduction as opposed to VAEs)

a sample $x_i \longrightarrow f_w(x_i)$



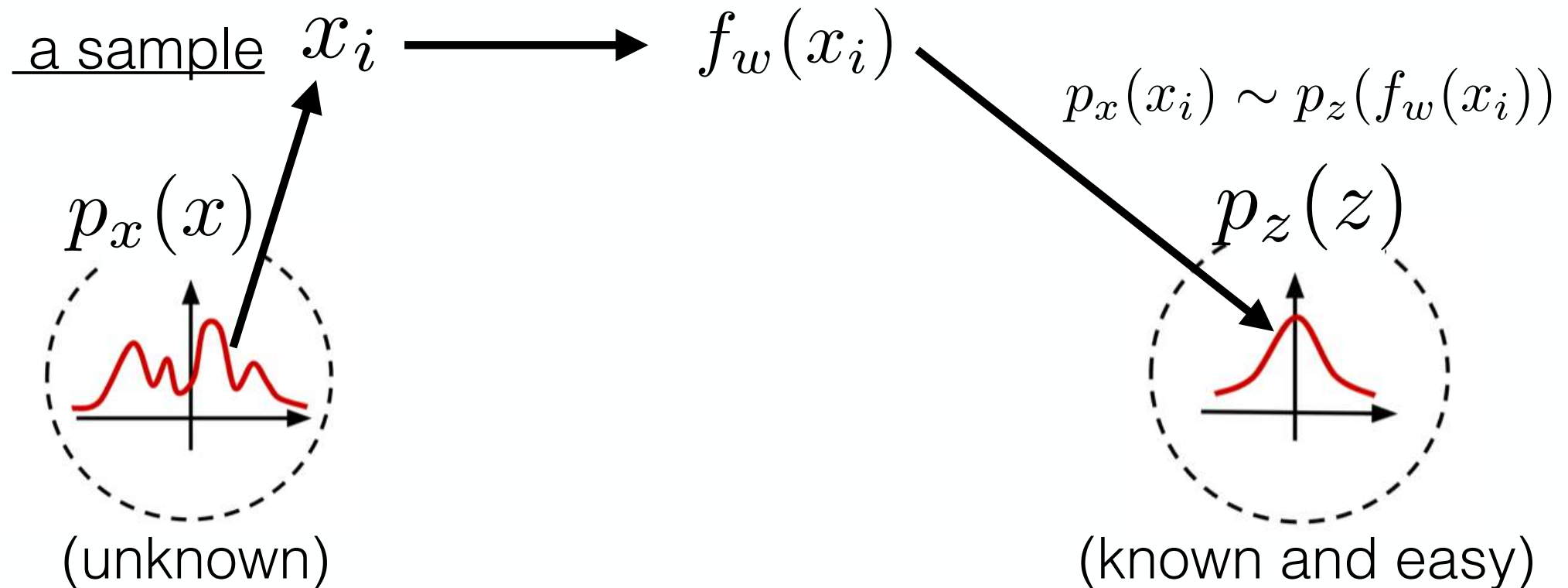
(unknown)



(known and easy)

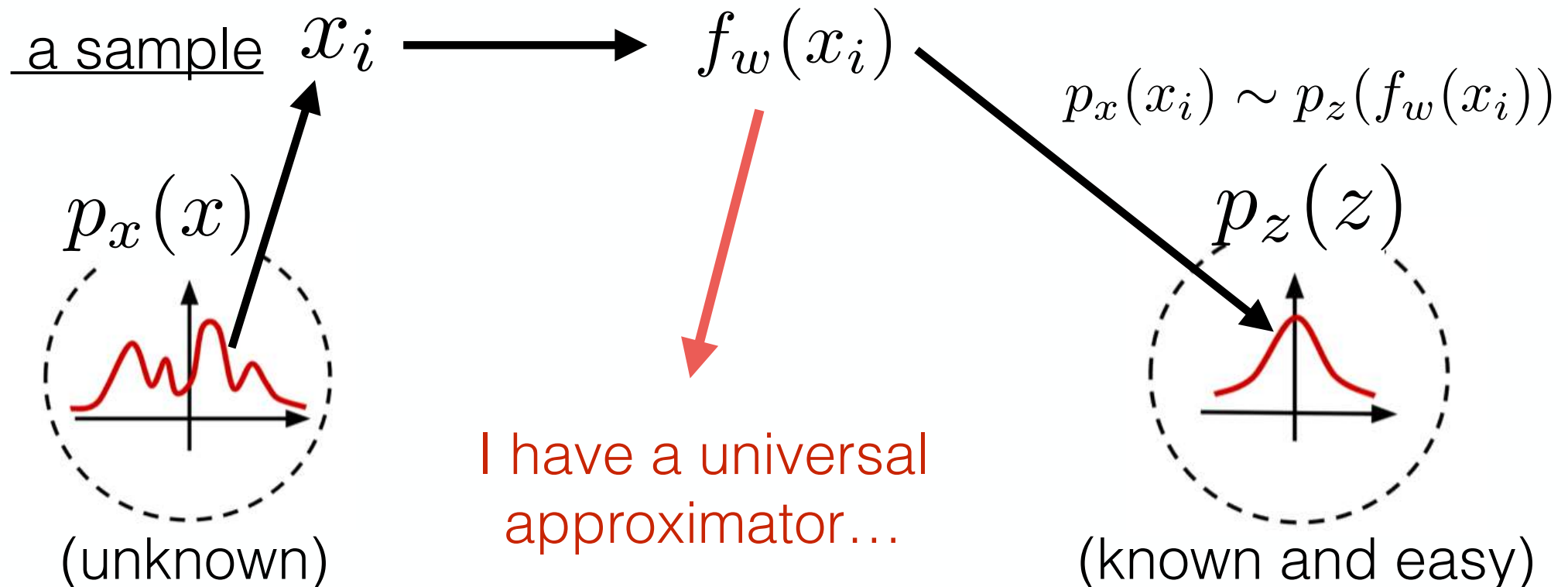
$$p_X(x) = p_Z(f(x)) |det J_{f(x)}|$$

Can represent $p_X(x)$ in terms of $f(x)$ and $p_Z(z)$
 (notice there is no dimensionality reduction as opposed to VAEs)

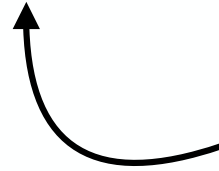


$$p_X(x) = p_Z(f(x)) | \det J_{f(x)} |$$

Can represent $p_X(x)$ in terms of $f(x)$ and $p_Z(x)$
 (notice there is no dimensionality reduction as opposed to VAEs)



Then $f(x)$ becomes: $f(x; W)$



NN learnable weights

Then $f(x)$ becomes: $f(x; W)$

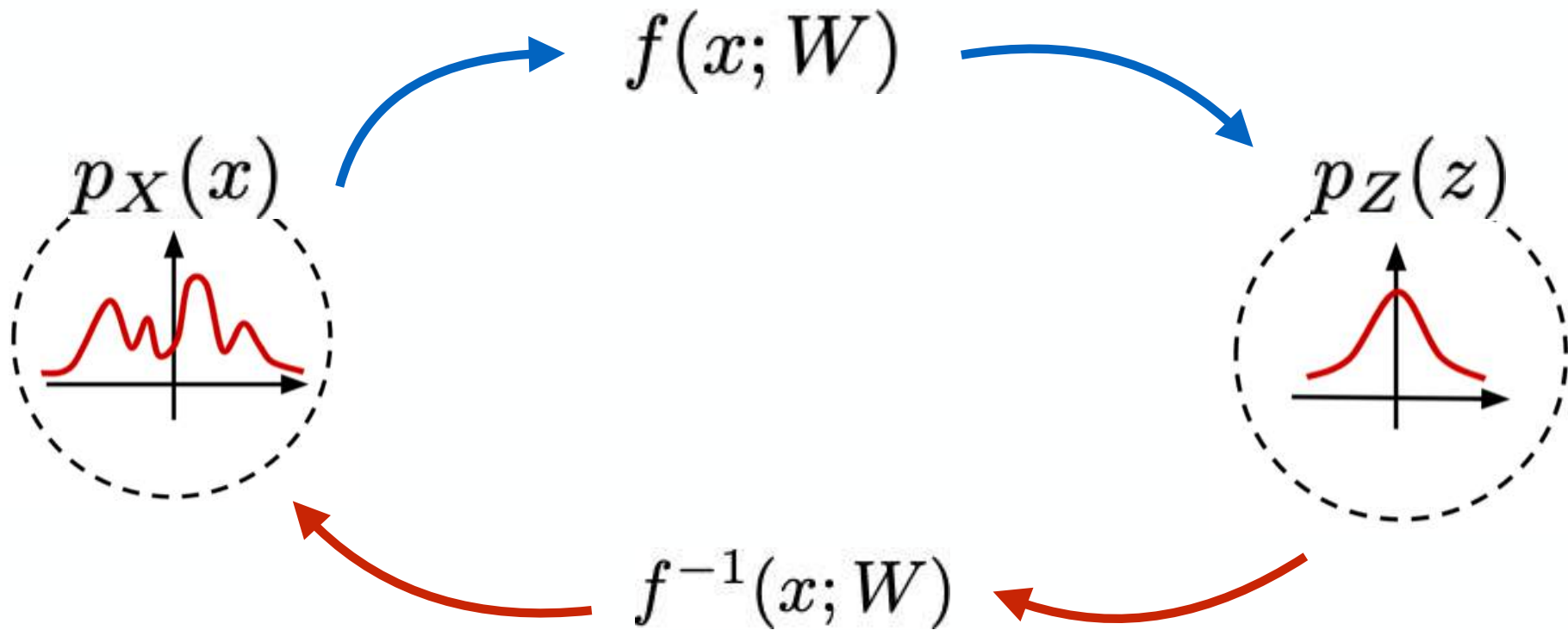
 NN learnable weights

And the loss function: log likelihood

$$\sum_{i=1}^N \log(p_Z(f(x_i; W)) + \log(|\det J_{f(x_i; W)}|))$$

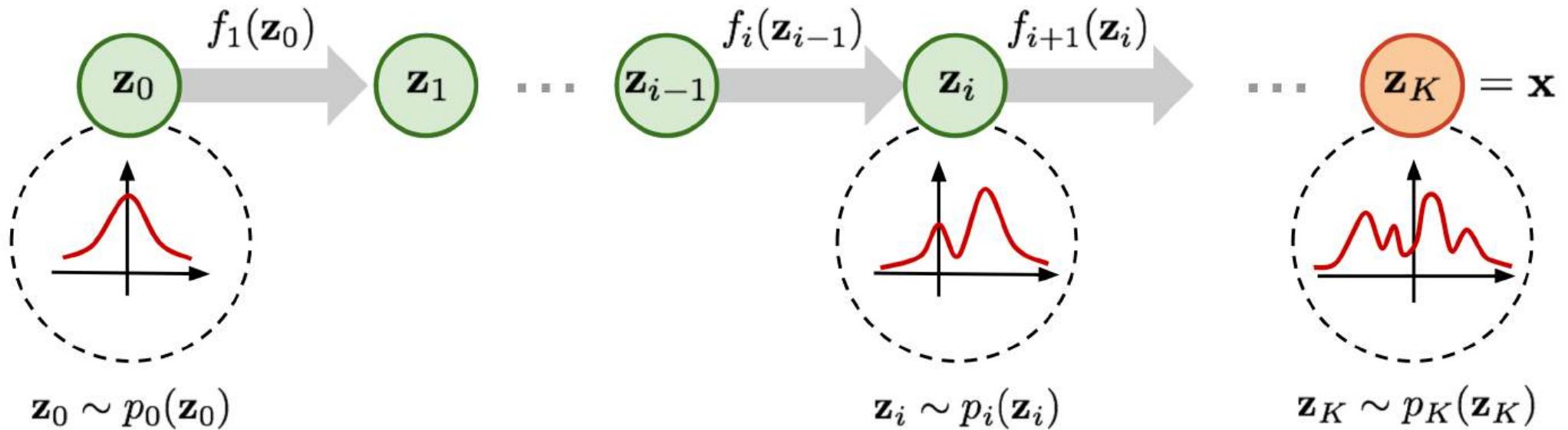
 This is a gaussian

Likelihood evaluation



Sampling

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

What functions satisfy these conditions?

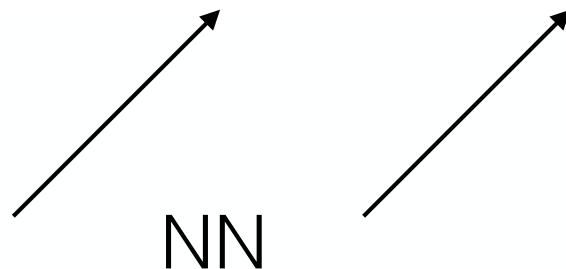
- 1. Easily invertible
- 2. Jacobian easy to compute

What functions satisfy these conditions?

- 1. Easily invertible
- 2. Jacobian easy to compute

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

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



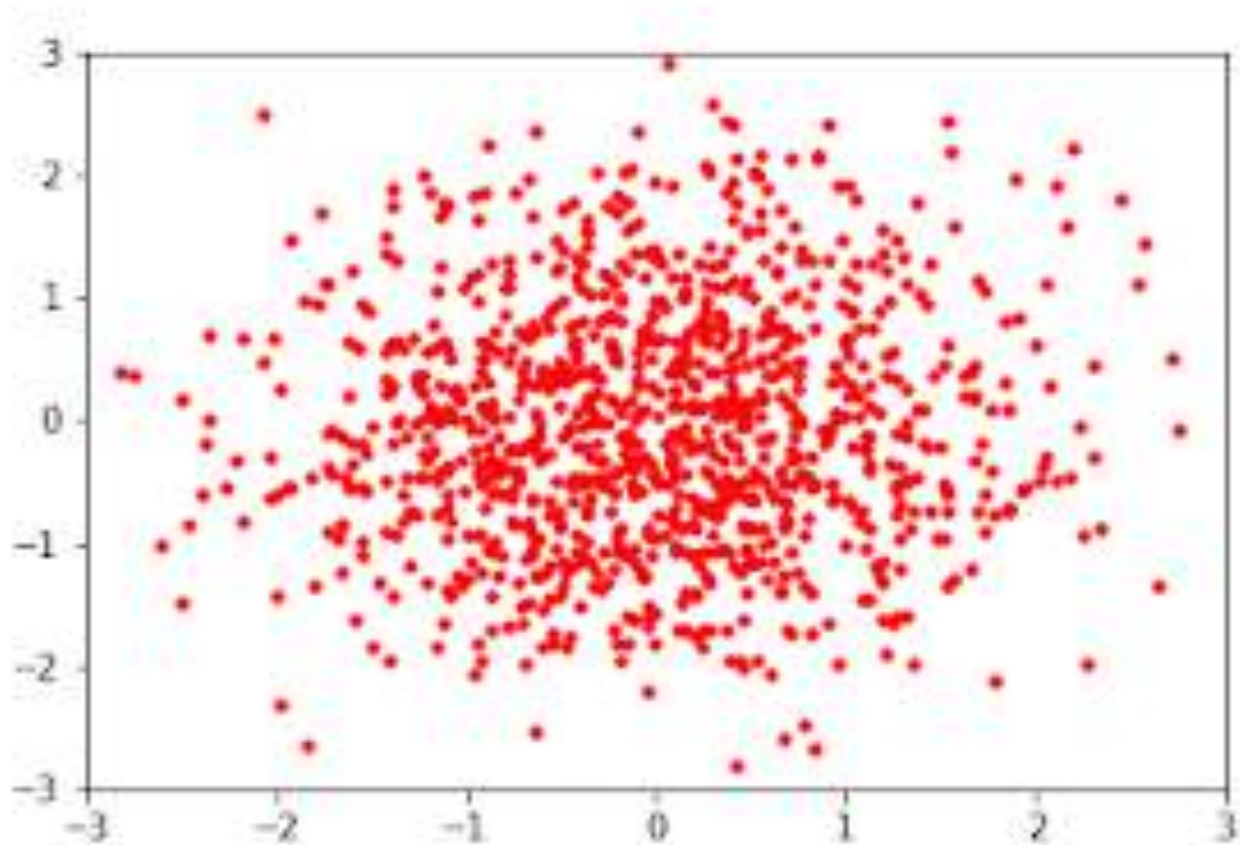
- Easily invertible:

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

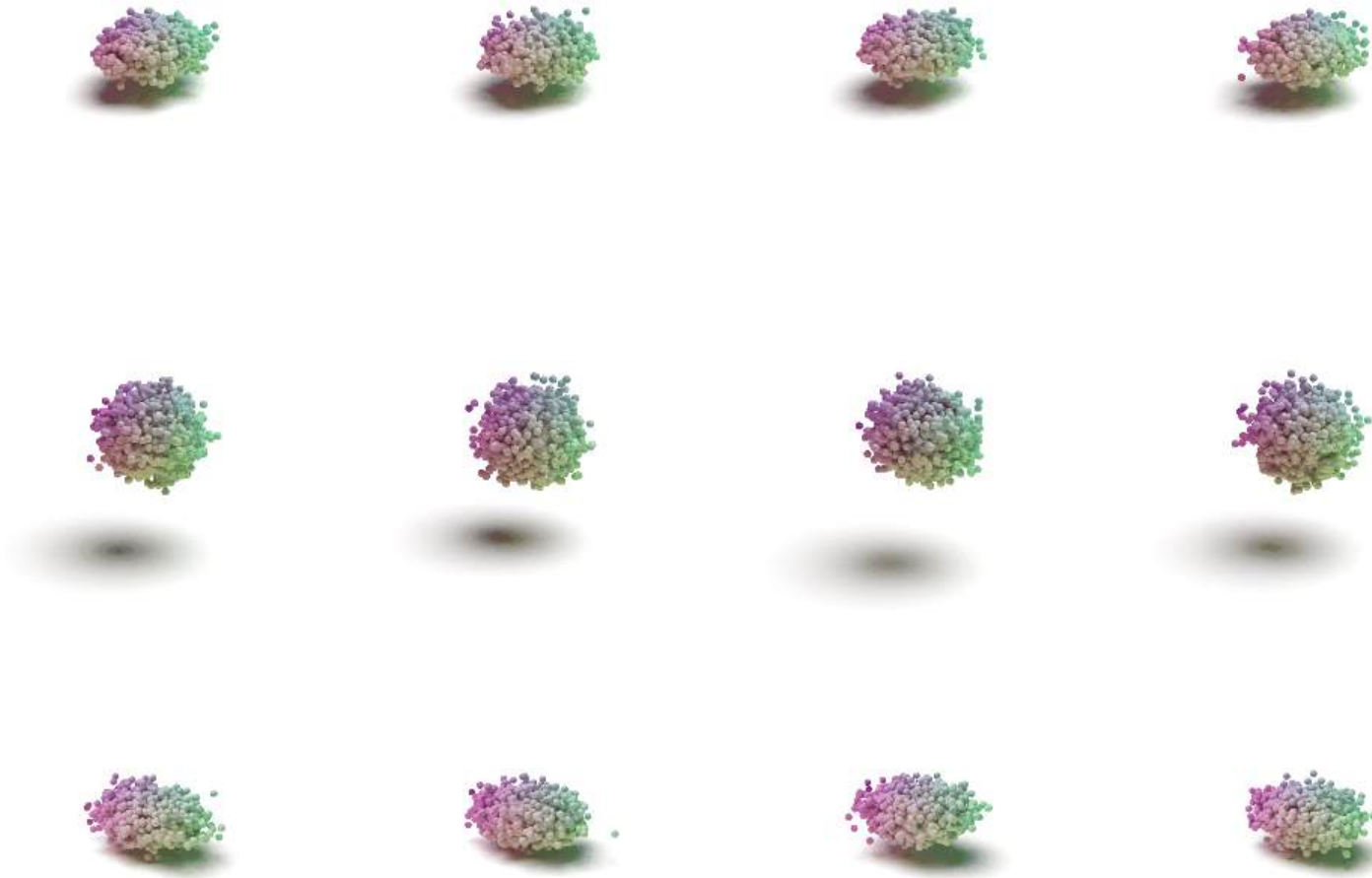
- Jacobian:

$$\mathbf{J} = \begin{bmatrix} \mathbb{I}_d & \mathbf{0}_{d \times (D-d)} \\ \frac{\partial \mathbf{y}_{d+1:D}}{\partial \mathbf{x}_{1:d}} & \text{diag}(\exp(s(\mathbf{x}_{1:d}))) \end{bmatrix} \quad \det(\mathbf{J}) = \prod_{j=1}^{D-d} \exp(s(\mathbf{x}_{1:d})_j) = \exp\left(\sum_{j=1}^{D-d} s(\mathbf{x}_{1:d})_j\right)$$

normalizing flows are generative models that are easy to evaluate and flexibly expressive

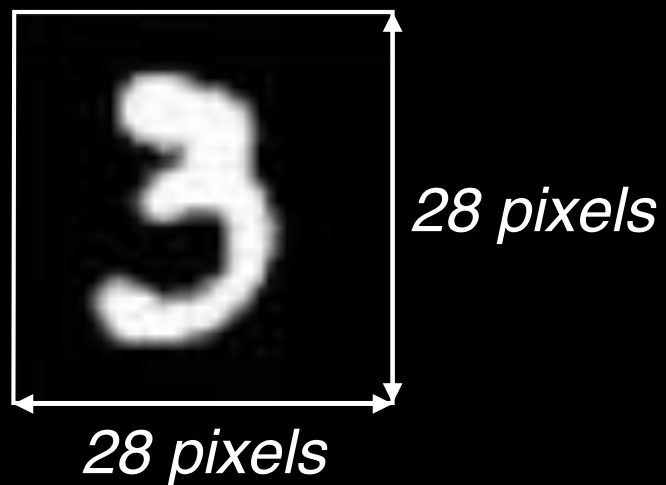


normalizing flows are generative models that are easy to evaluate and flexibly expressive



3	4	2	1	9	5	6	2	1	8
8	9	1	2	5	0	0	6	6	4
6	7	0	1	6	3	6	3	7	0
3	7	7	9	4	6	6	1	8	2
2	9	3	4	3	9	8	7	2	5
1	5	9	8	3	6	5	7	2	3
9	3	1	9	1	5	8	0	8	4
5	6	2	6	8	5	8	8	9	9
3	7	7	0	9	4	8	5	4	3
7	9	6	4	7	0	6	9	2	3

train a generative model on MNIST — *estimate*
 $p(\text{pixels}) \approx q_{\phi}(\text{pixels})$

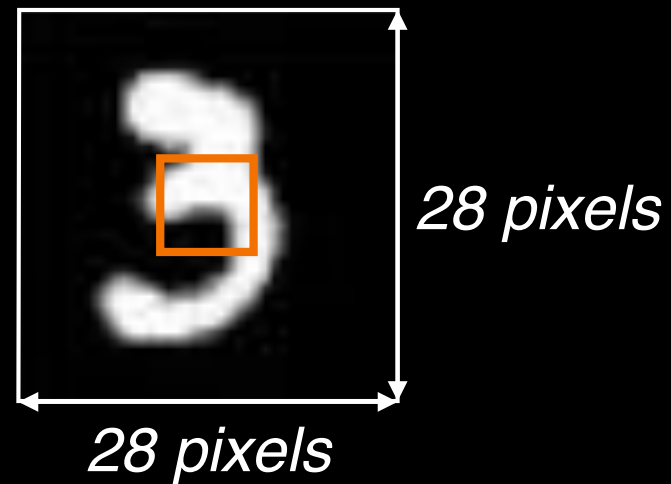


sample q_ϕ (pixels)

sample $q_\phi(\text{pixels})$

3	8	3	3	0	1	4	1	6	3
5	0	5	7	5	1	5	1	0	9
5	1	0	3	1	5	9	6	4	4
6	4	9	8	7	5	0	1	2	7
4	5	2	5	5	3	4	9	8	7
0	4	4	1	7	7	9	7	3	0
7	2	2	1	8	8	8	8	0	4
3	9	0	7	3	3	4	1	3	9
3	9	7	7	3	1	1	0	6	3
7	0	2	2	5	9	2	5	3	5

train a generative model to estimate **conditional distributions**



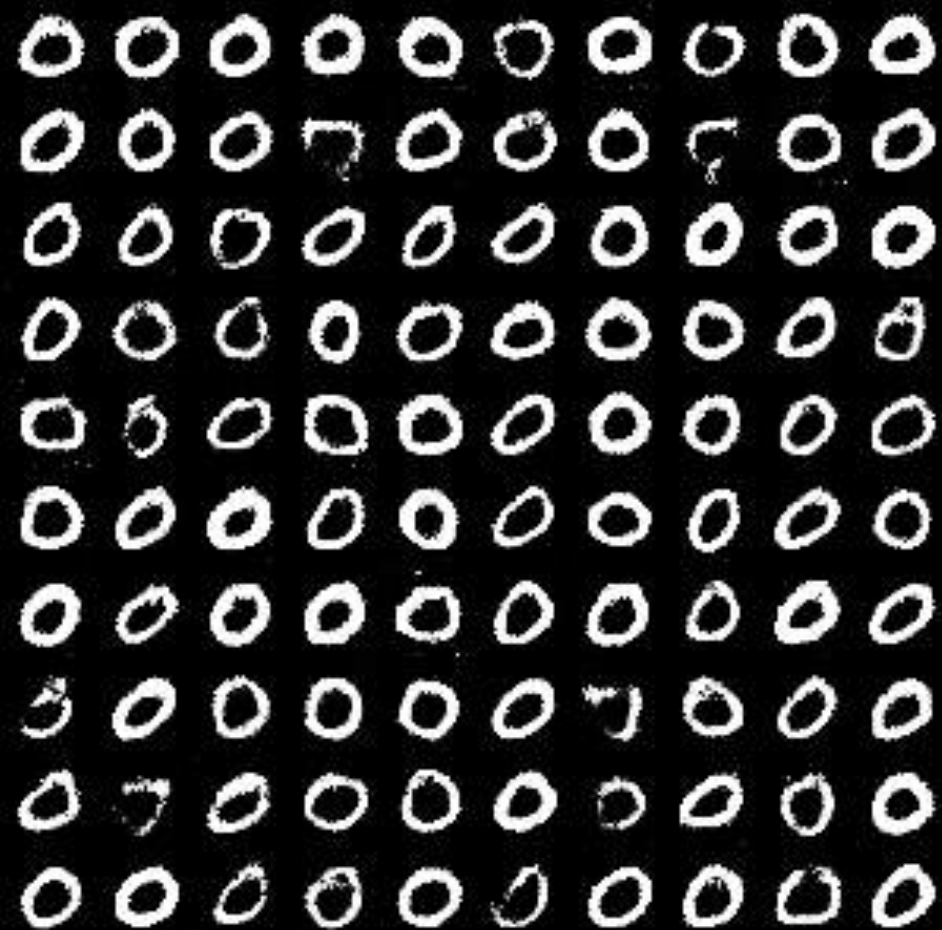
$$p(\text{pixels} \mid \text{central pixels}) \approx q_{\phi}(\text{pixels} \mid \text{central pixels})$$

train a generative model to estimate **conditional distributions**



samples from $p(\text{pixels} \mid \text{central pixels} = 0)$ *should* just be 0s

samples drawn from $q_{\phi}(\text{pixels} \mid \text{central pixels} = 0)$



$$p(\text{pixels} \mid \text{central pixels}) \approx q_\phi$$

θ

$$p(\text{pixels} | \text{central pixels}) \approx q_\phi$$

X

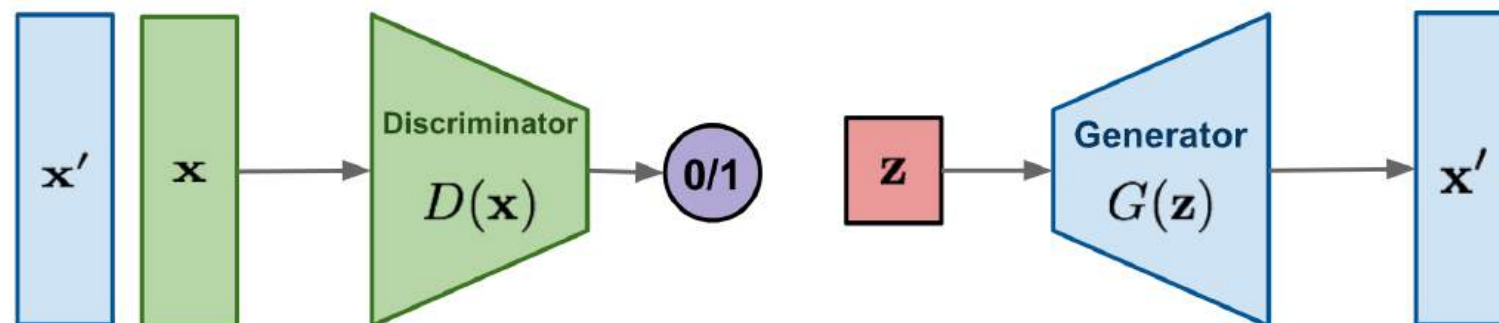
GANs AND VAEs ARE VERY POWERFUL BUT DO NOT PROVIDE AN EXPLICIT LIKELIHOOD

Method	Train on data	One-pass Sampling	Exact log-likelihood	Free-form Jacobian
Variational Autoencoders	✓	✓	✗	✓
Generative Adversarial Nets	✓	✓	✗	✓
Likelihood-based Autoregressive	✓	✗	✓	✗

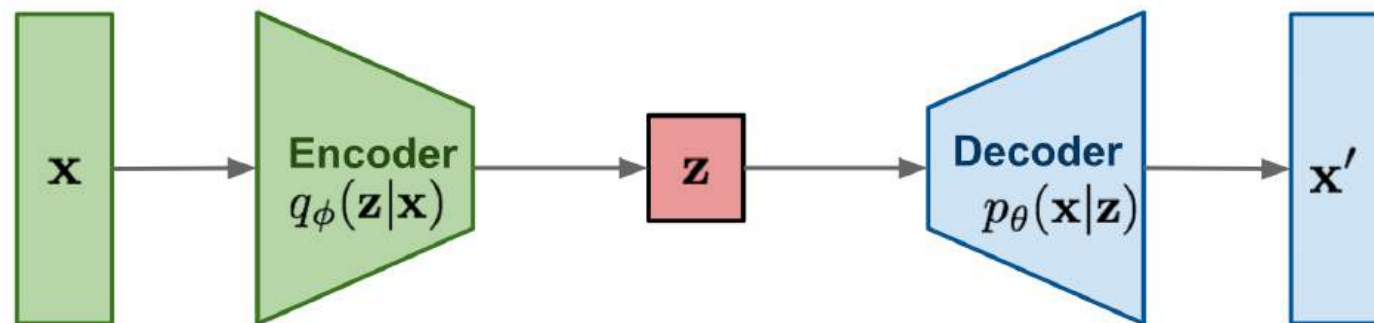
Grathwohl+18

* there is a third member of the family now: score based models which try to estimate the gradient of the likelihood (score) instead of the likelihood

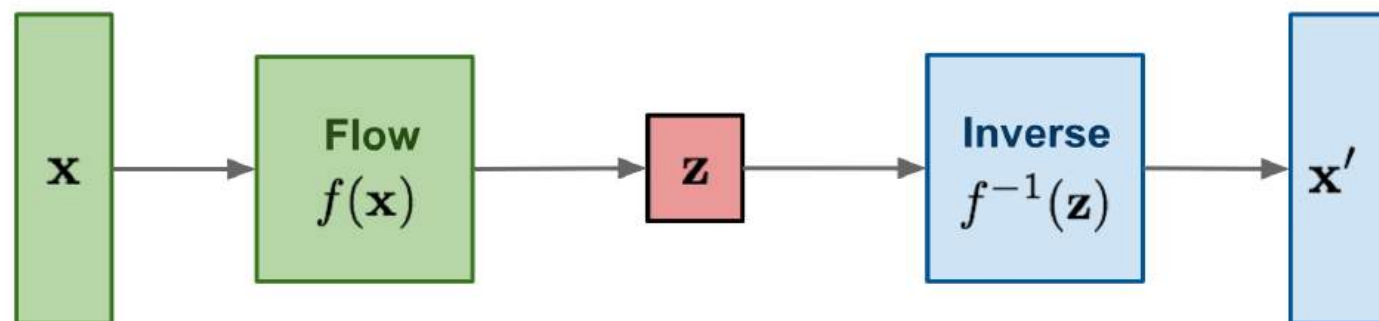
GAN: minimax the classification error loss.



VAE: maximize ELBO.



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



(+ score based diffusion models)

4. Simulation Based Inference (SBI)

“standard” bayesian inference

goal of inference

posterior

likelihood *prior*

$$p(\theta | X) = \frac{p(X | \theta) p(\theta)}{p(X)}$$

$$\ln p(X | \theta) = \ln \mathcal{L} = (X - m(\theta))^T \mathbf{C}^{-1} (X - m(\theta))$$

↑
covariance matrix

↑
model

“standard” bayesian inference

$$\ln p(\theta | X) = \ln p(X | \theta) + \ln p(\theta) + \text{const} .$$

sample $p(\theta | X)$ using montecarlo sampling to estimate posterior
(*e.g. MCMC, HMC, nested sampling*)

Gaussian likelihood is often **an incorrect assumption (not enough data)**

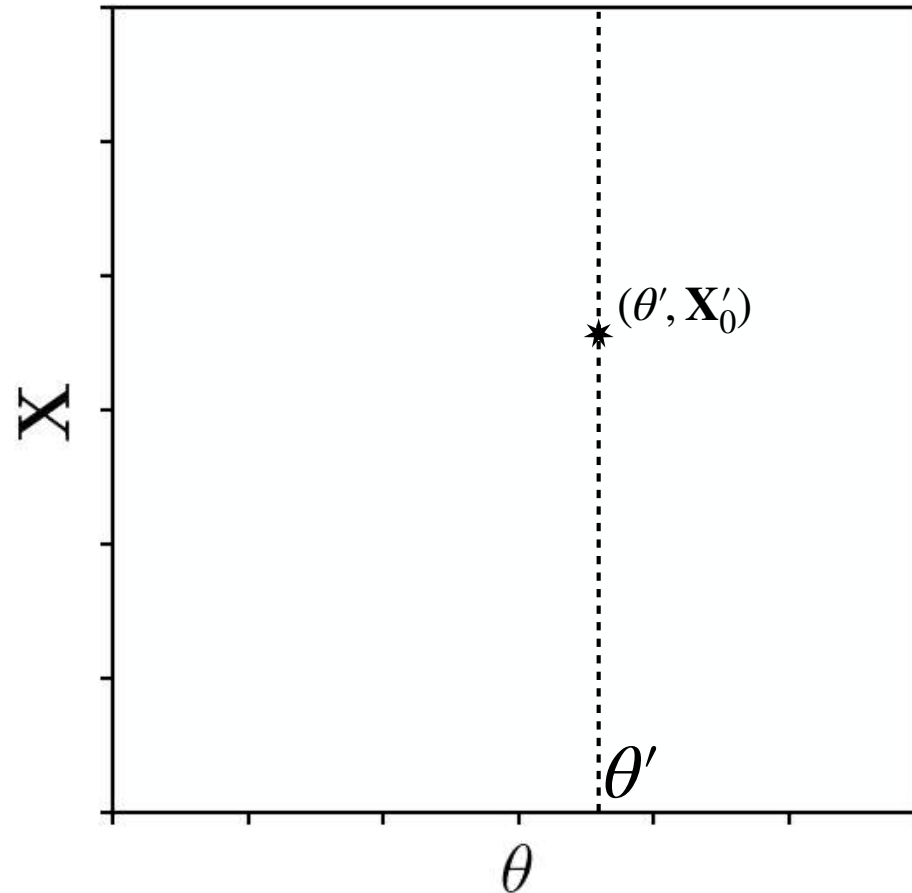
$$\ln p(X | \theta) = \ln \mathcal{L} \neq (X - m(\theta))^T \mathbf{C}^{-1} (X - m(\theta))$$

Sometimes the likelihood is intractable. e.g. Baryonic properties of DM haloes?

MCMC is slow! amortized inference to infer *billions* of observations

forward model F that can simulate observations:

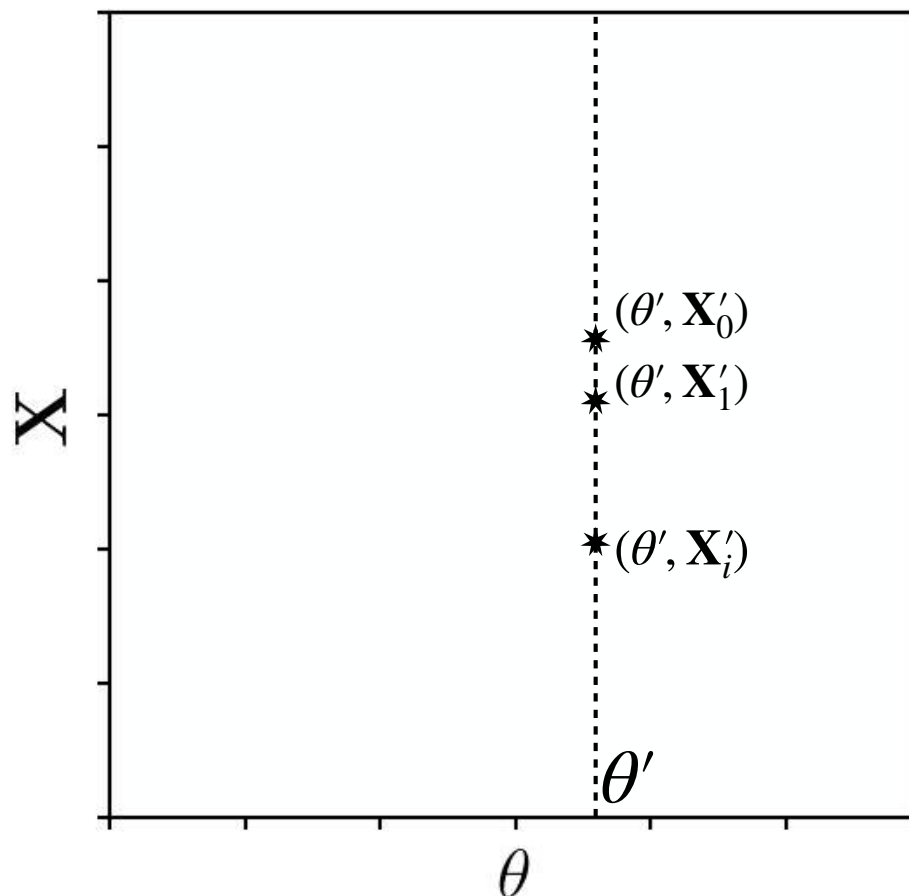
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F must **include noise model** and **observational systematics**

forward model F that can simulate observations:

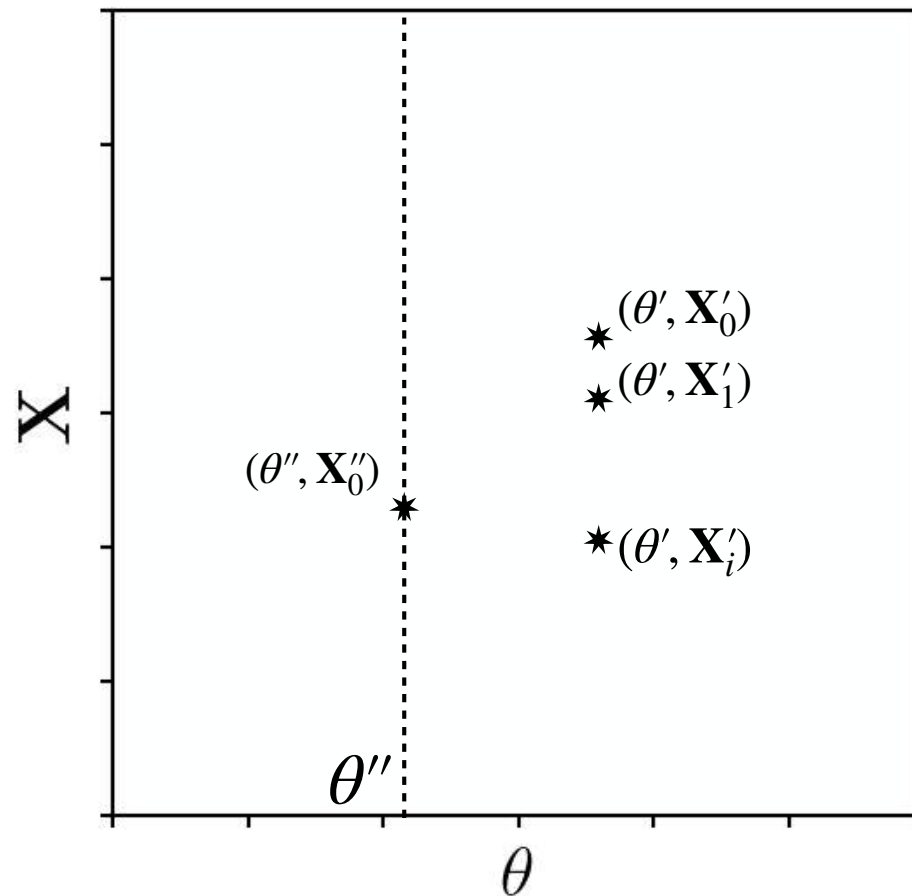
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F must **include noise model**

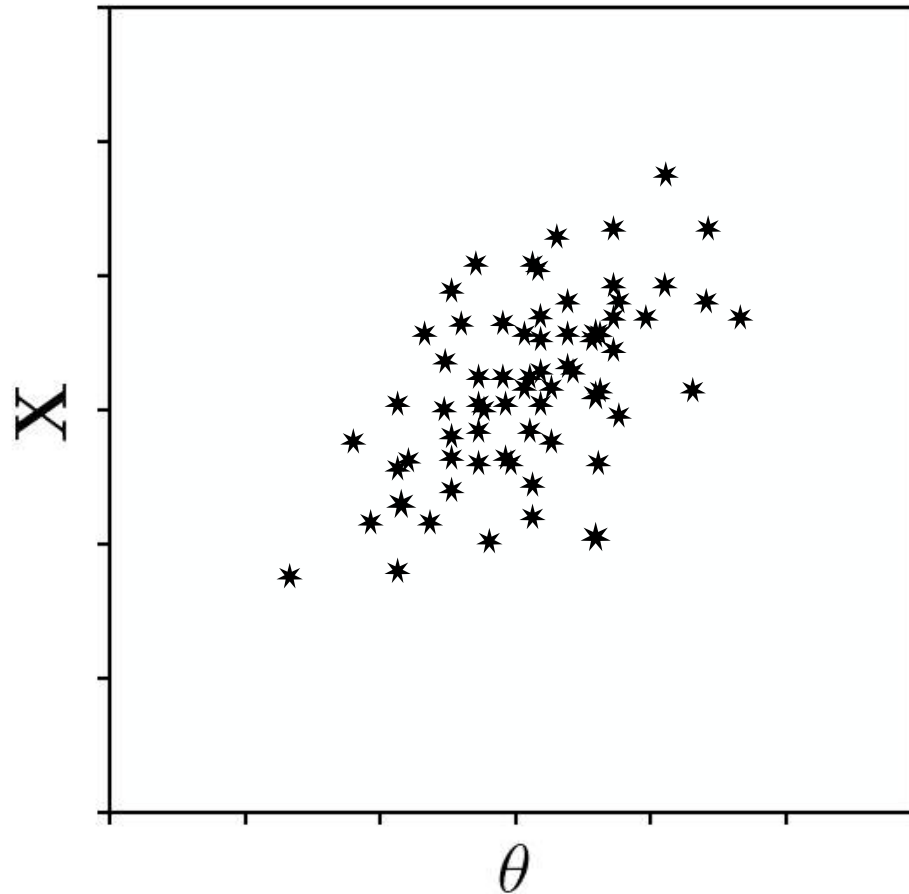
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



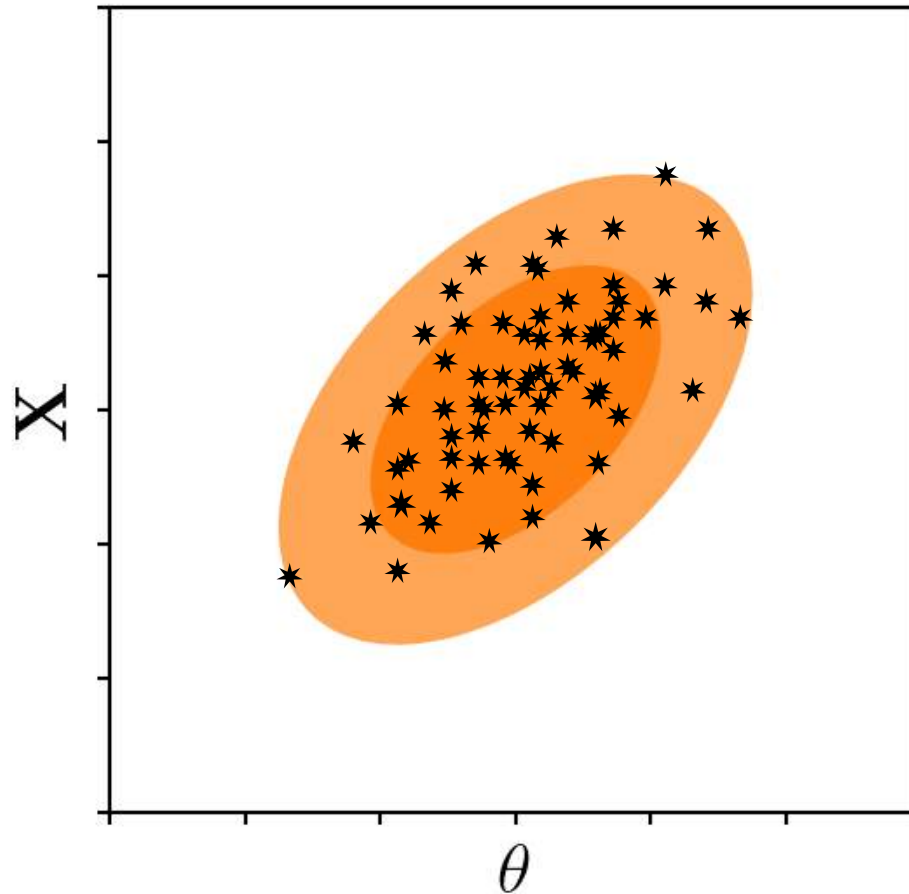
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



forward model F that can simulate observations:

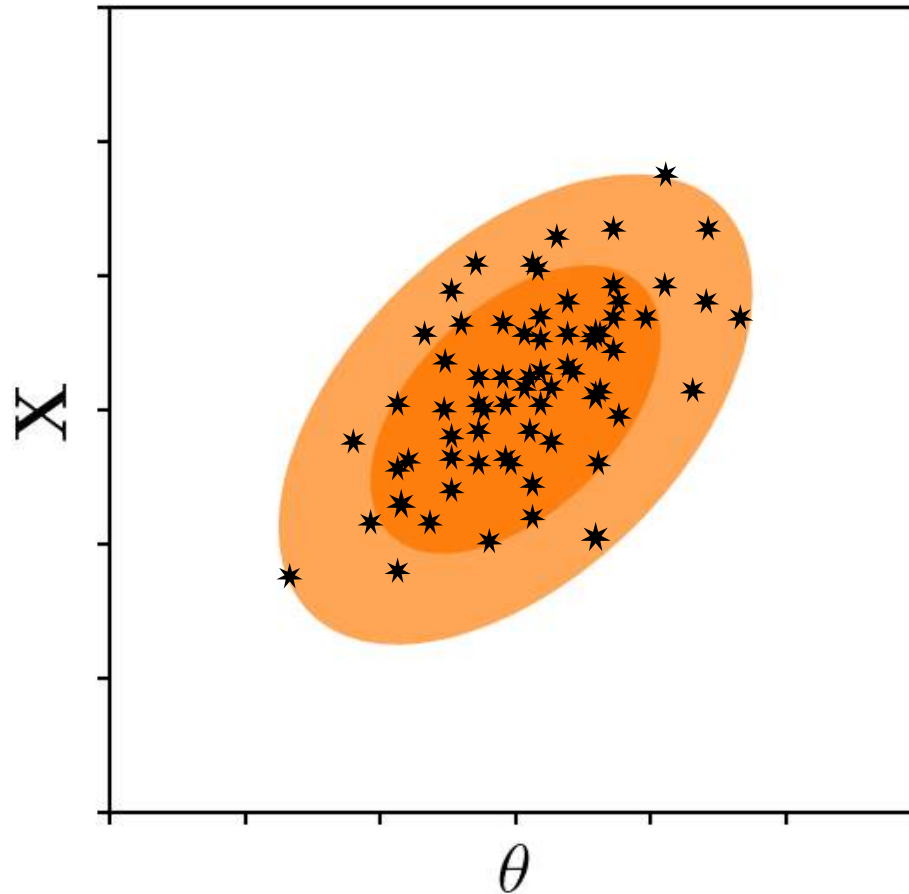
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F allows us to sample $(\theta', X') \sim p(\theta, X)$ \longleftarrow Unknown

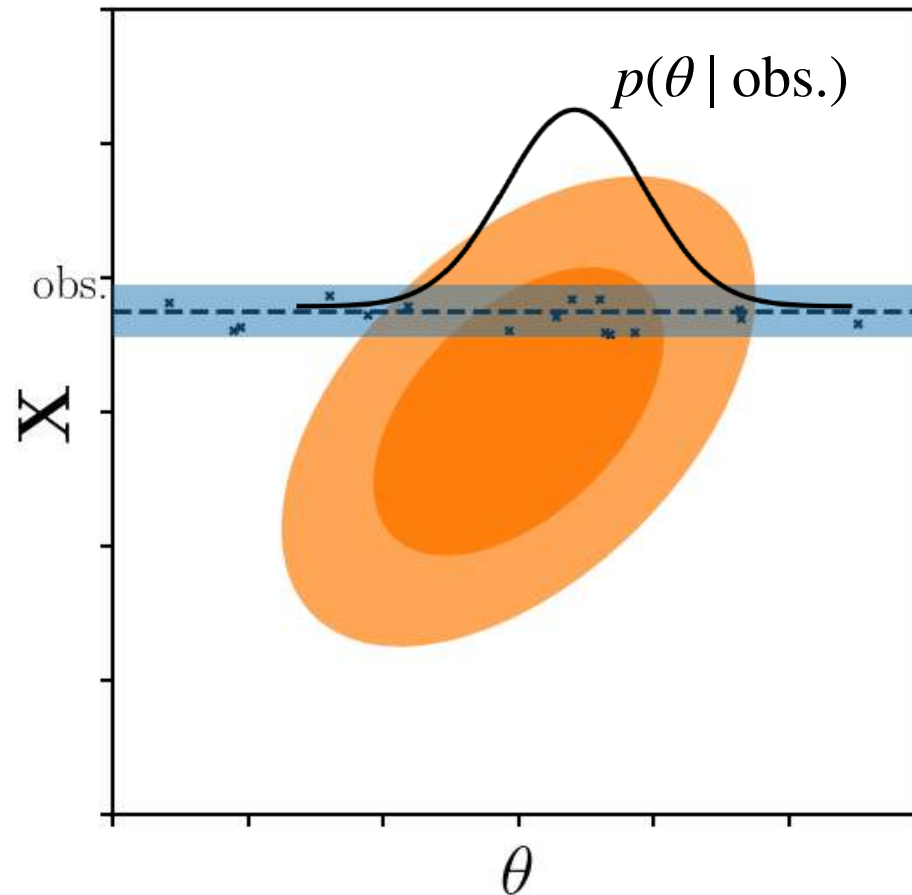
forward model F that can simulate observations:

$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



F allows us to sample $(\theta', X') \sim p(\theta, X)$ \longleftarrow Unknown

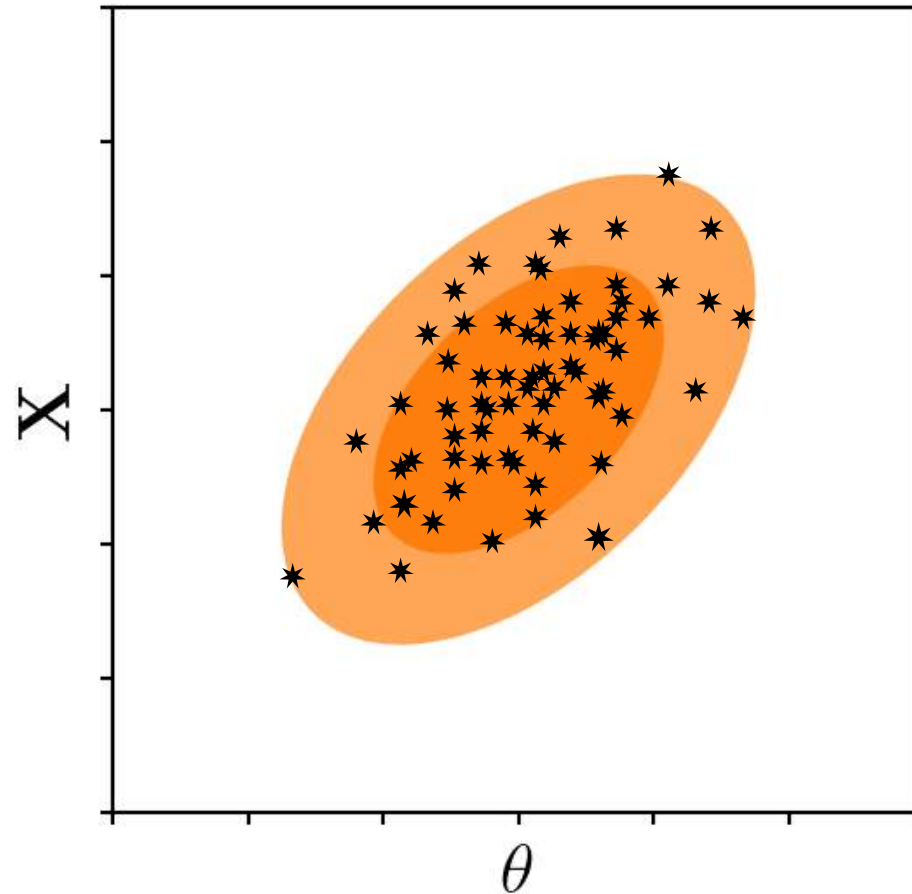
with F one way to infer the posterior is through *brute force* — **approximate bayesian computation (ABC)**



only keep the simulations “close” to the observation

(This is called simulation-based inference, likelihood-free inference or implicit inference)

simulation-based inference is a **density estimation problem!**



F allows us to sample $(\theta', X') \sim p(\theta, X)$ ← Unknown

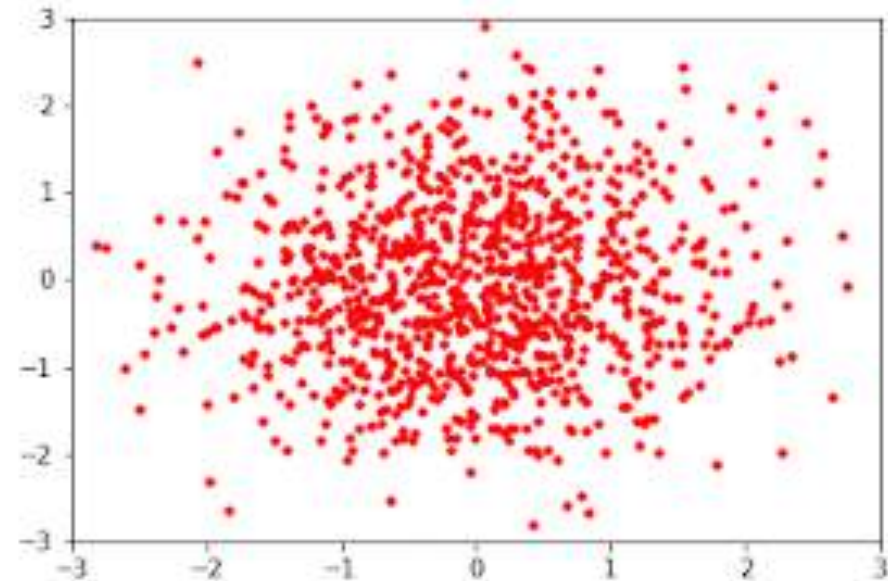
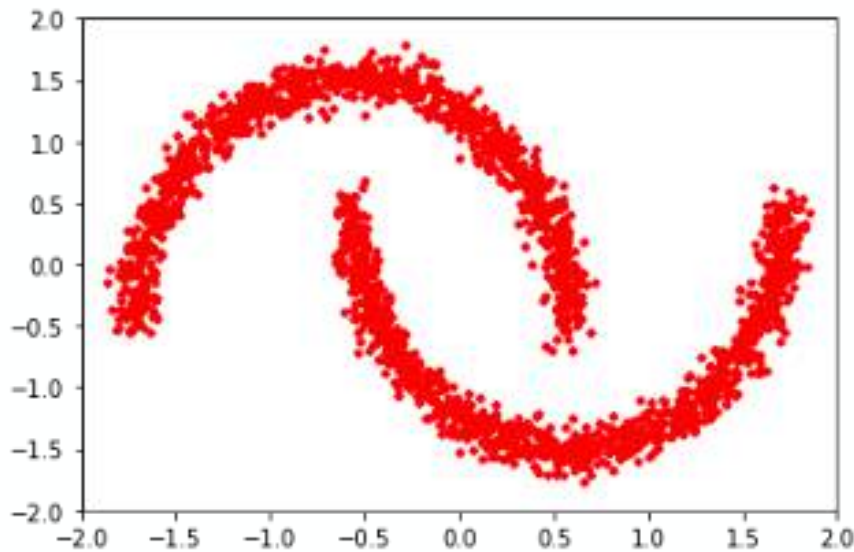
Pixel space!

Pixel space

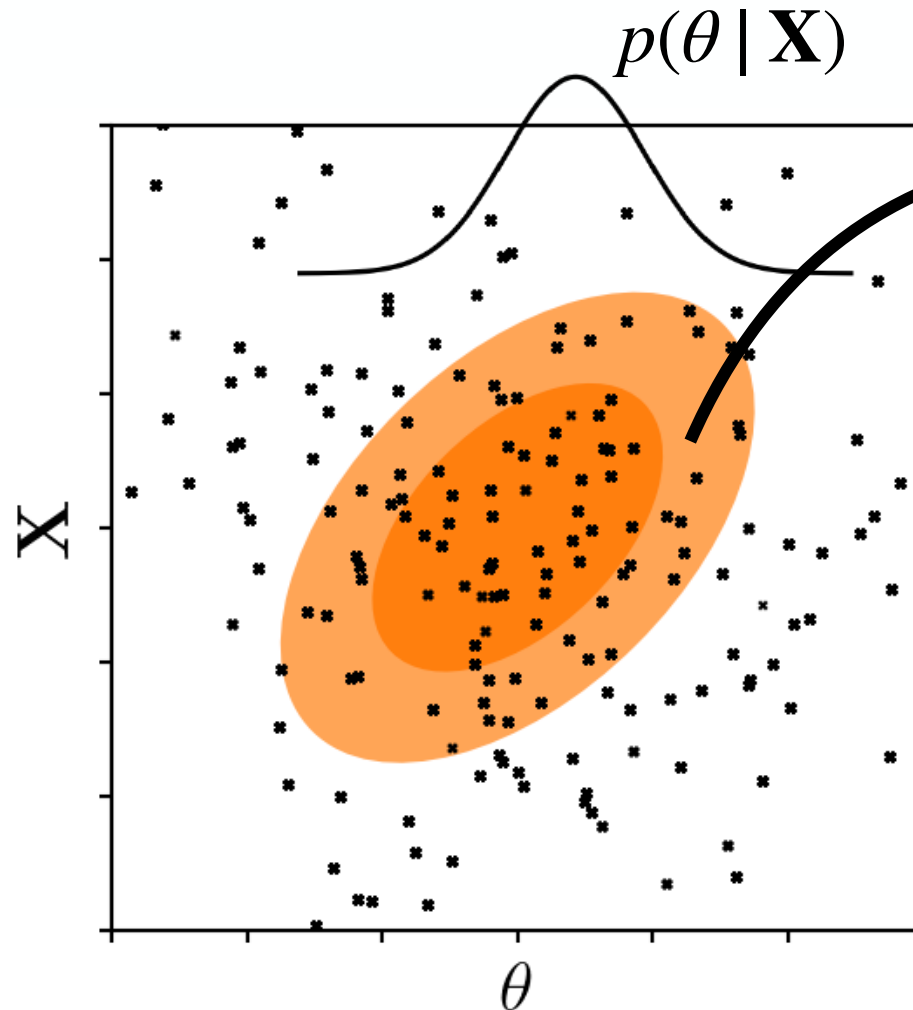
The deep learning revolution
has enabled fast density
estimation at very high
dimension

$$p(\theta | x)$$

Not Gaussian



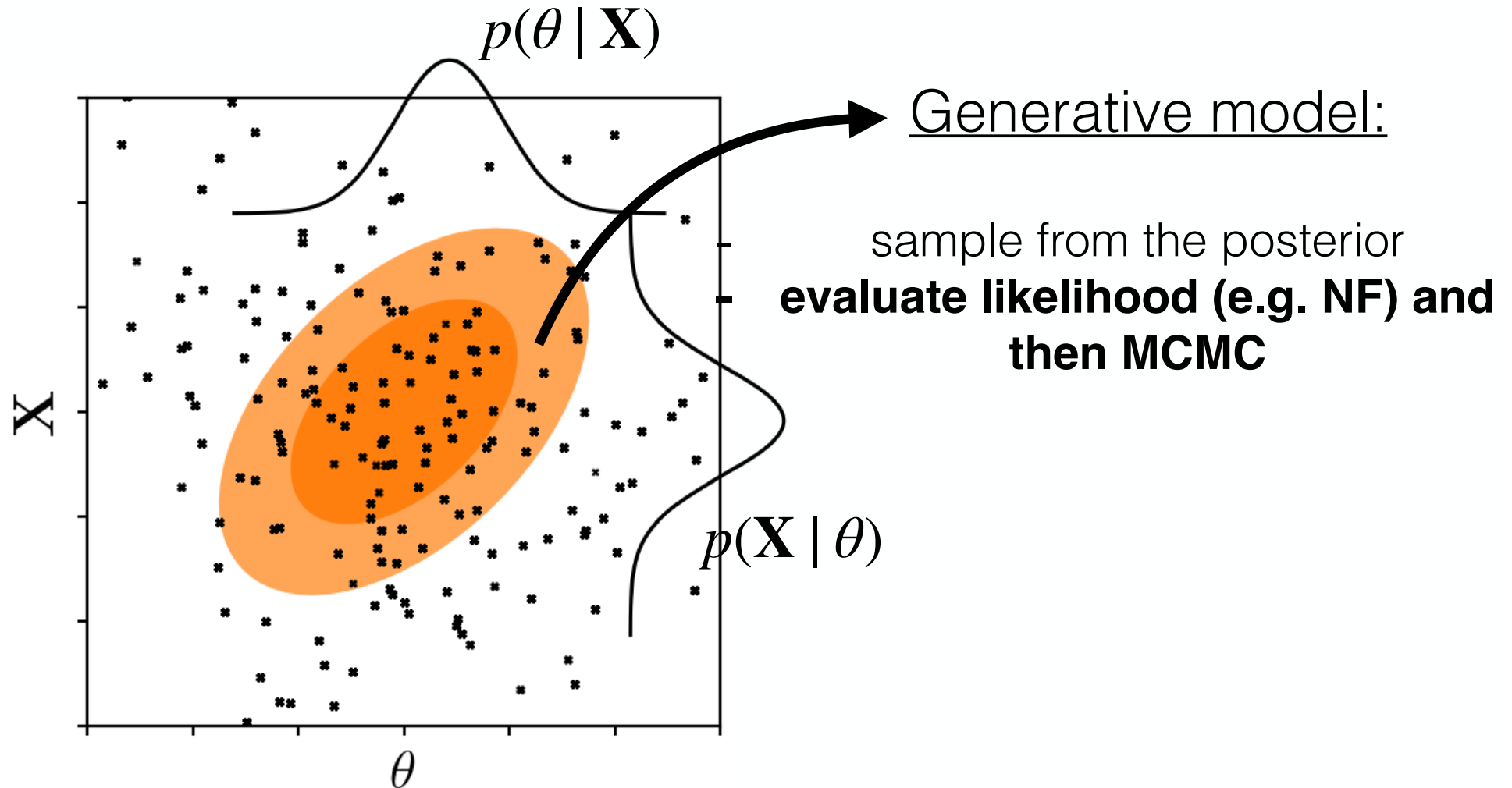
simulation-based inference is a **density estimation problem!**



Generative model:

- sample from the posterior (amortized)

simulation-based inference is a **density estimation problem!**



- Assuming we have a perfect simulator, we have (\mathbf{x}, θ) pairs. How do we do inference?

$$\underbrace{P(\boldsymbol{\theta}|\mathbf{x})}_{\text{"Posterior"}} \propto \underbrace{P(\mathbf{x}|\boldsymbol{\theta})}_{\text{"Likelihood"}} \underbrace{P(\boldsymbol{\theta})}_{\text{"Prior"}}$$

Neural Posterior Estimation

Neural Likelihood Estimation

- Assuming we have a perfect simulator, we have (\mathbf{x}, θ) pairs. How do we do inference?

$$\underbrace{P(\boldsymbol{\theta}|\mathbf{x})}_{\text{"Posterior"}} \propto \underbrace{P(\mathbf{x}|\boldsymbol{\theta})}_{\text{"Likelihood"}} \underbrace{P(\boldsymbol{\theta})}_{\text{"Prior"}}$$

Neural Posterior Estimation

Neural Likelihood Estimation

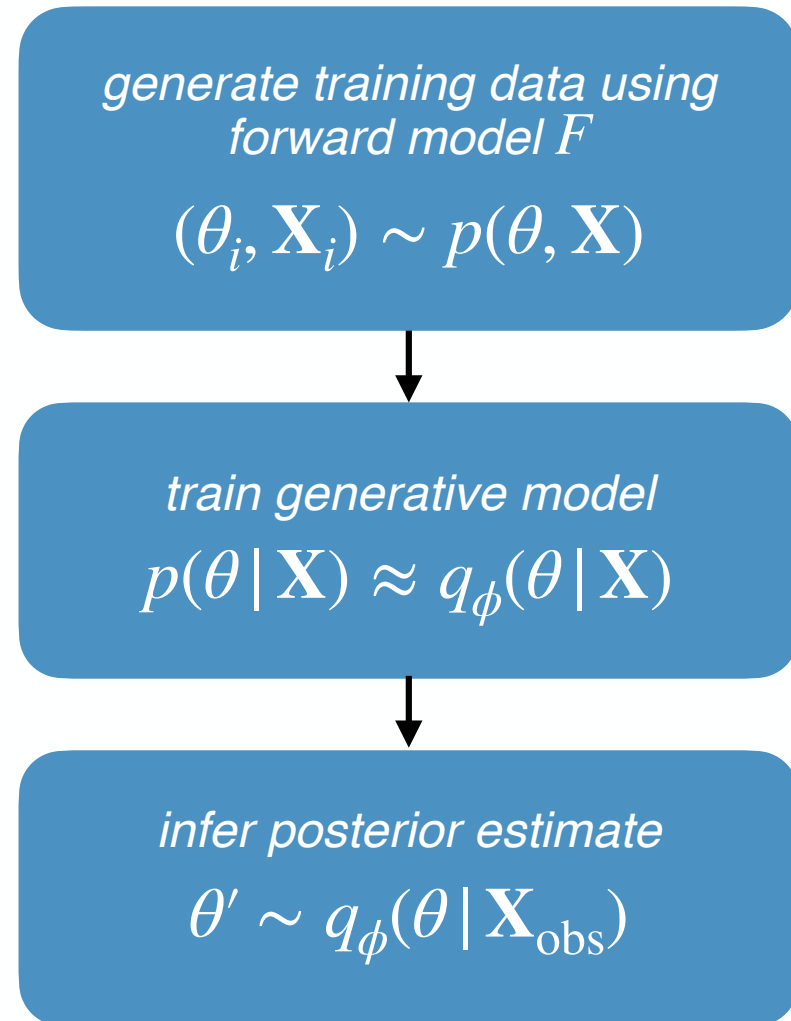
Amortized

Requires sampling (e.g. MCMC)

(Need to sample many posteriors for many \mathbf{x} 's)

(Need to build model into hierarchical sampling scheme)

simulation-based inference flowchart for amortized inference



simulation-based inference flowchart for amortized inference

assumes
the forward model
perfectly matches reality



*generate training data using
forward model F*

$$(\theta_i, \mathbf{X}_i) \sim p(\theta, \mathbf{X})$$



train generative model

$$p(\theta | \mathbf{X}) \approx q_\phi(\theta | \mathbf{X})$$



infer posterior estimate

$$\theta' \sim q_\phi(\theta | \mathbf{X}_{\text{obs}})$$

simulation-based inference flowchart for amortized inference

assumes
the forward model
perfectly matches reality



*generate training data using
forward model F*

$$(\theta_i, \mathbf{X}_i) \sim p(\theta, \mathbf{X})$$



train generative model

$$p(\theta | \mathbf{X}) \approx q_\phi(\theta | \mathbf{X})$$



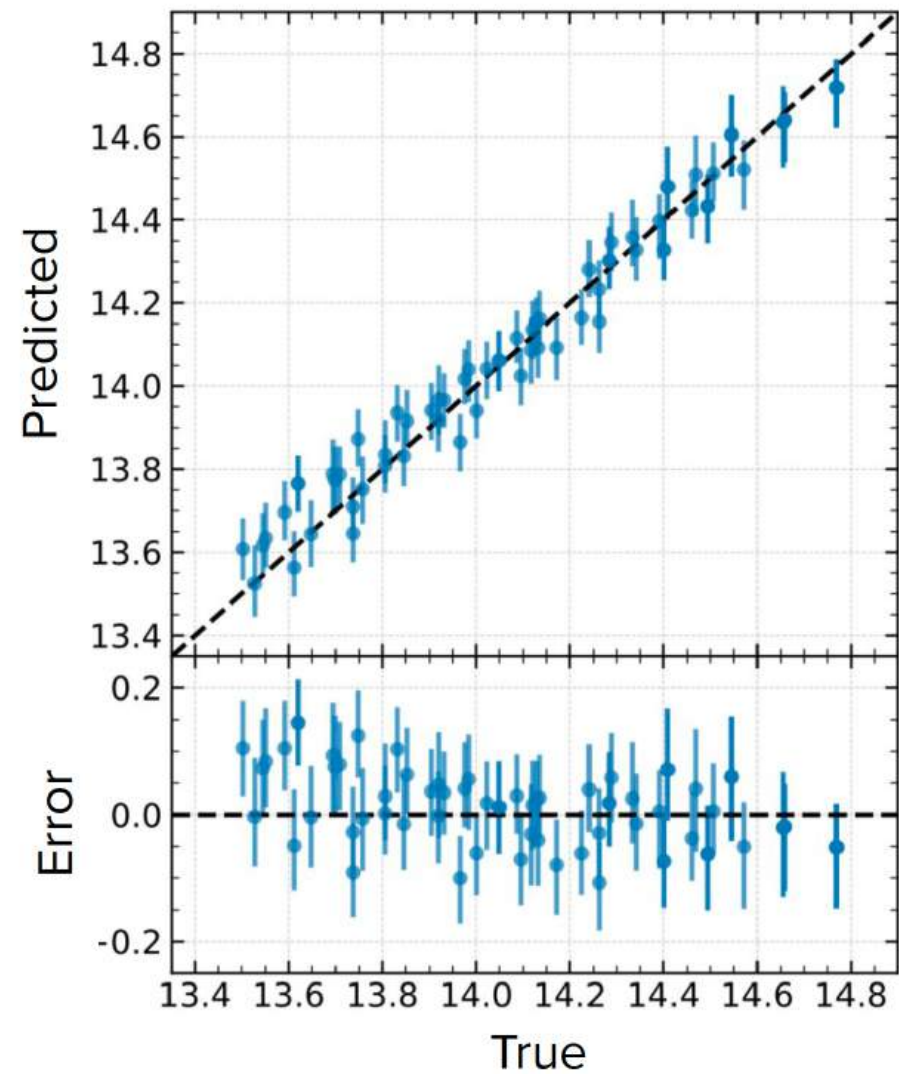
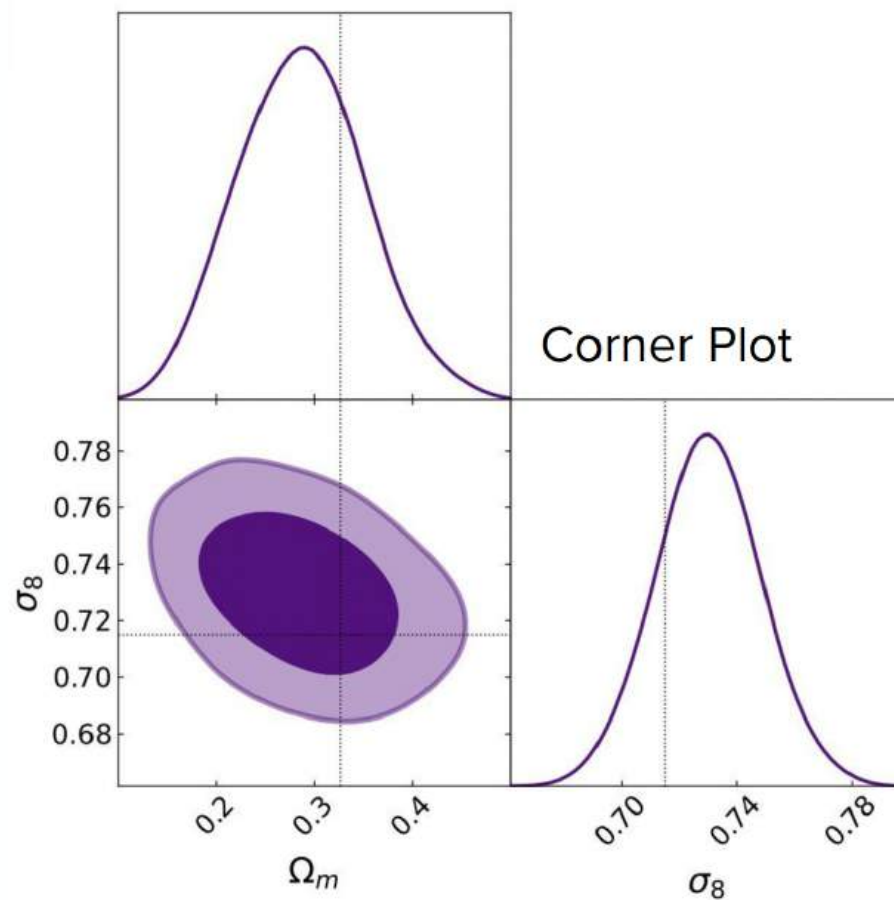
infer posterior estimate

$$\theta' \sim q_\phi(\theta | \mathbf{X}_{\text{obs}})$$

**no guarantee it converges
to the true posterior**



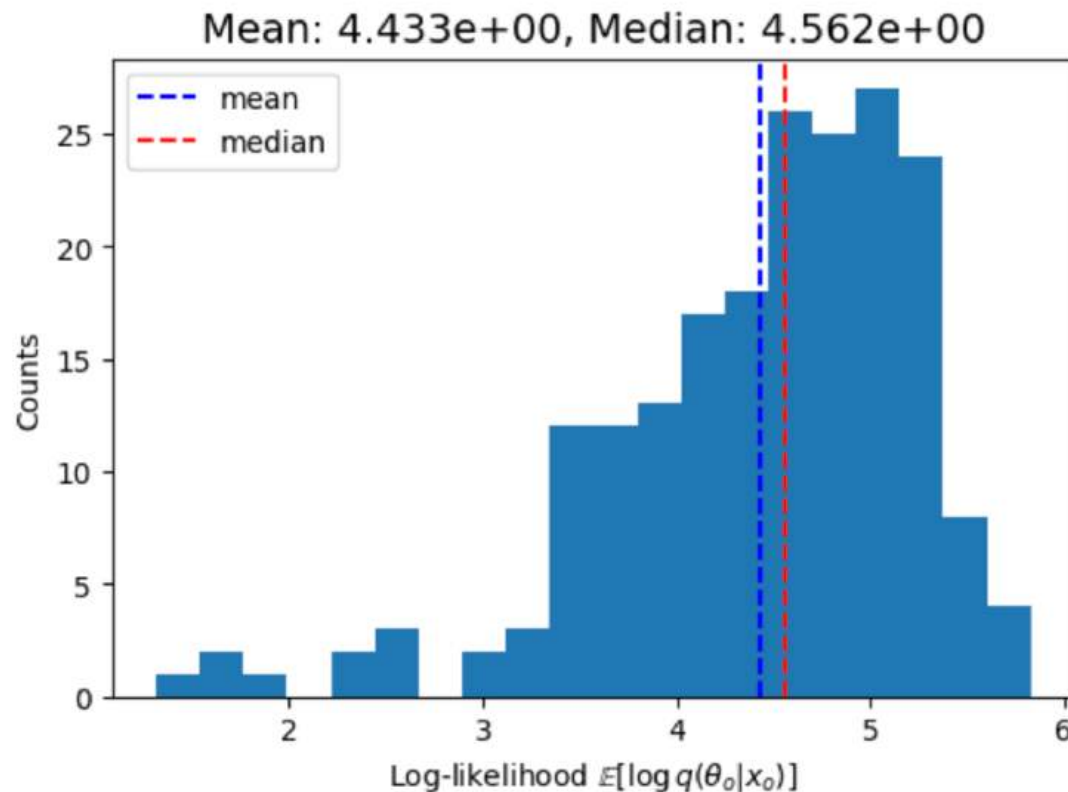
Validation



*adapted from M. Ho (SBI conference)

Cumulative Likelihood of the Test Dataset

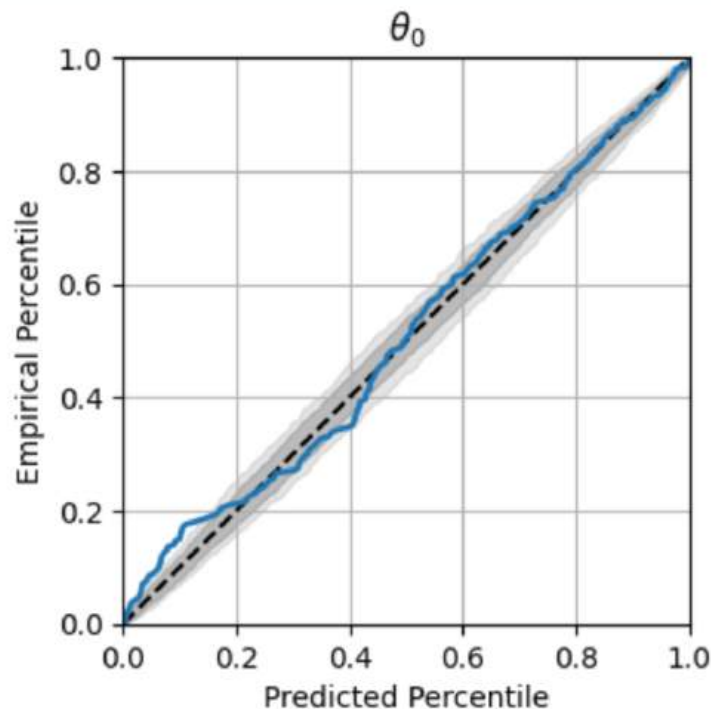
$$\prod_{i=1}^{N_{test}} \hat{P}(\theta_i | x_i)$$



Large values
mean good
predictive power

P-P plots

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



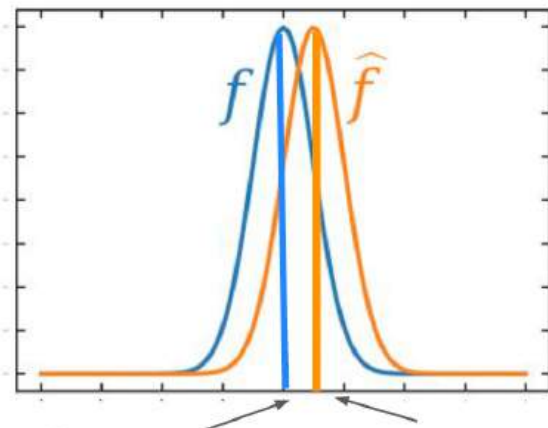
“What percentile level is my posterior model assigning to the true value, and with what frequency does this occur in the test set, e.g. **we should predict the true value below the 50-th percentile 50% of the time**”

Ho+24

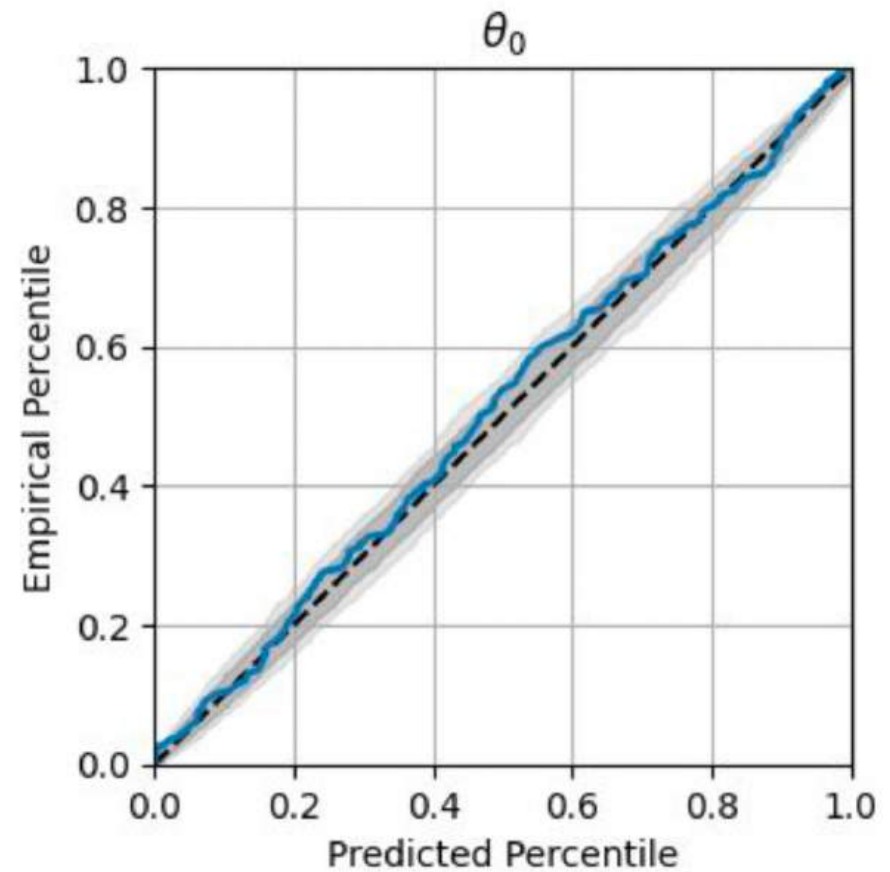
Validation

PIT tests / P-P plots

- Tests of marginal coverage



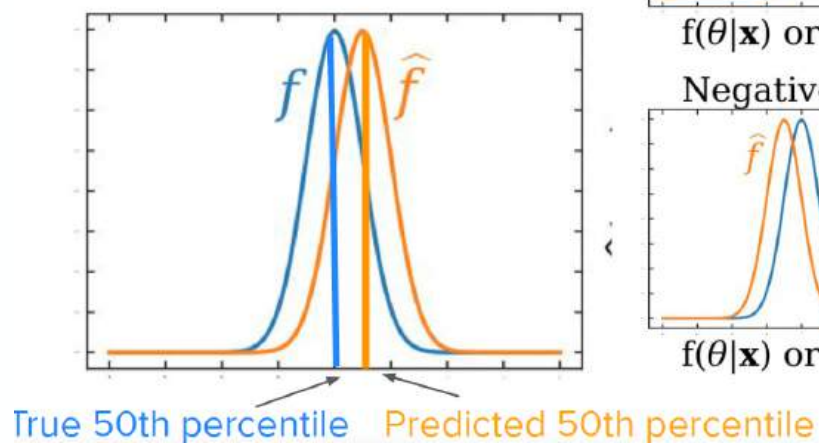
True 50th percentile Predicted 50th percentile



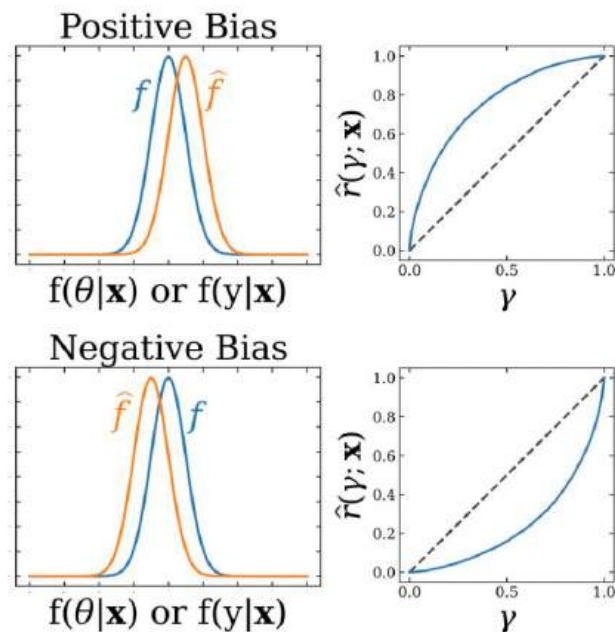
*adapted from M. Ho (SBI conference)

PIT tests / P-P plots

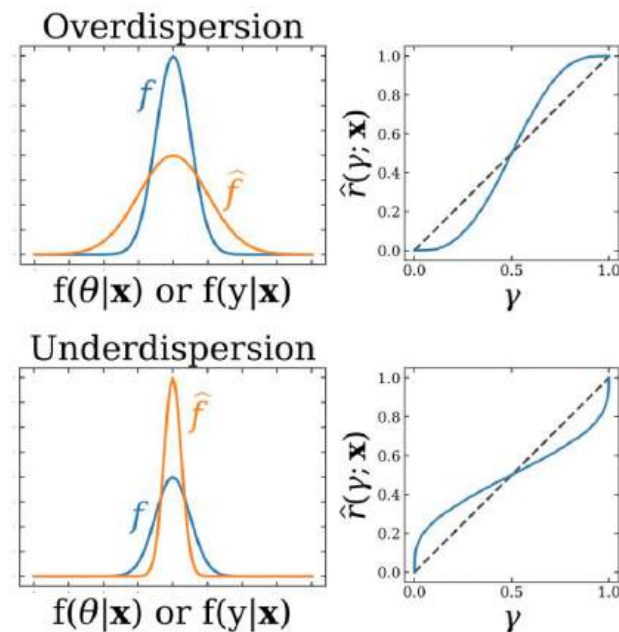
- Tests of marginal coverage



BIAS

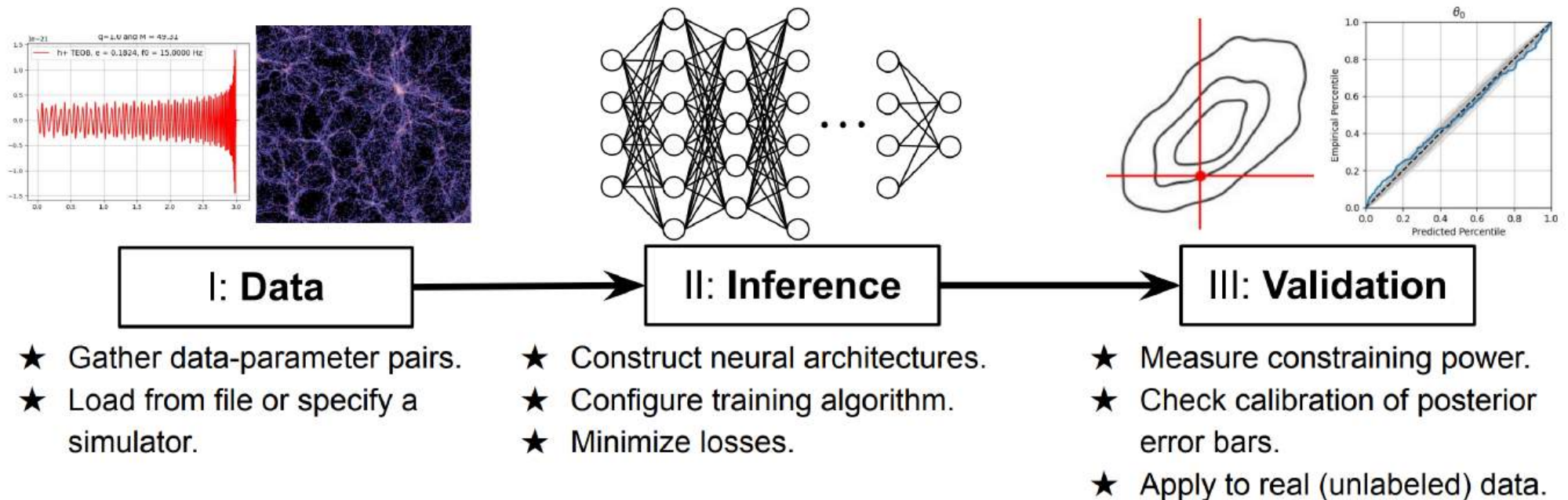


DISPERSION

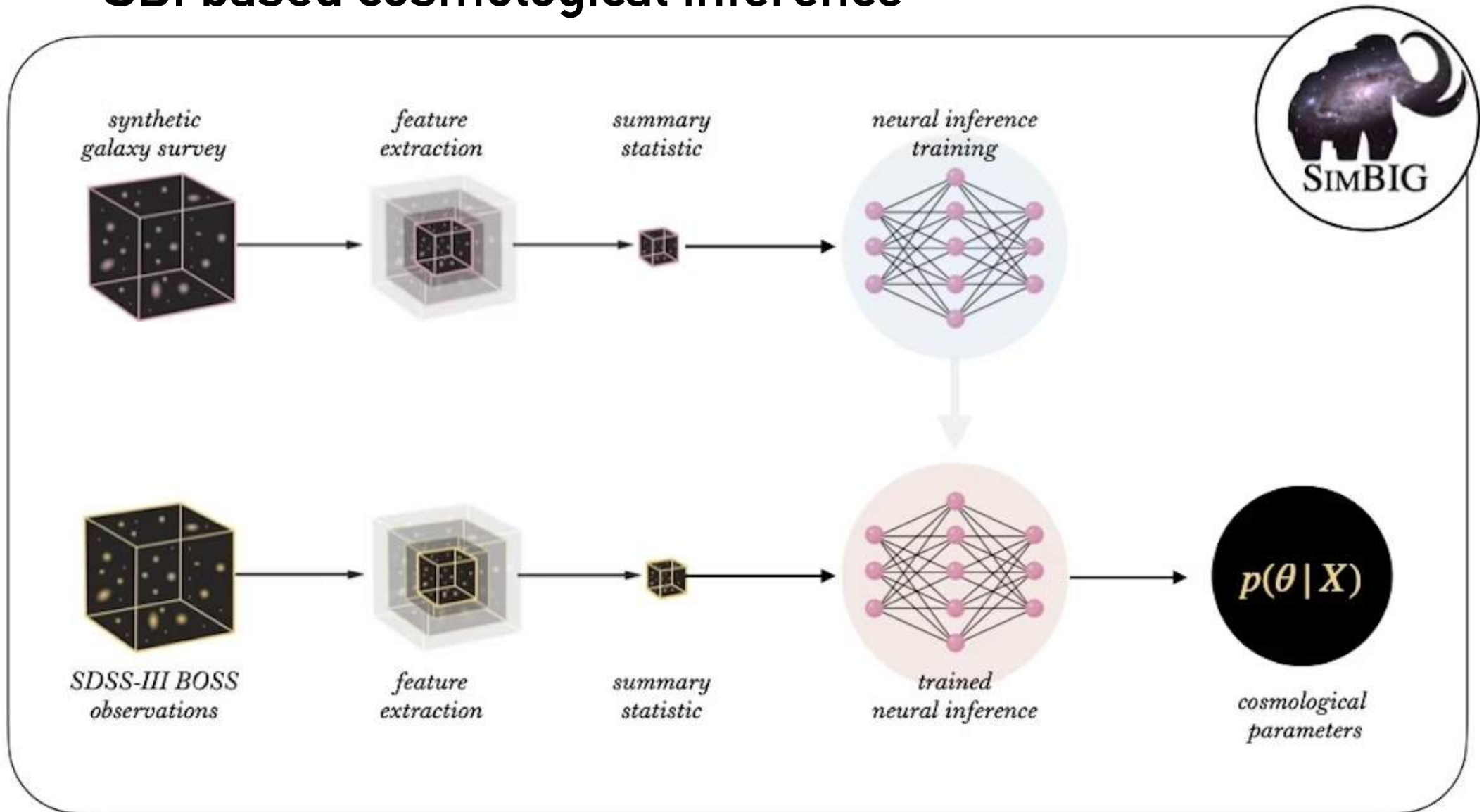


*adapted from M. Ho (SBI conference)

SBI pipeline

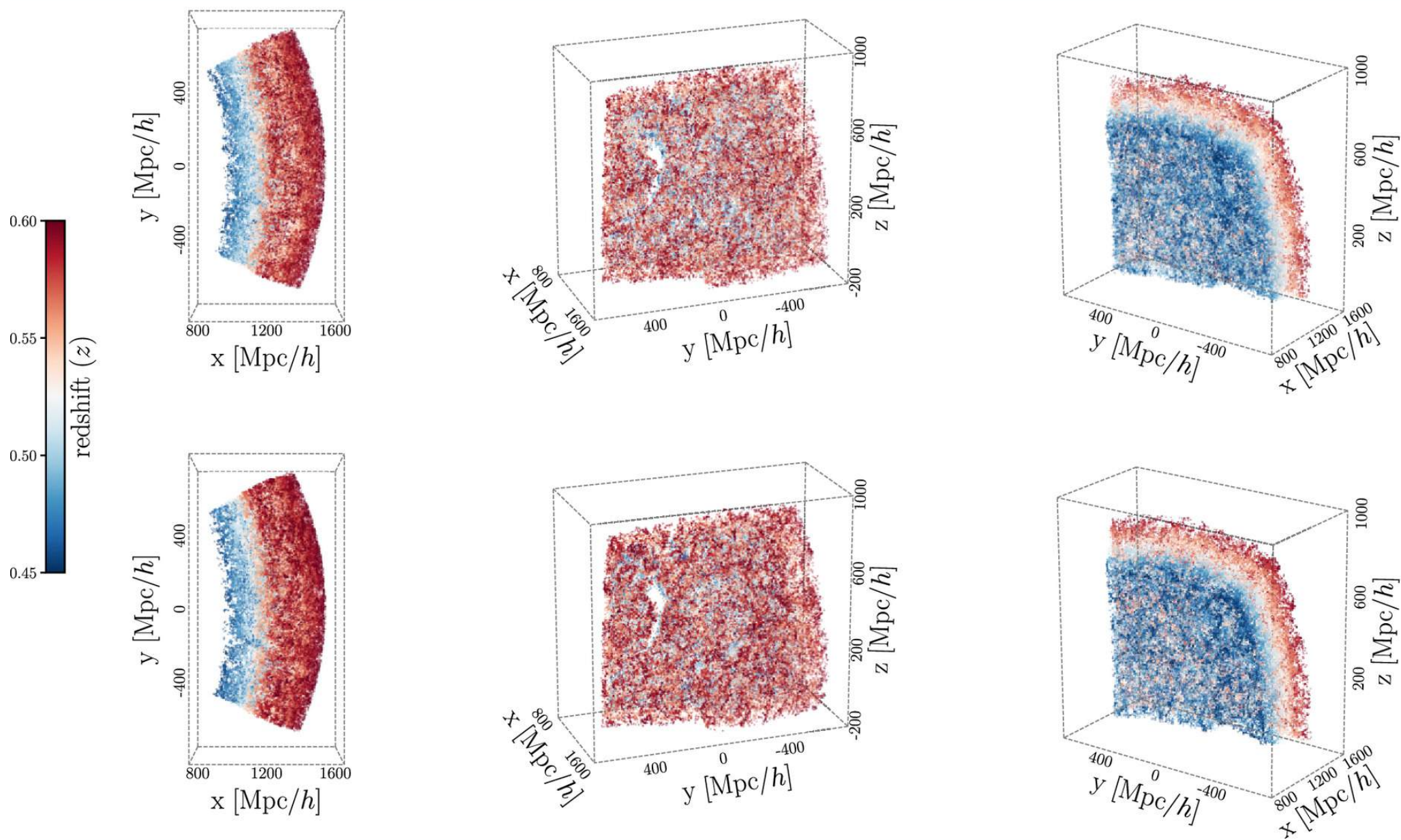


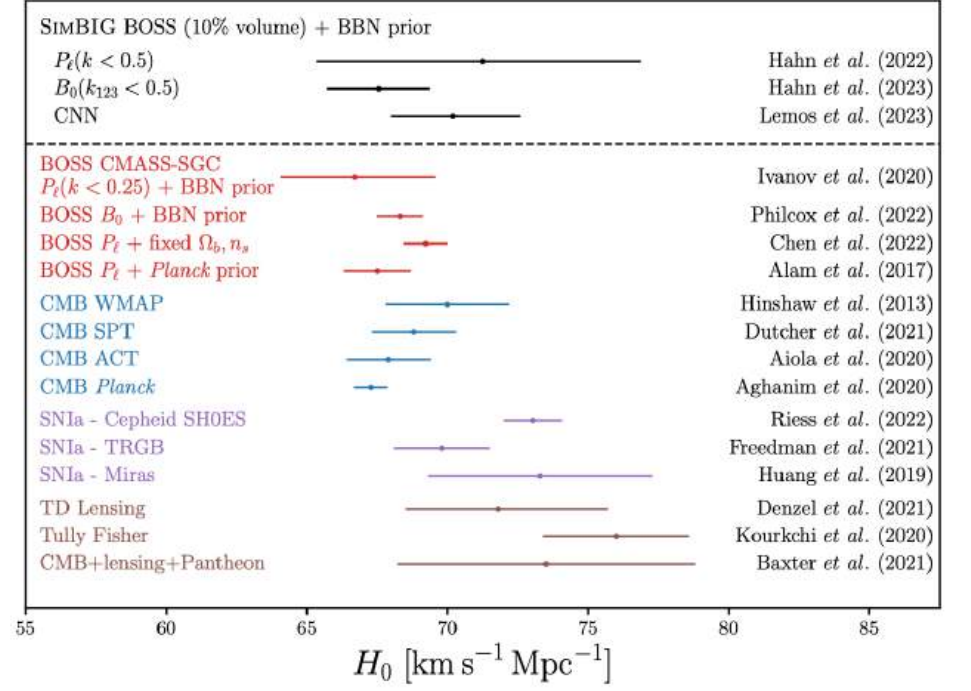
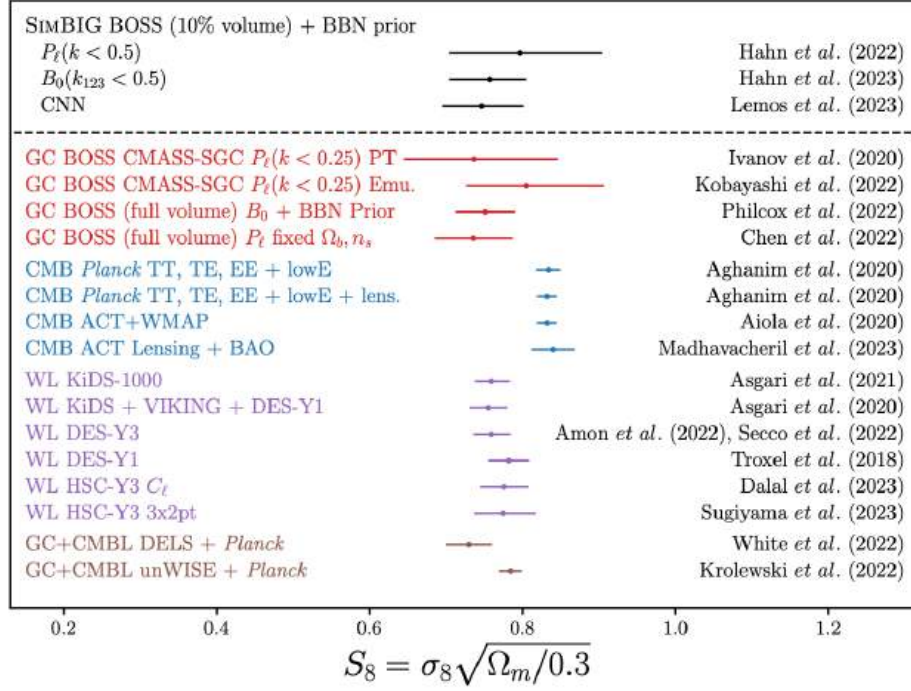
SBI based cosmological inference



“If you can simulate, you can do inference”

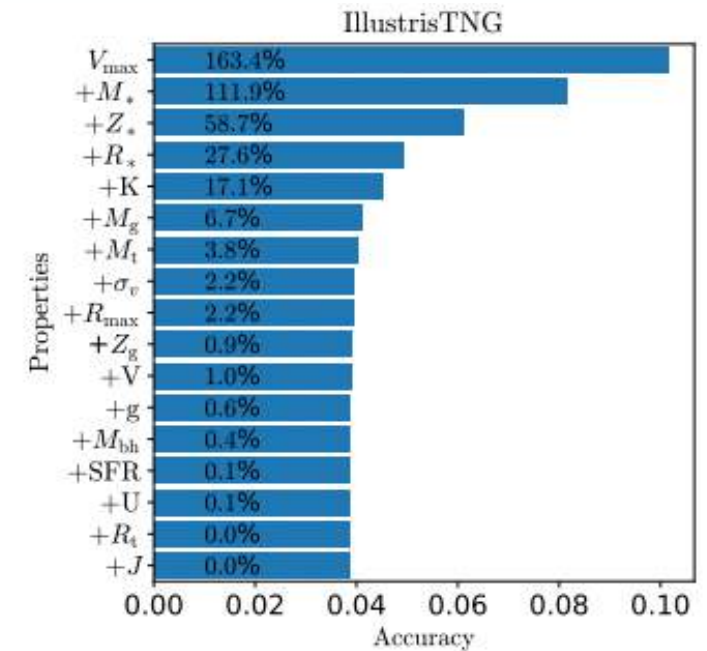
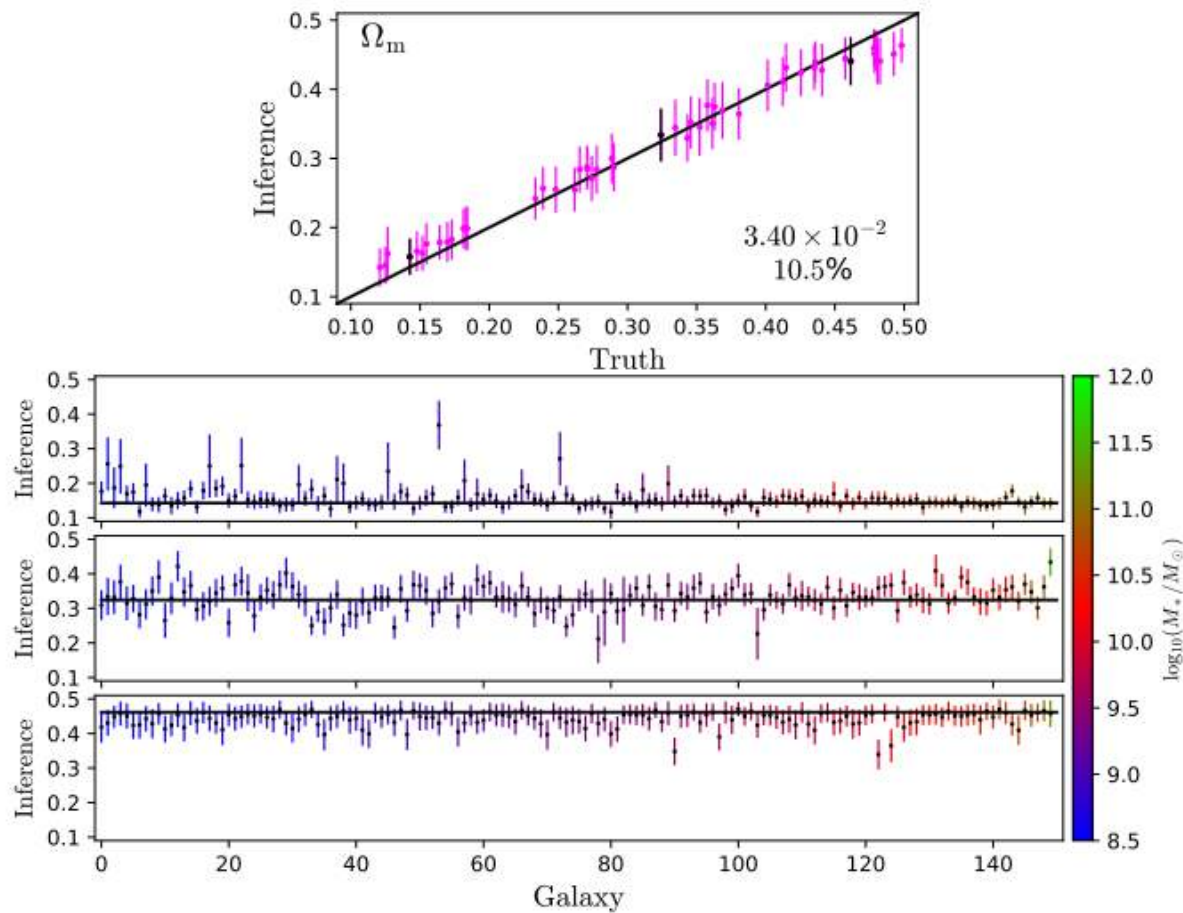
Towards Field Level Cosmological Inference





Hahn+23

Towards Field Level Cosmological Inference



(Cosmology and Astrophysics with Machine Learning Simulations)

CAMELS

Villaescusa-Navarro+22

Extra Material

GANs AND VAEs ARE VERY POWERFUL BUT DO NOT PROVIDE AN EXPLICIT LIKELIHOOD

Method	Train on data	One-pass Sampling	Exact log-likelihood	Free-form Jacobian
Variational Autoencoders	✓	✓	✗	✓
Generative Adversarial Nets	✓	✓	✗	✓
Likelihood-based Autoregressive	✓	✗	✓	✗

Grathwohl+18

* there is a third member of the family now: score based models which try to estimate the gradient of the likelihood (score) instead of the likelihood

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

***refer to the original paper by Song and Ermon 2019
(Much clearer than any blogpost)**

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_θ : Score network, i.e. Neural Network that approximates the gradient of p

One key ingredient is Langevin Dynamics:

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

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

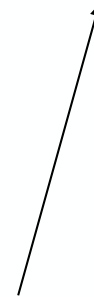
$$\tilde{x}_0 \sim \pi(x)$$

For epsilon small, T large:

$$\tilde{x}_T \quad \text{Is exact sample of} \quad p(x)$$

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

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



Perturbation with gaussian noise

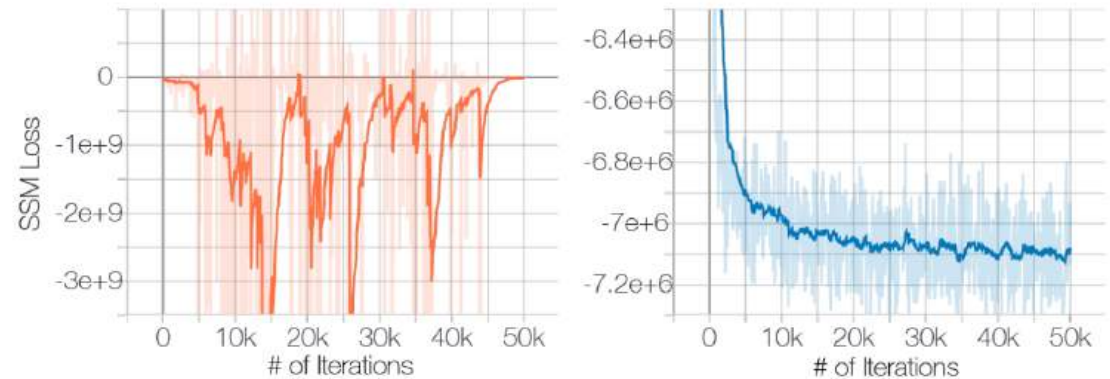
$$s_{\theta}^*(x) = \nabla_x \log q_{\sigma}(x) \simeq \nabla_x \log p_{\text{data}}(x)$$

If noise is small

Estimating the score is difficult

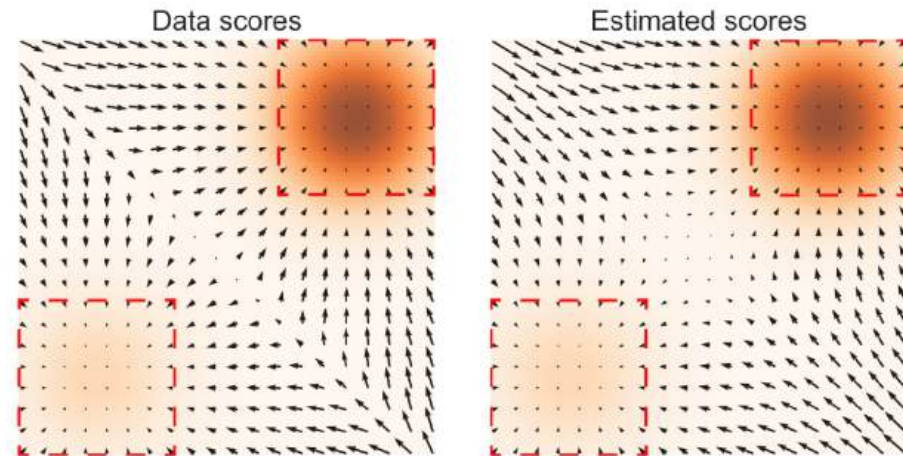
1. Manifold Hypothesis

If the support of x is not the whole space,
score estimation techniques fail



2. Low Density Regions

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



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:

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

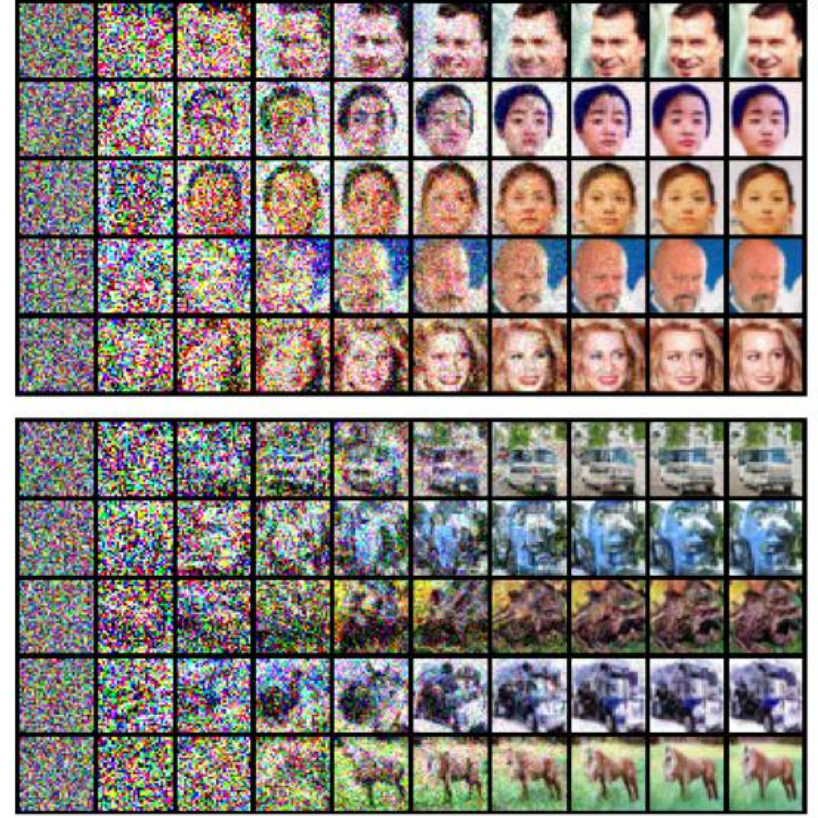
And then combined:

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

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: $\alpha_i \leftarrow \epsilon \cdot \sigma_i^2 / \sigma_L^2$ $\triangleright \alpha_i$ 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 $q_{\sigma_L}(\mathbf{x})$.