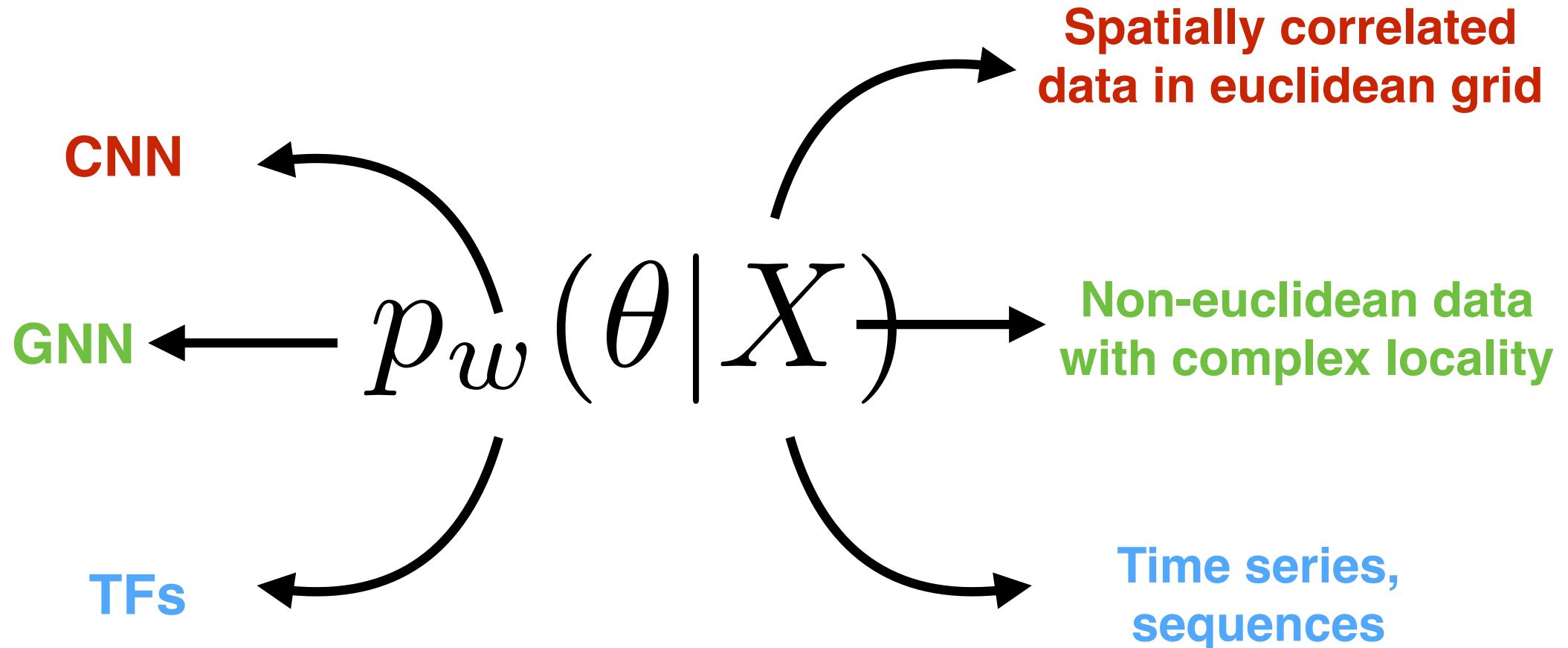
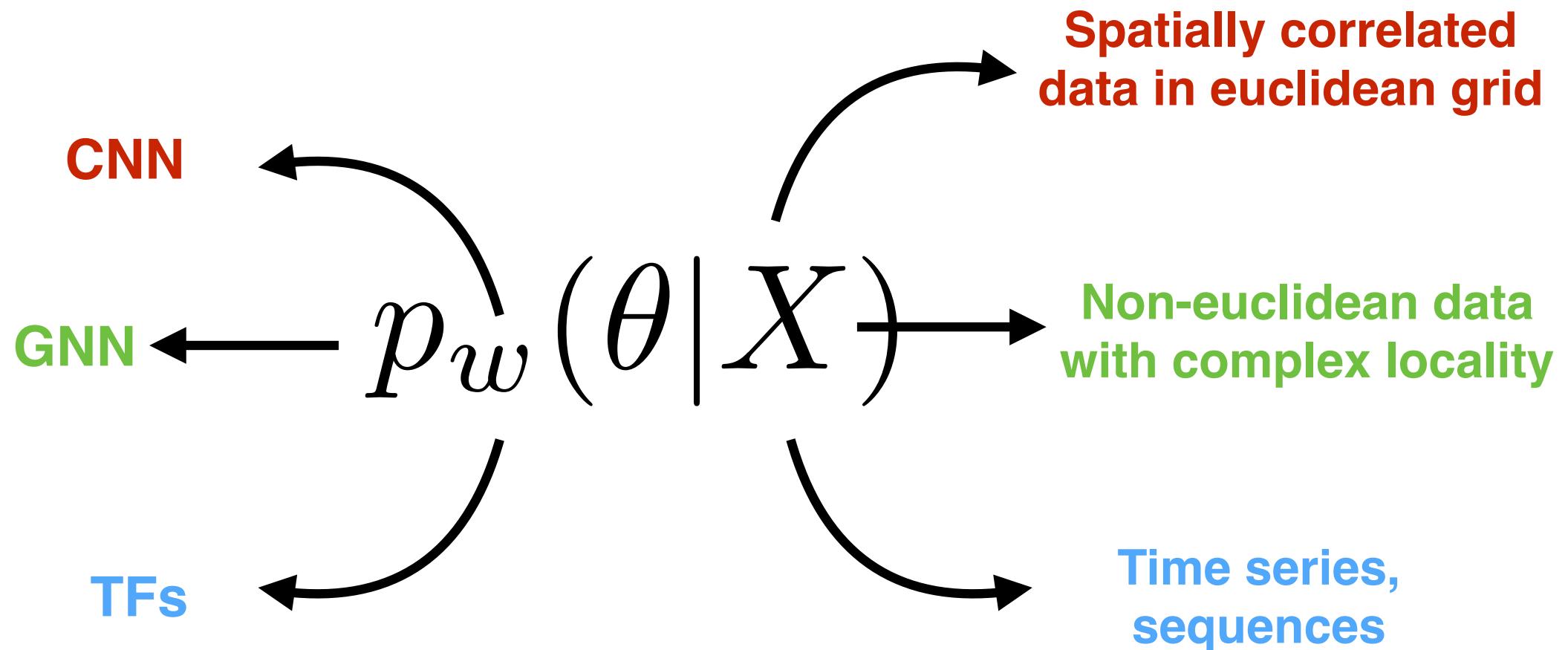


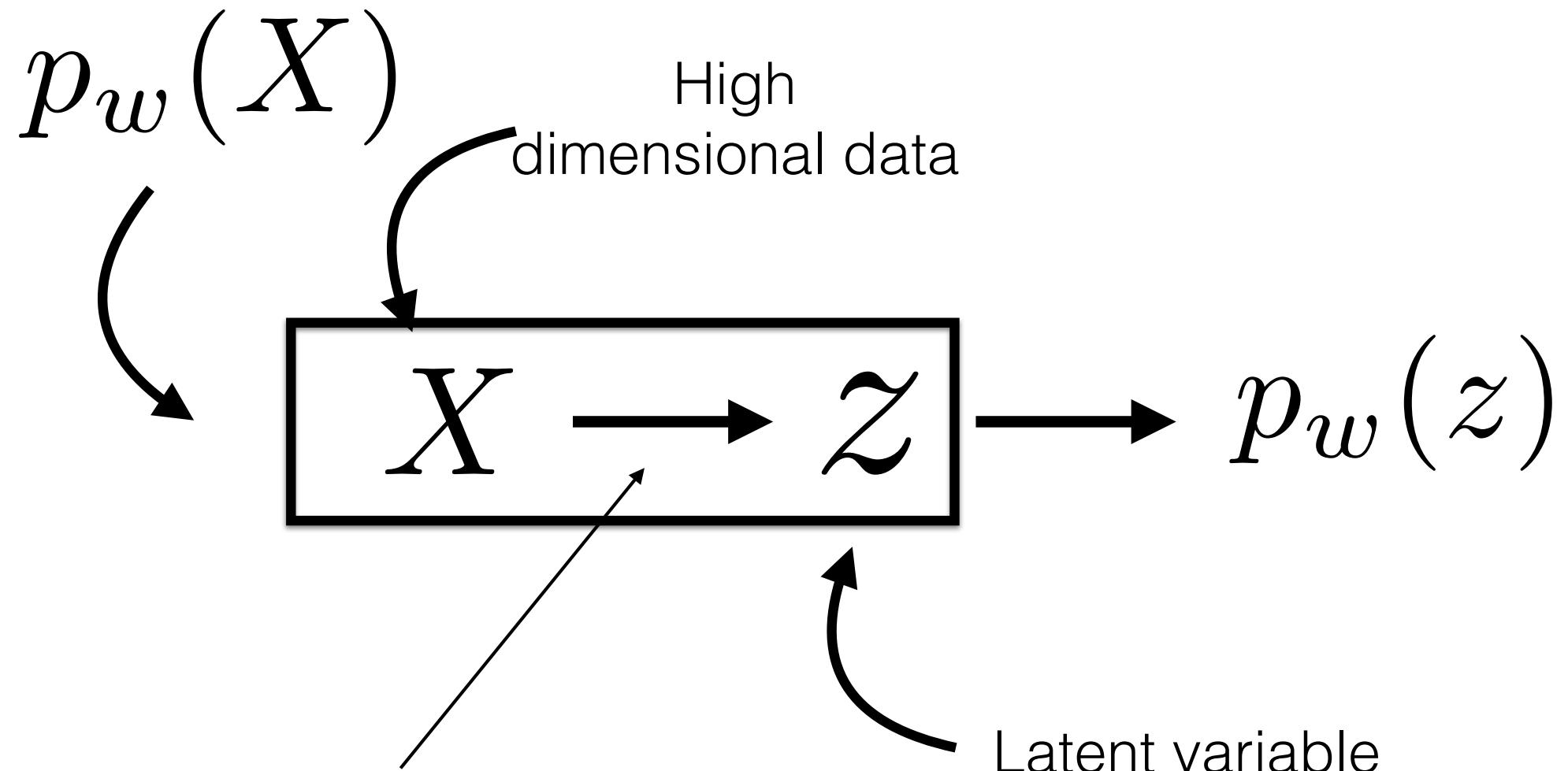
## RECAP5a:



## RECAP5b:

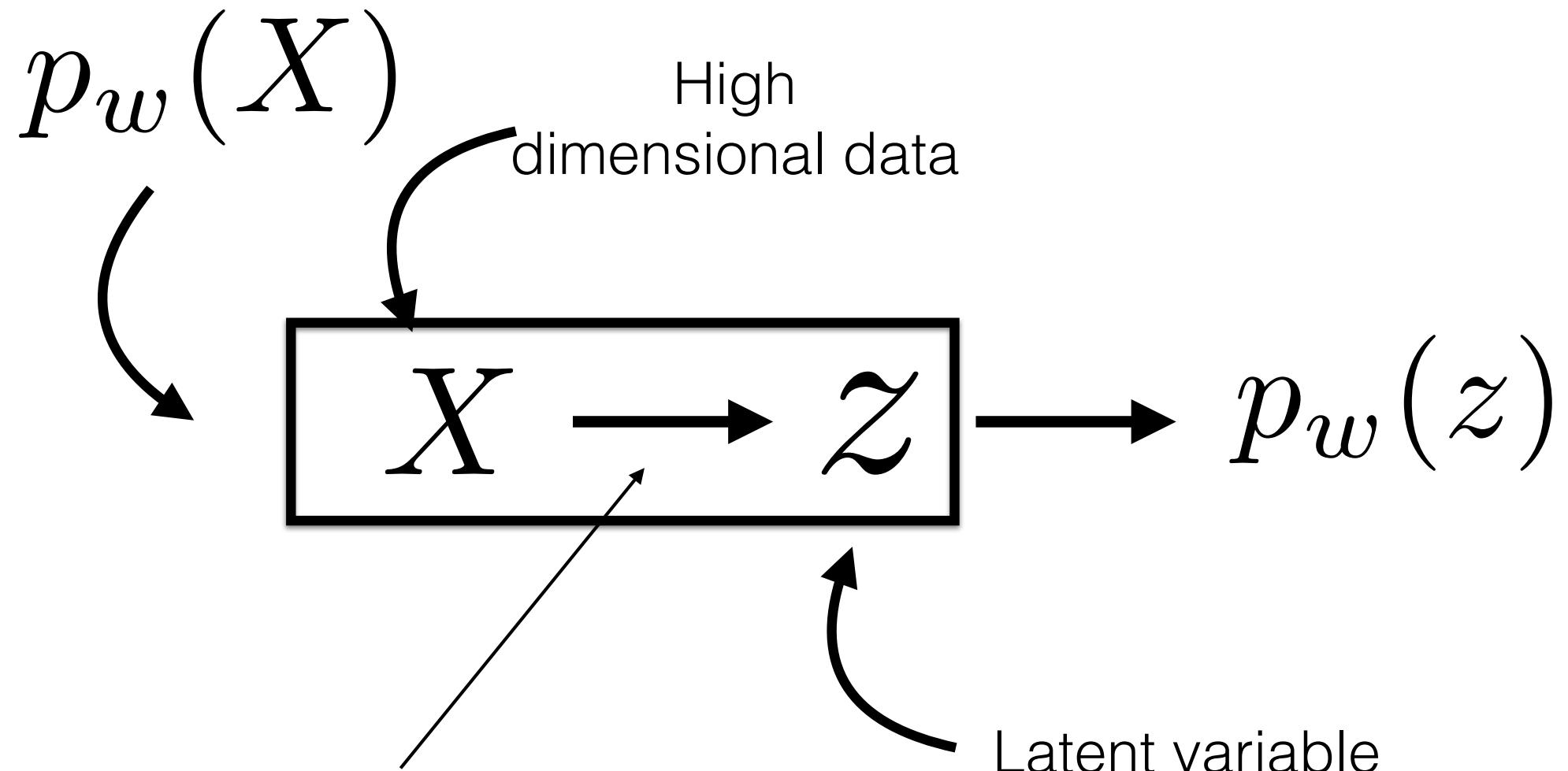


## RECAP5b:



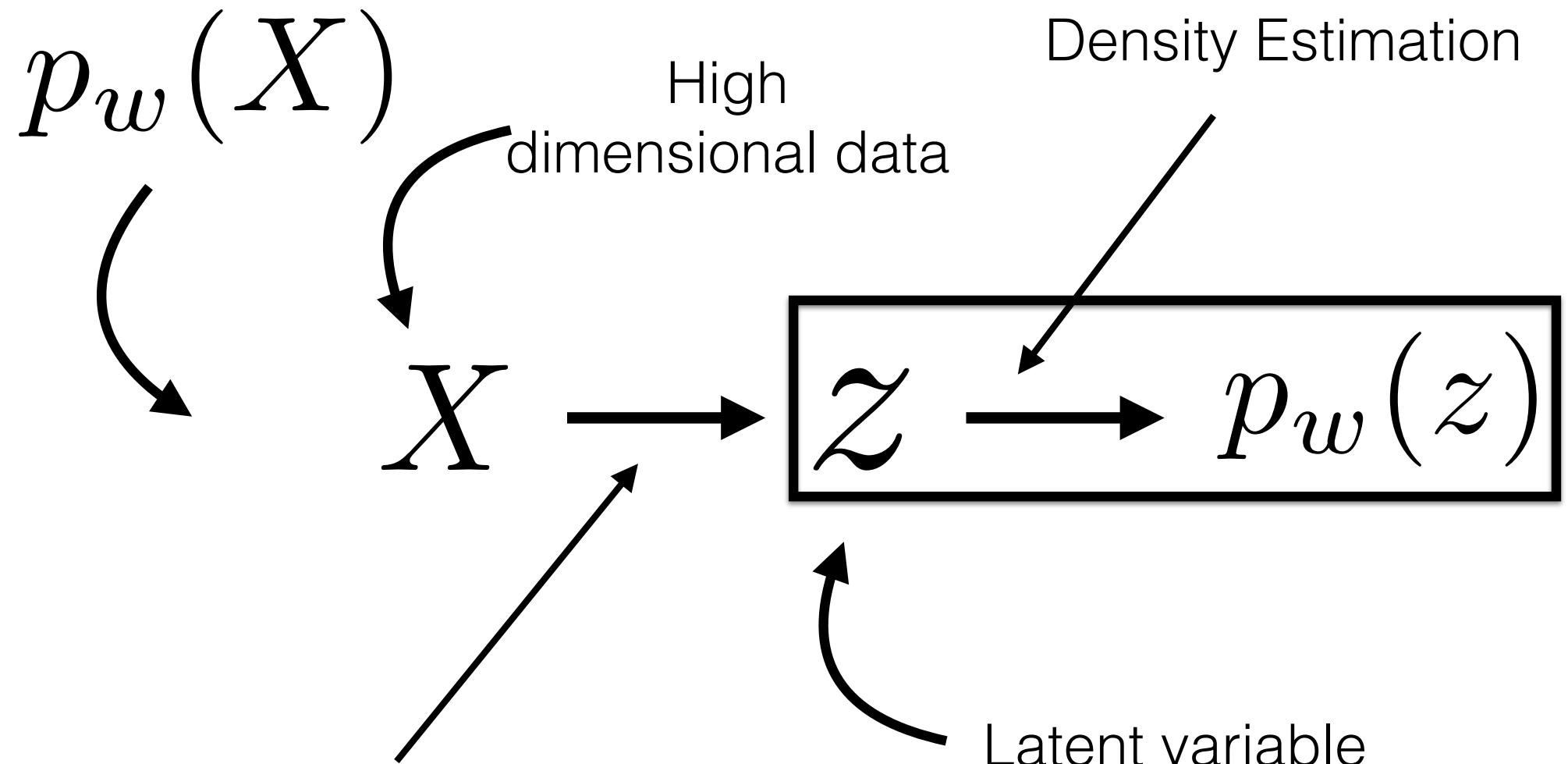
Representation learning /  
dimensionality reduction

## RECAP5b:



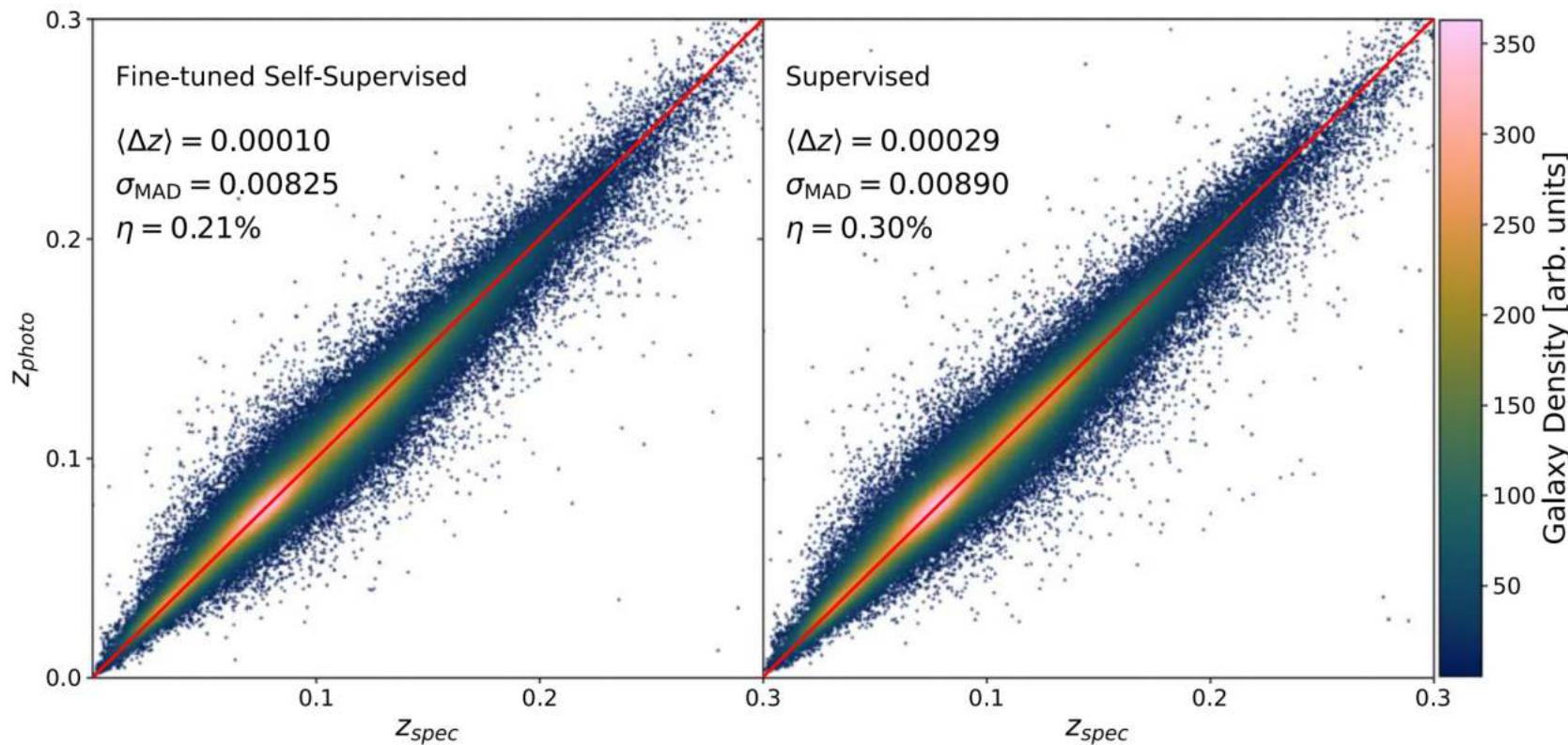
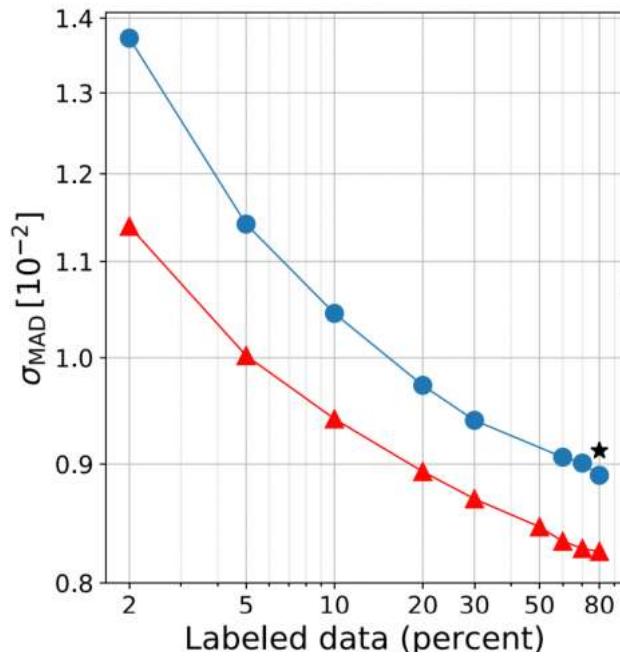
Representation learning /  
dimensionality reduction

## RECAP5b:

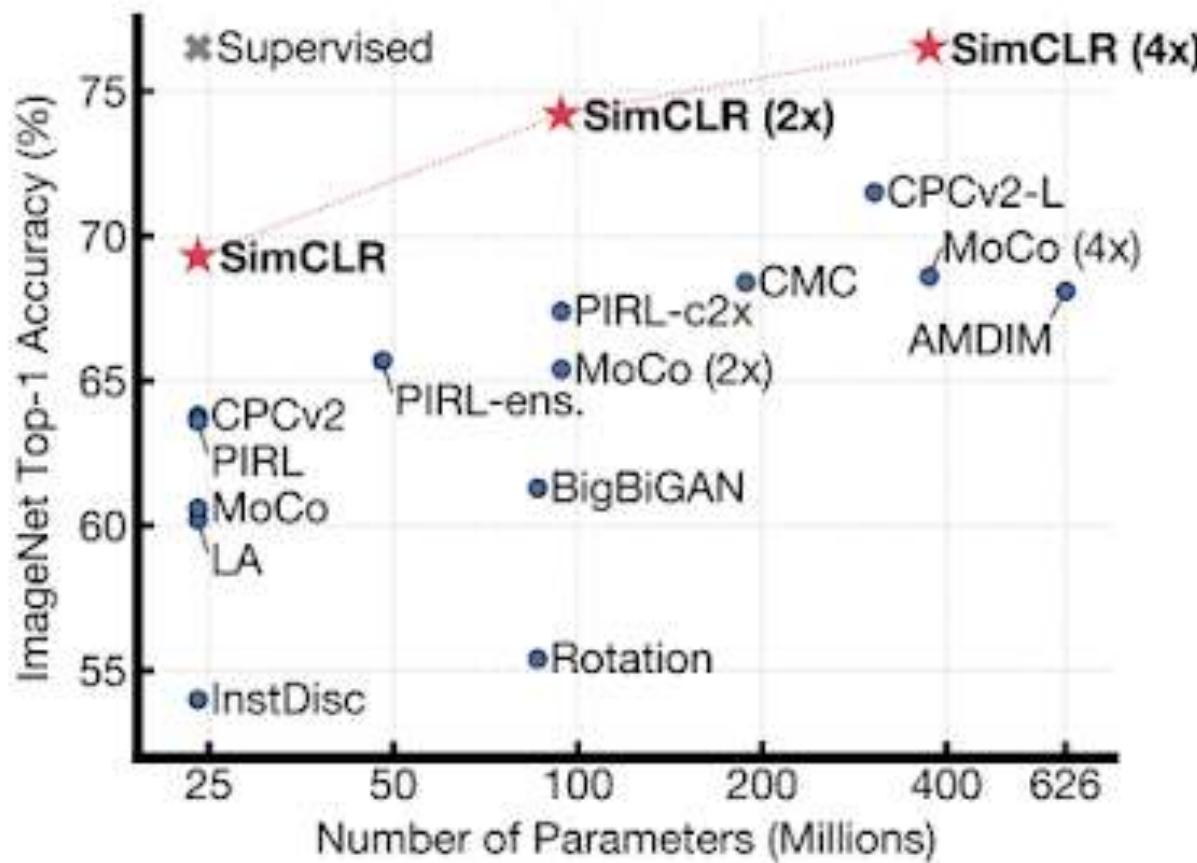


Representation learning /  
dimensionality reduction

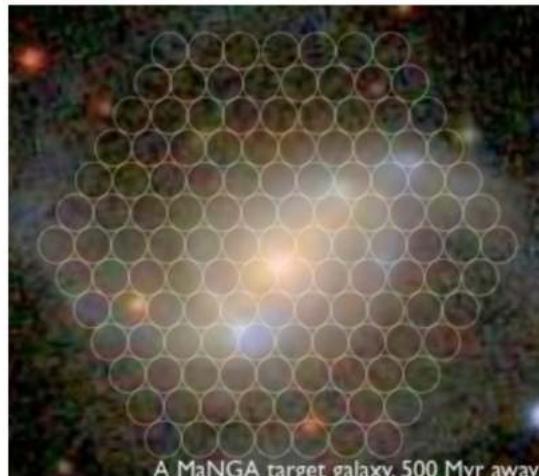
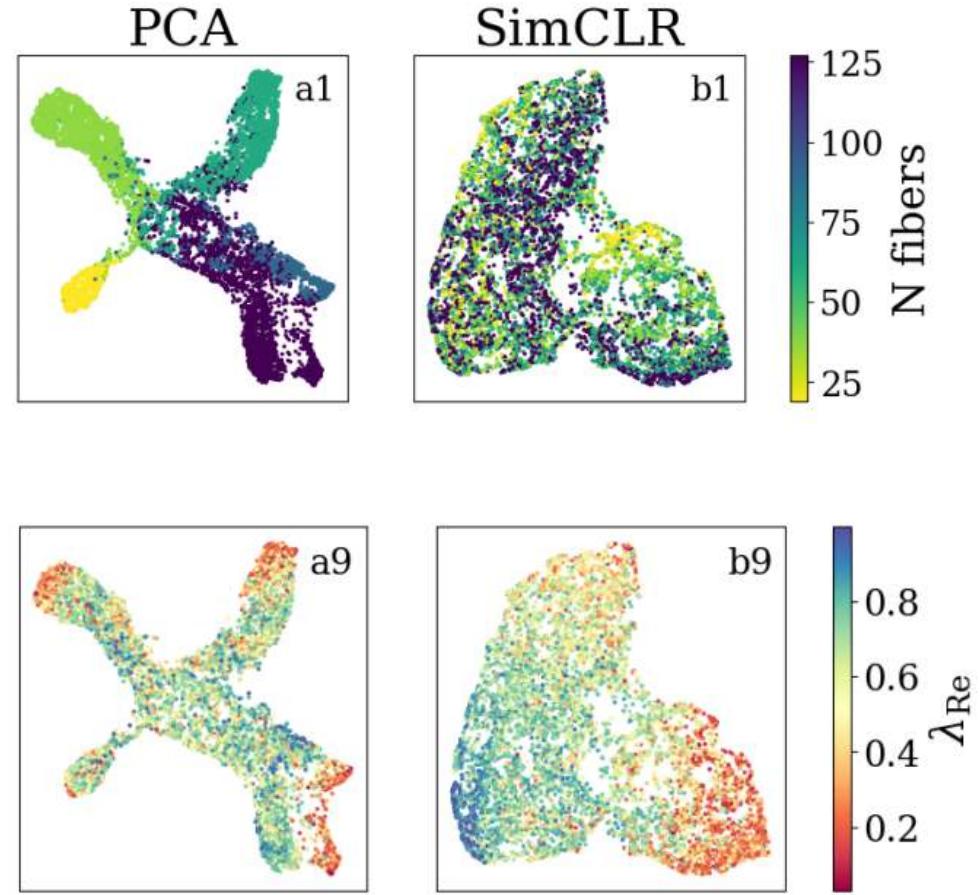
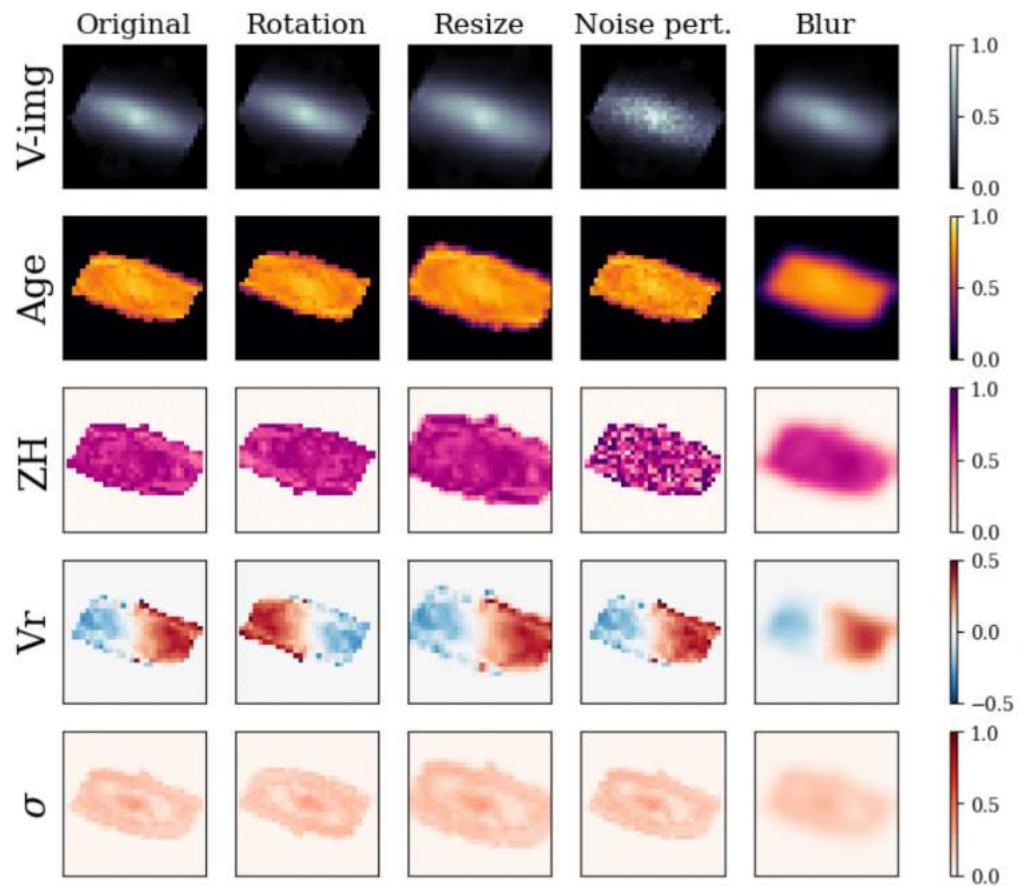
# Hayat+21



# SELF-SUPERVISED LEARNING REACHES COMPARABLE ACCURACY TO FULLY SUPERVISED APPROACHES...



## 2. Sampling



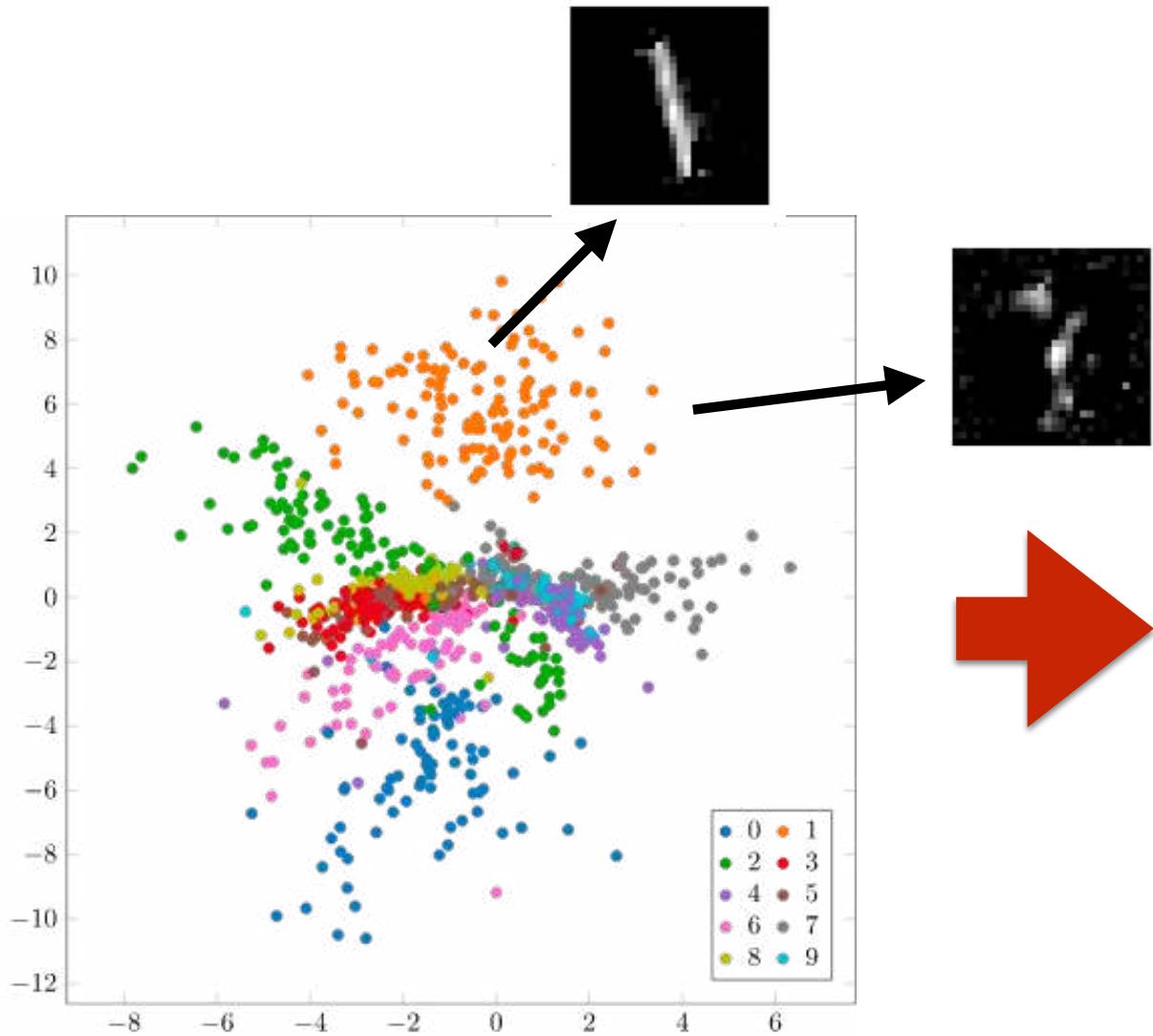
## Contrastive learning representation of Manga galaxies

Sarmiento+21

HOW MIGHT YOU SOLVE THESE  
RELATED PROBLEMS ONCE YOU HAVE  
THE LATENT SPACE COORDINATES FOR  
YOUR TRAINING SAMPLE?

GENERATE A RANDOM SAMPLE DRAWN FROM THE INPUT  
DISTRIBUTION (“**SAMPLING**”)

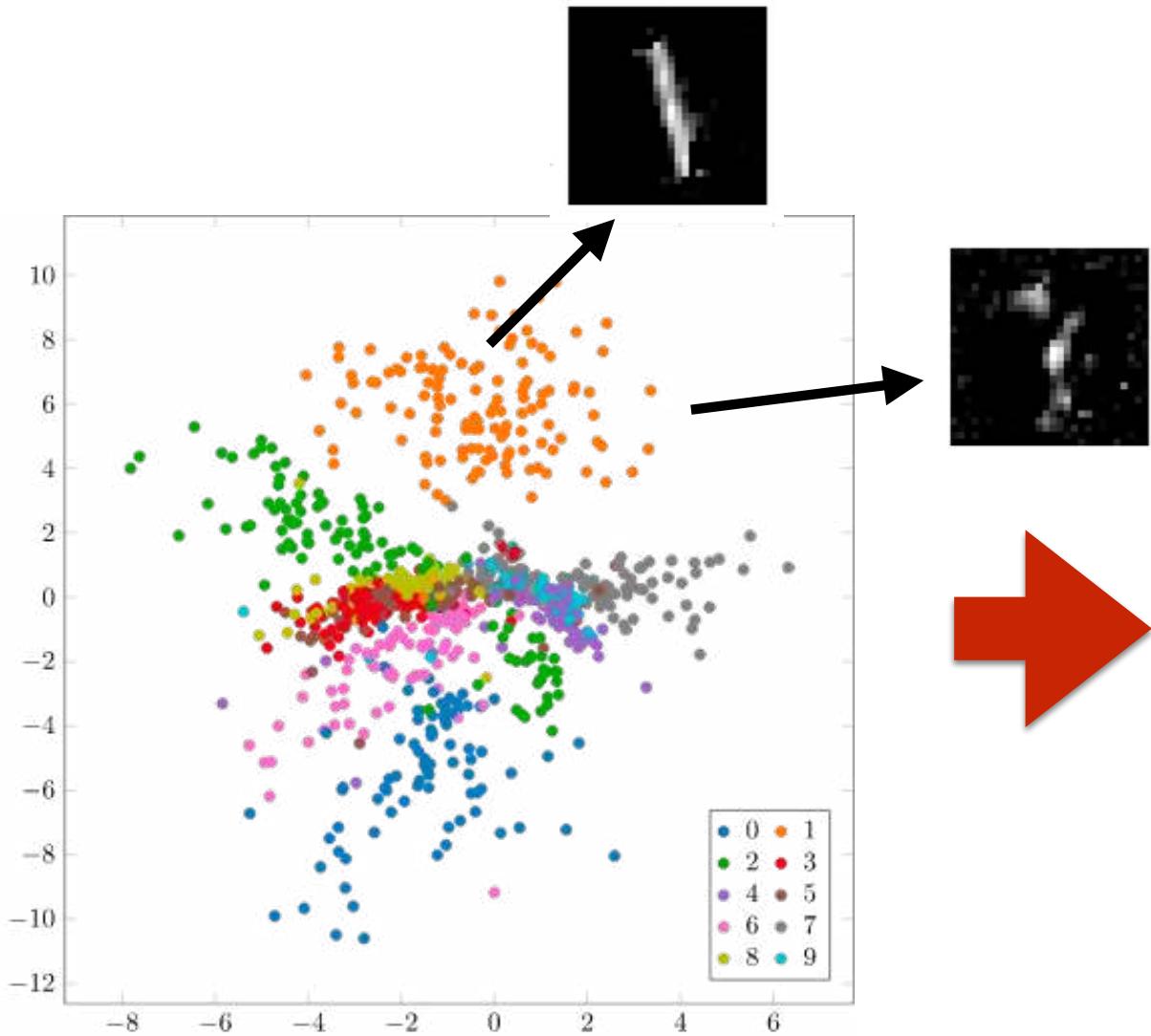
ESTIMATE THE PROBABILITY DENSITY OF AN ARBITRARY  
INPUT, RELATIVE TO THE INPUT DISTRIBUTION (“**DENSITY  
ESTIMATION**”)



z space (latent space)

HOW CAN I  
ESTIMATE  $P(X)$ ?

$(P(z|x))$



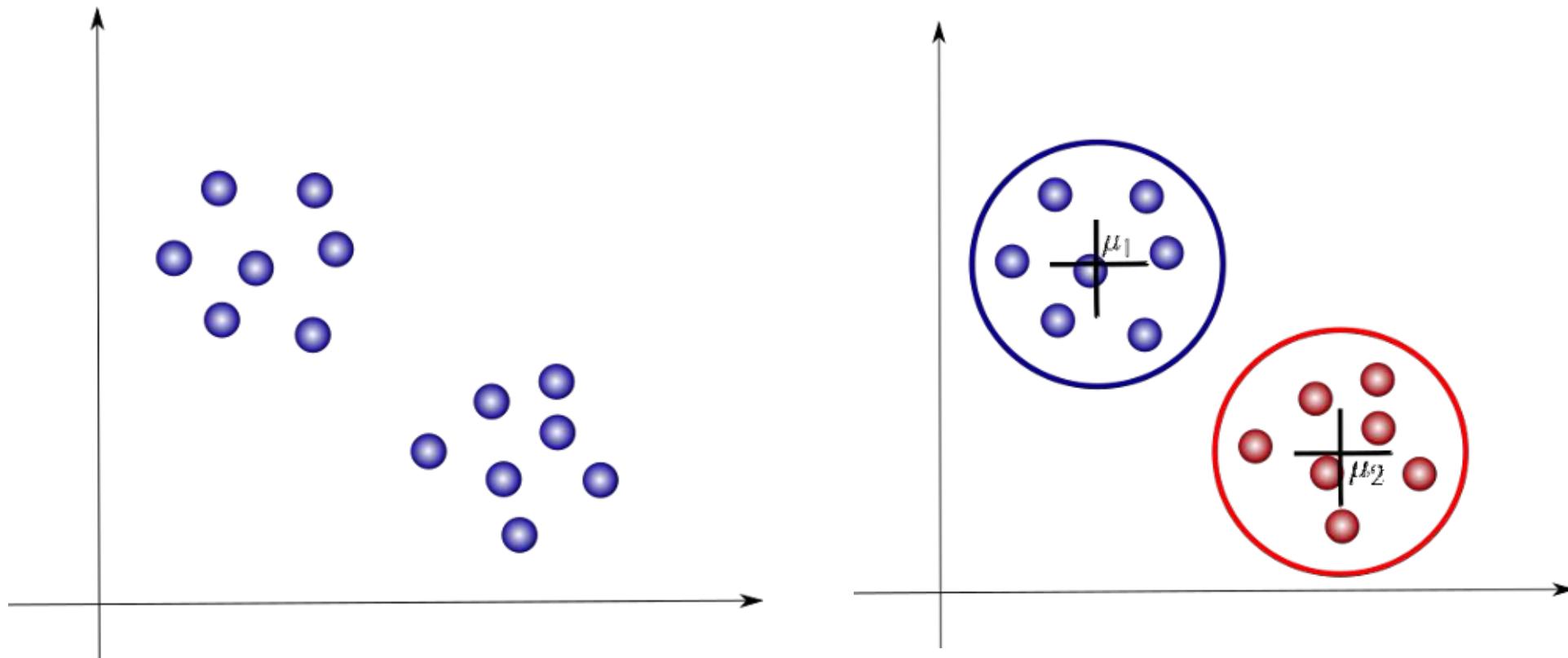
HOW CAN I  
ESTIMATE  $P(X)$ ?

$$(P(z|x))$$

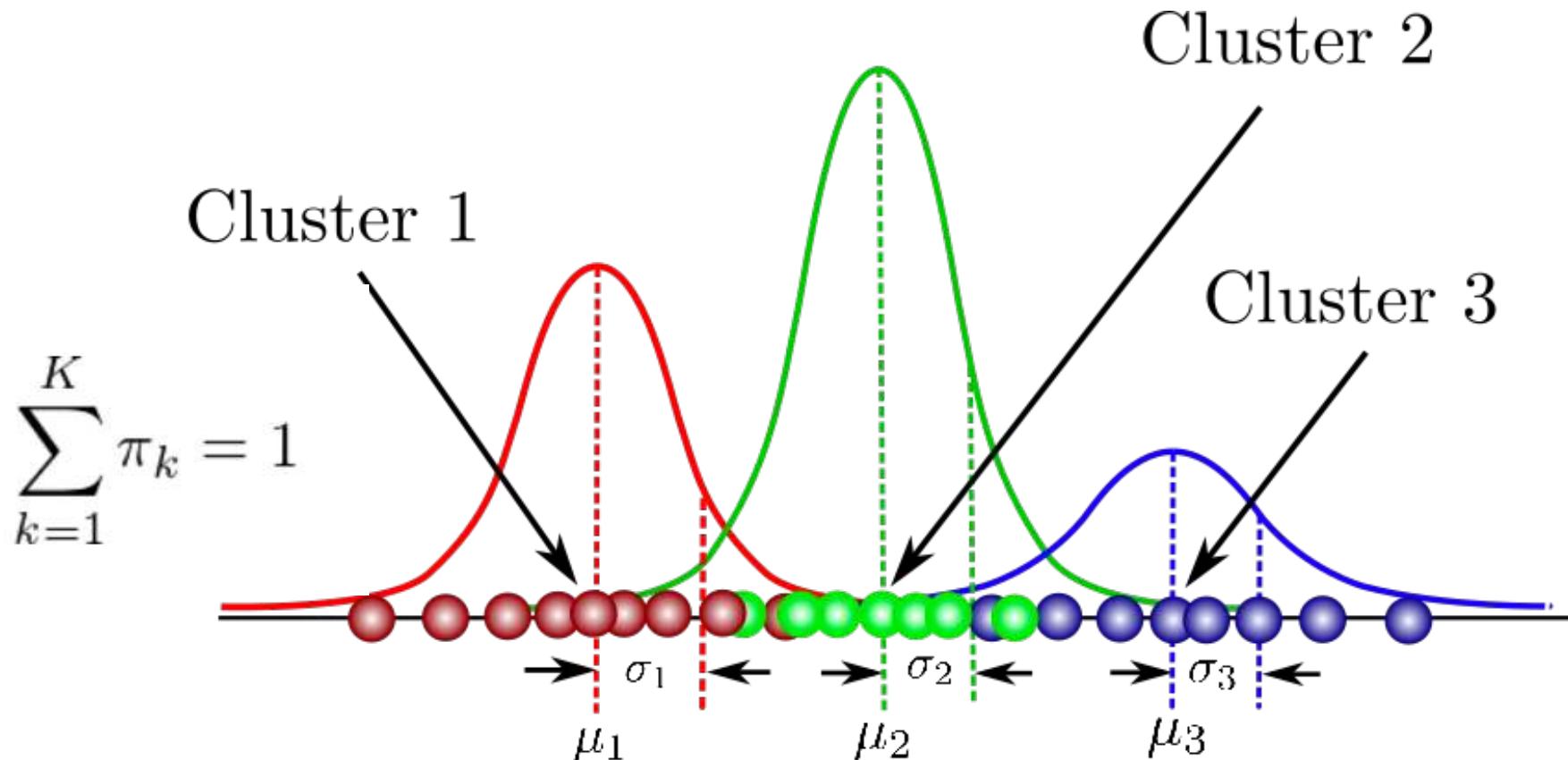
z space (latent space)

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

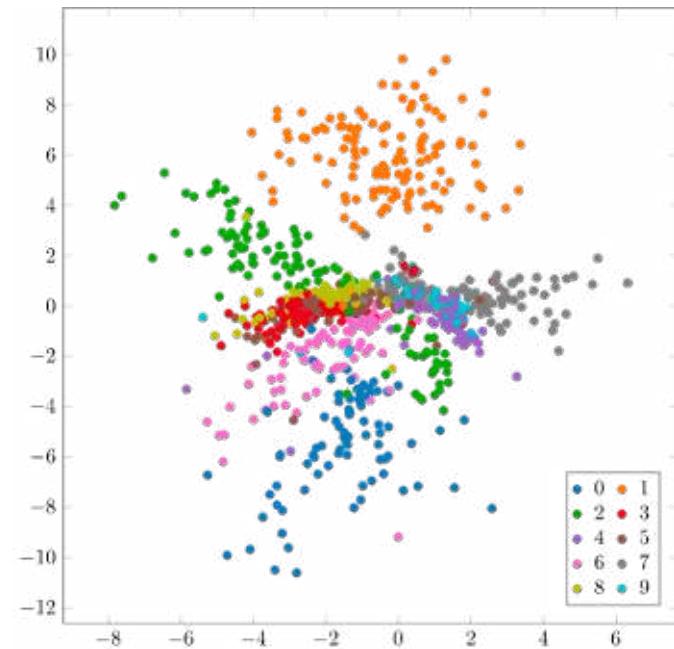
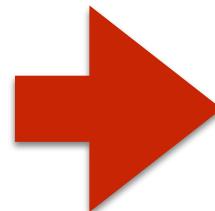
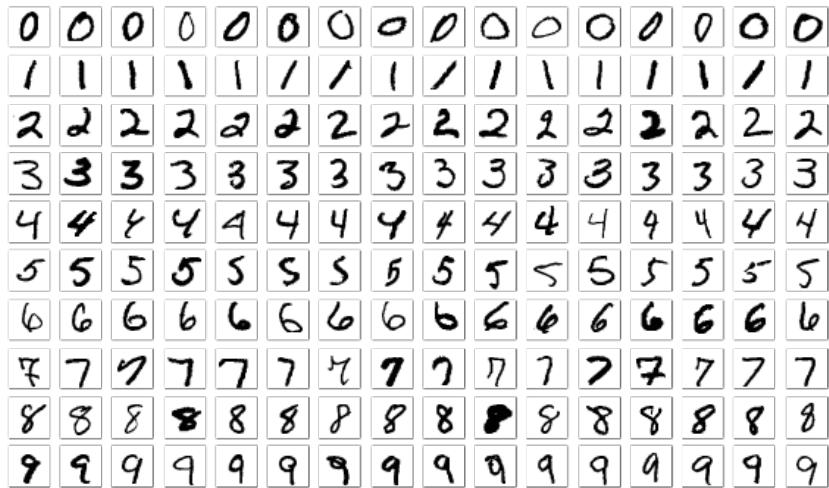


# GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION

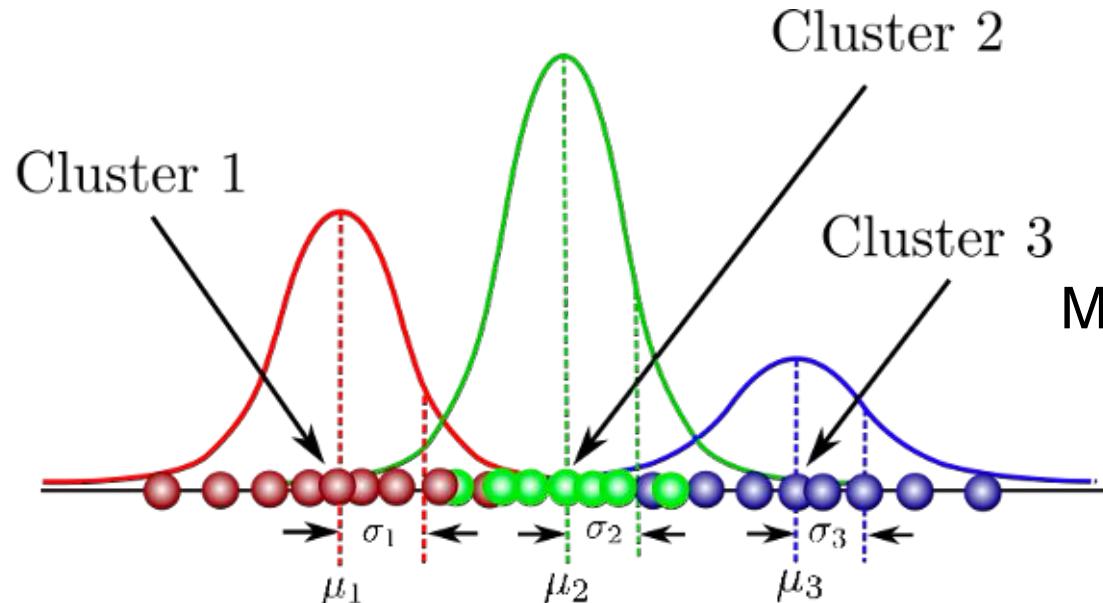


means, sigmas and scale factors of each gaussian are free parameters

# DATA



REPRESENTATION  
(AUTOENCODER -  $z$ )



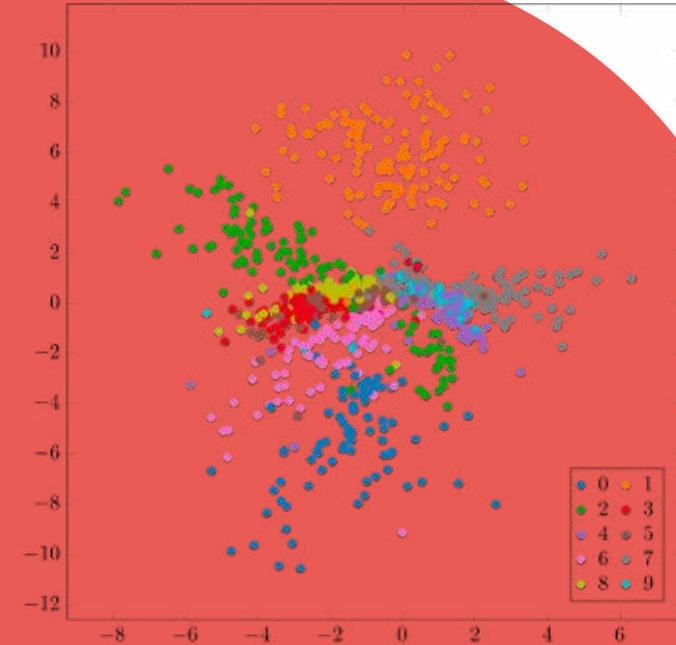
MODELLING OF  $P(z|x)$  WITH GMMs

# 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

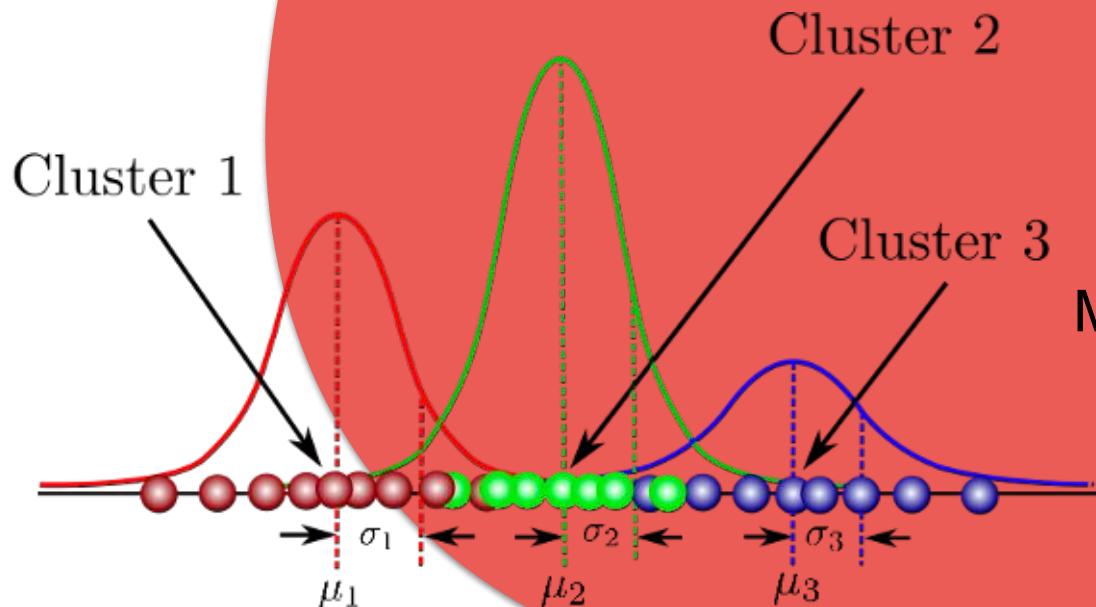
# DATA

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



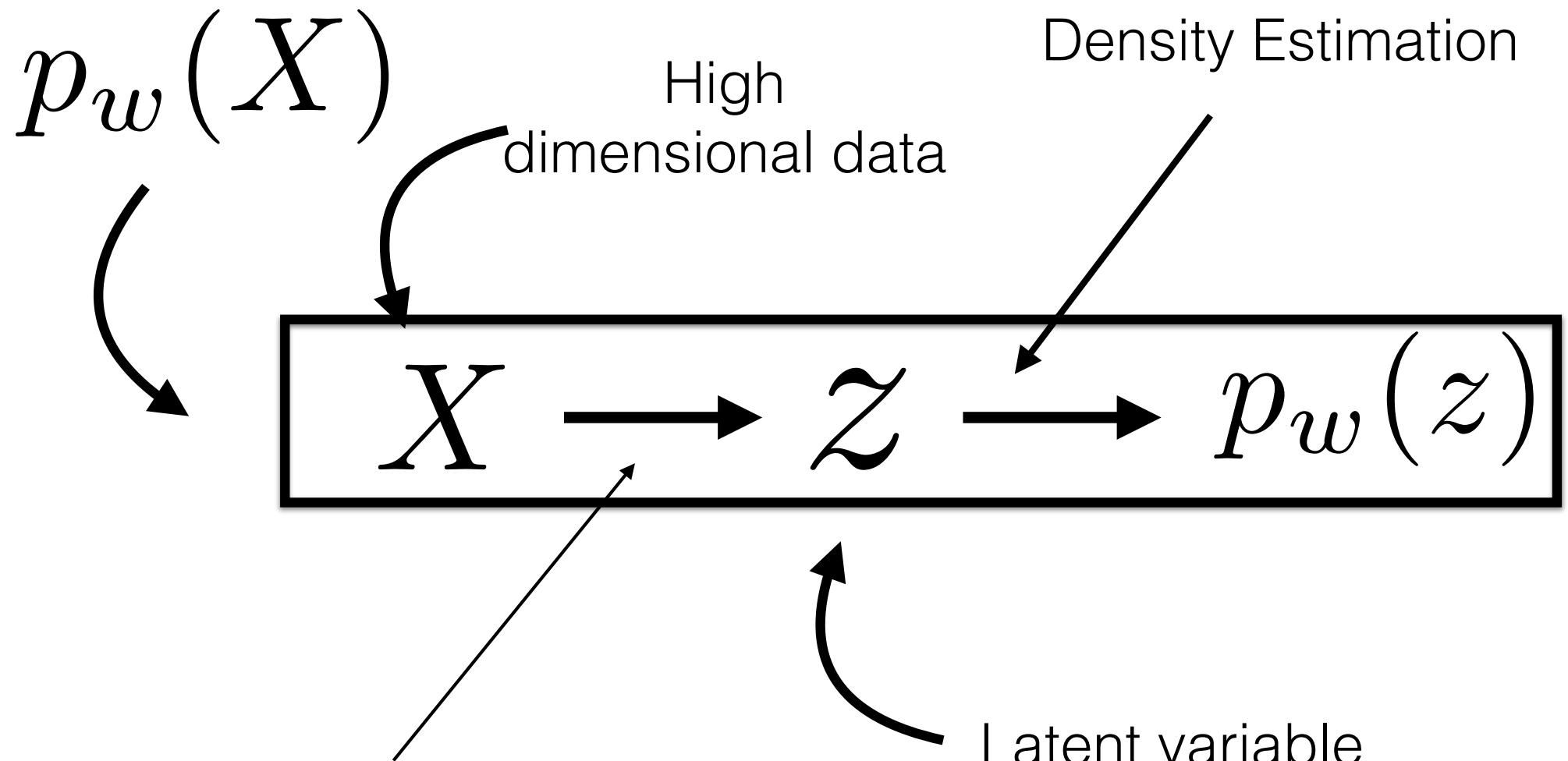
CAN WE COMBINE  
THESE 2 STEPS?

REPRESENTATION  
(AUTOENCODER)



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

## RECAP5b:



Representation learning /  
dimensionality reduction

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 NETWRK)**

**NF-ARF**  
**(NORMALIZING FLOW, AUTOREGRESSIVE FLOW)**

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

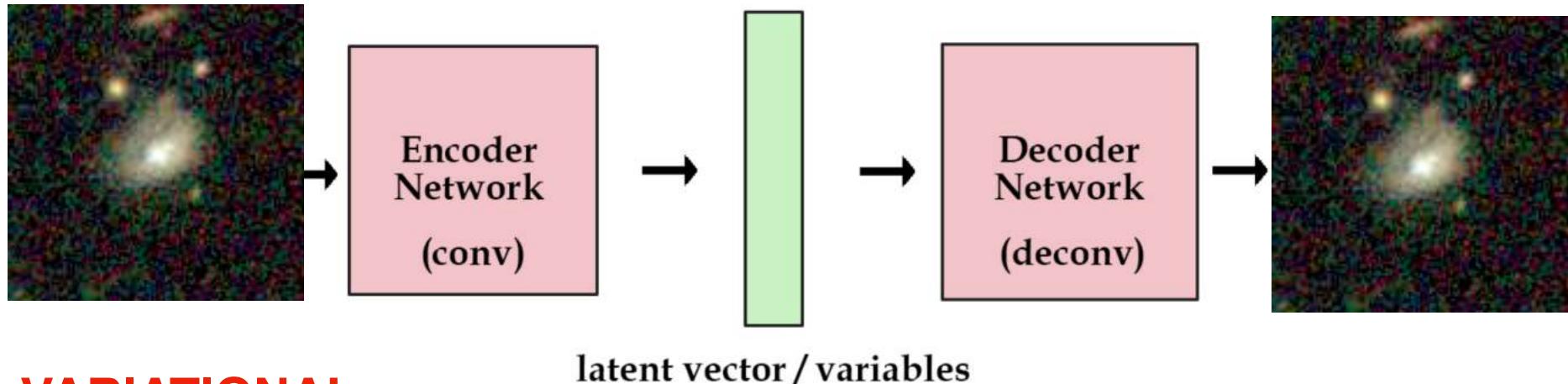
Much of the recent progress in unsupervised deep learning has been to invent network architectures that are capable of solving either or both of these related problems directly, without resorting to any auxiliary methods

**VAE**  
**(VARIATIONAL AUTOENCODER)**

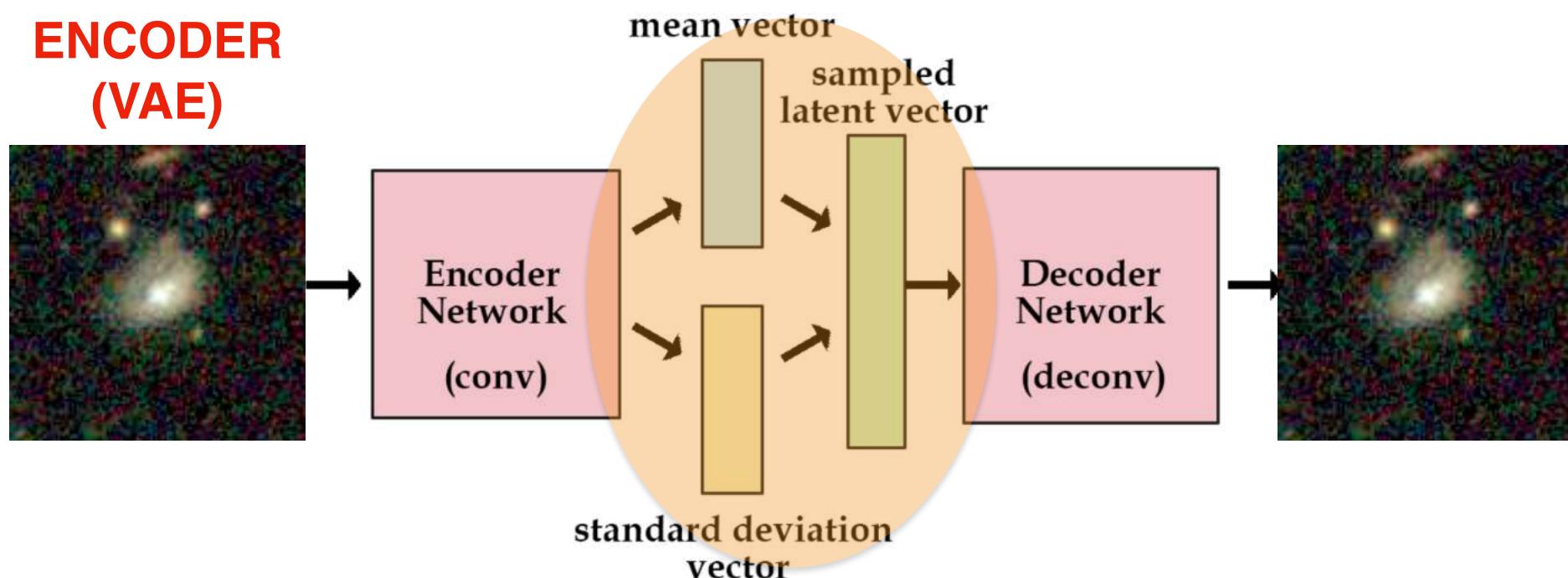
**GAN**  
**(GENERATIVE ADVERSARIAL NETWRK)**

**NF-ARF**  
**(NORMALIZING FLOW, AUTOREGRESSIVE FLOW)**

# AUTO-ENCODER

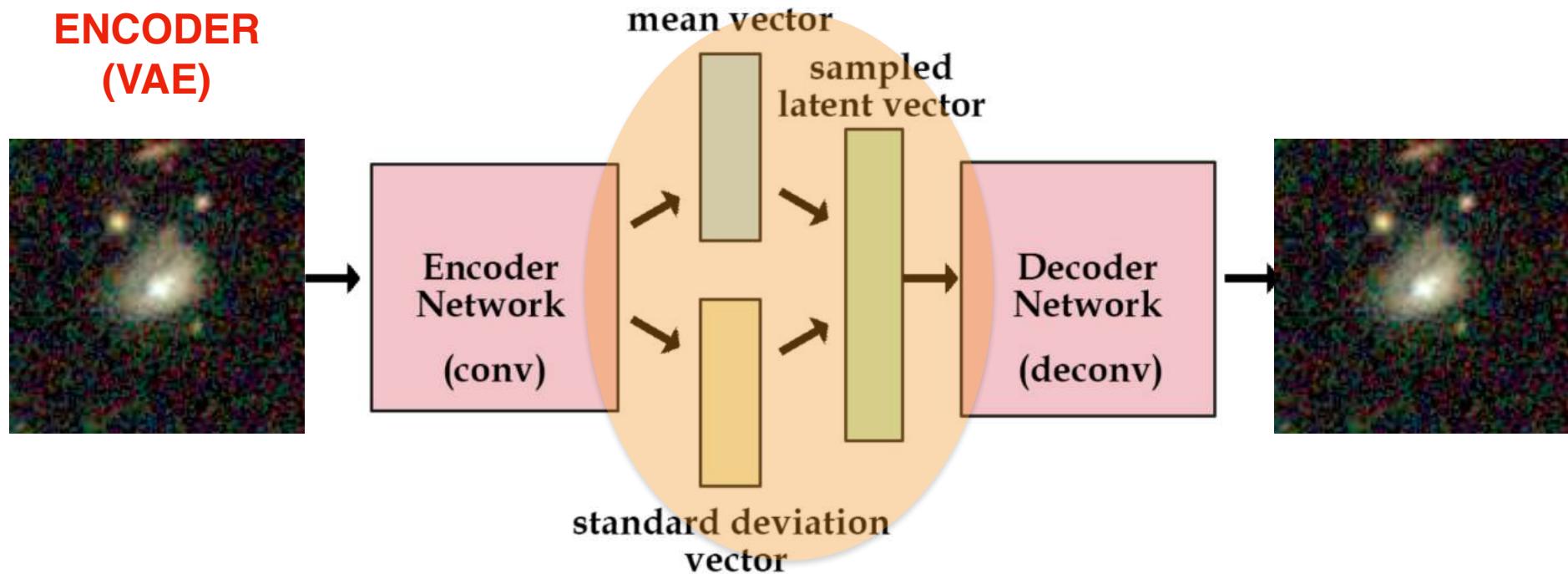


# VARIATIONAL AUTO-ENCODER (VAE)



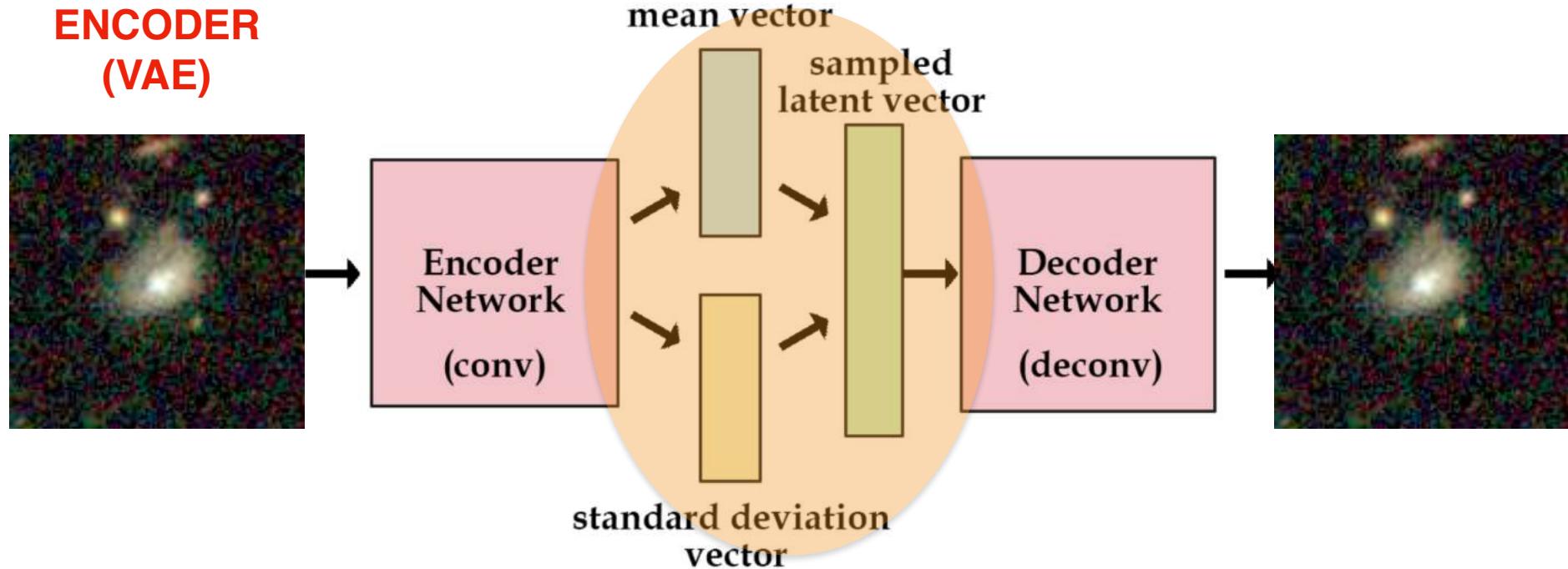
LET'S MODEL THE LATENT SPACE WITH A MIXTURE OF GAUSSIANS (reparametrization trick from the second lecture)

## VARIATIONAL AUTO- ENCODER (VAE)



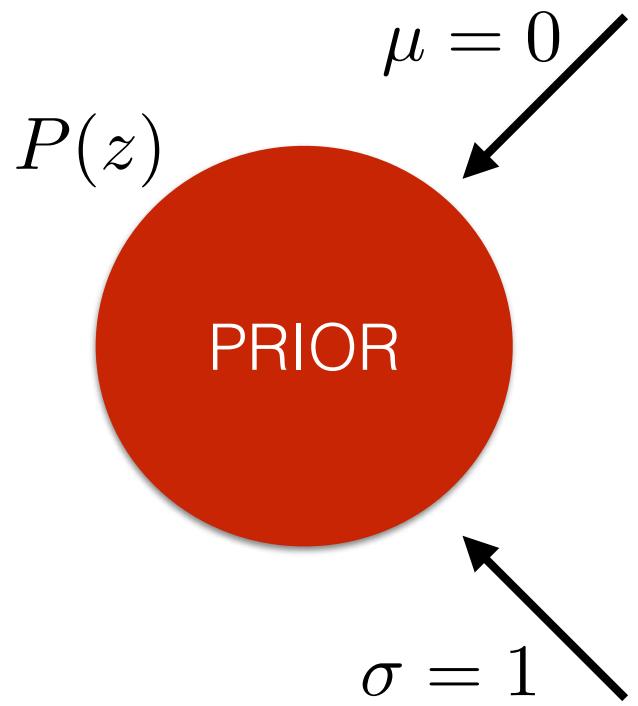
HOWEVER, NOTHING GUARANTEES US THAT  $P(z)$  CAN BE ACCURATELY MODELLED BY A MIXTURE OF GAUSSIANS....

## VARIATIONAL AUTO- ENCODER (VAE)

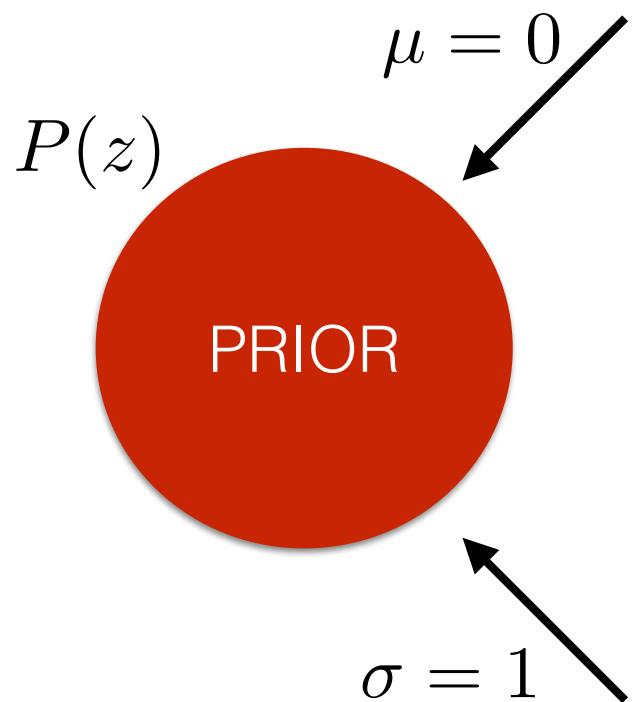


HOWEVER, NOTHING GUARANTEES US THAT THE LATENT  
SPACE CAN BE MODELLED BY A MIXTURE OF  
GAUSSIANS....

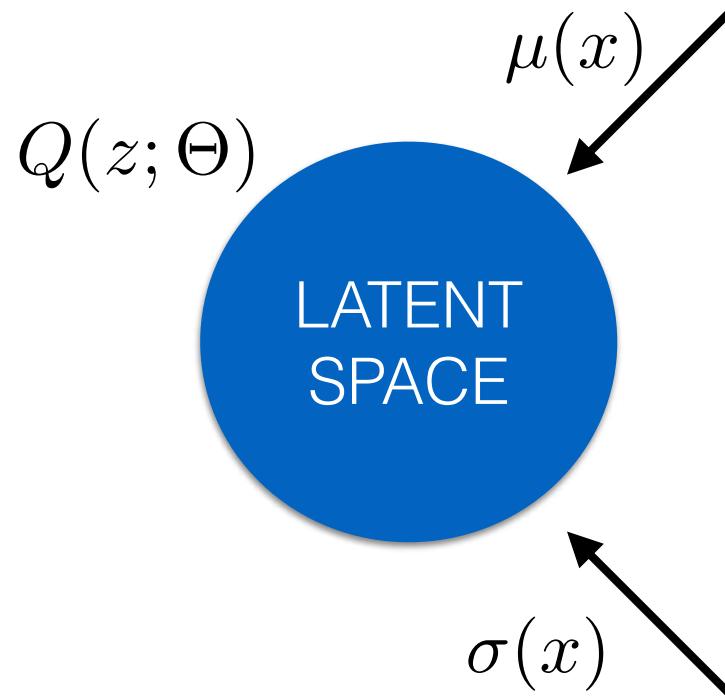
... LET'S FORCE IT TO BE GAUSSIAN LIKE!



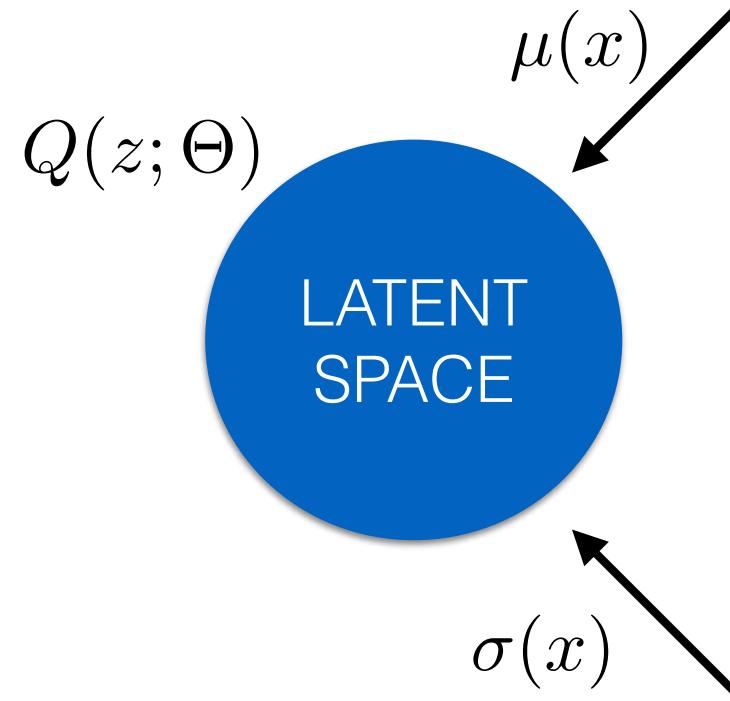
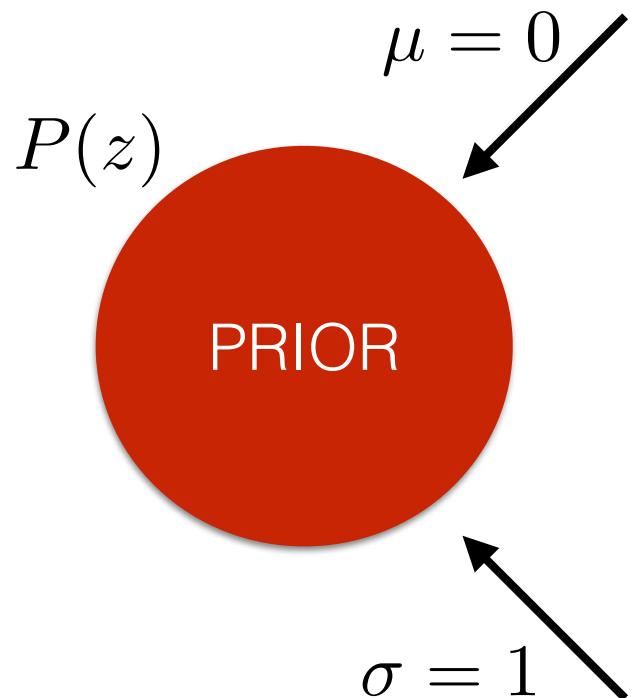
WE ASSUME A SIMPLE PRIOR



WE ASSUME A SIMPLE PRIOR



LATENT SPACE MODELING

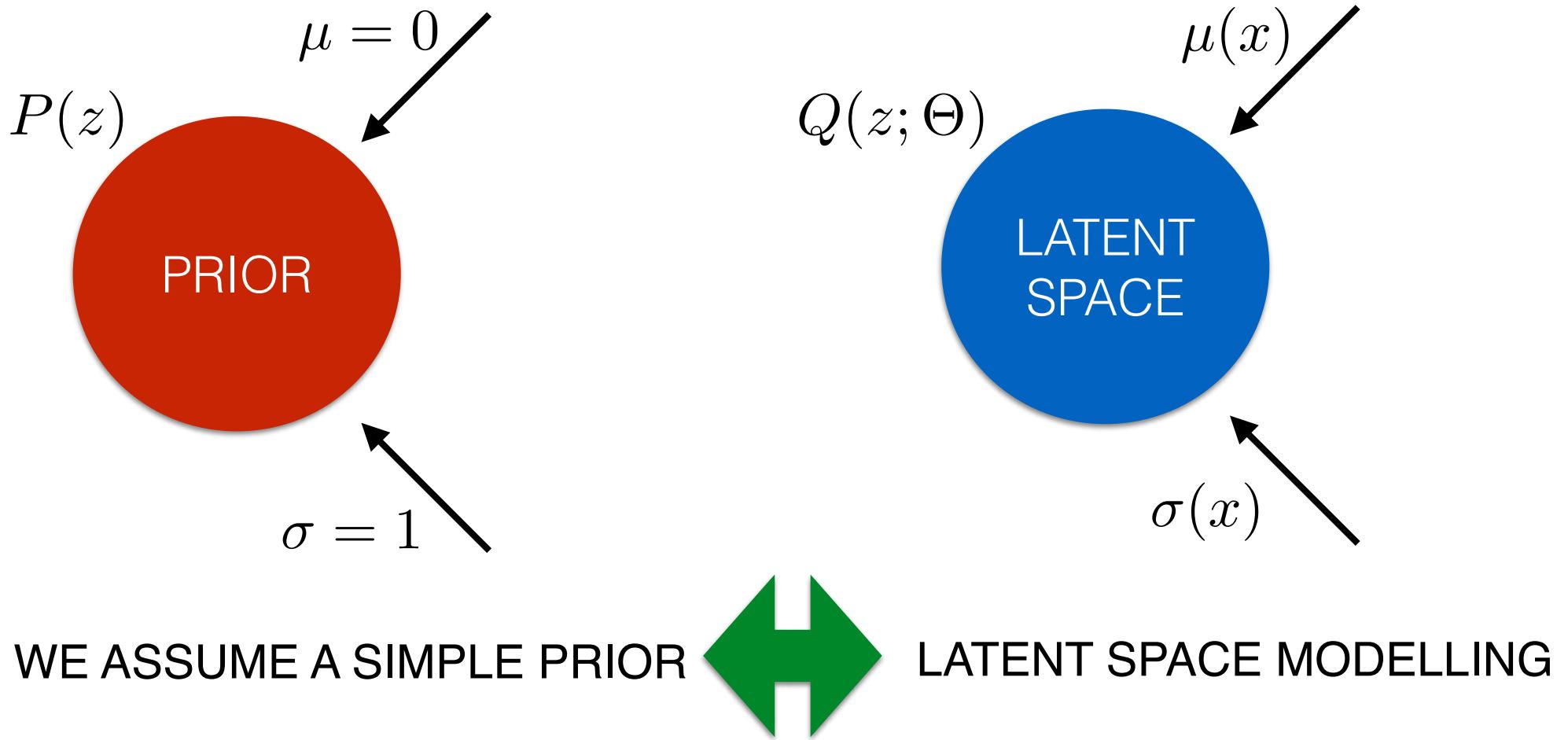


WE ASSUME A SIMPLE PRIOR



LATENT SPACE MODELLING

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



WE 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?

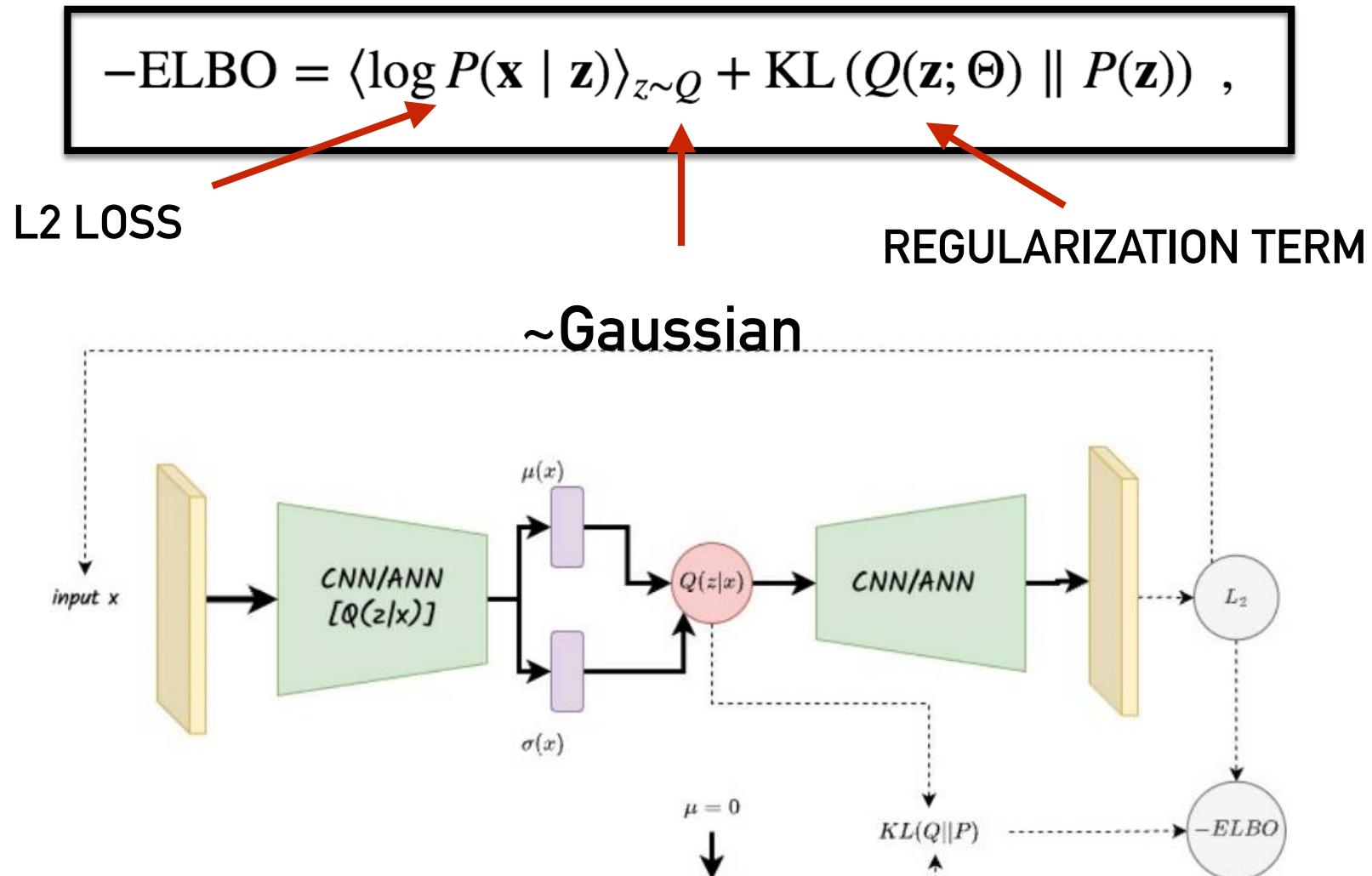
The key insight of VAE is that we are actually performing variational inference here, which then tells us what the loss function should be...

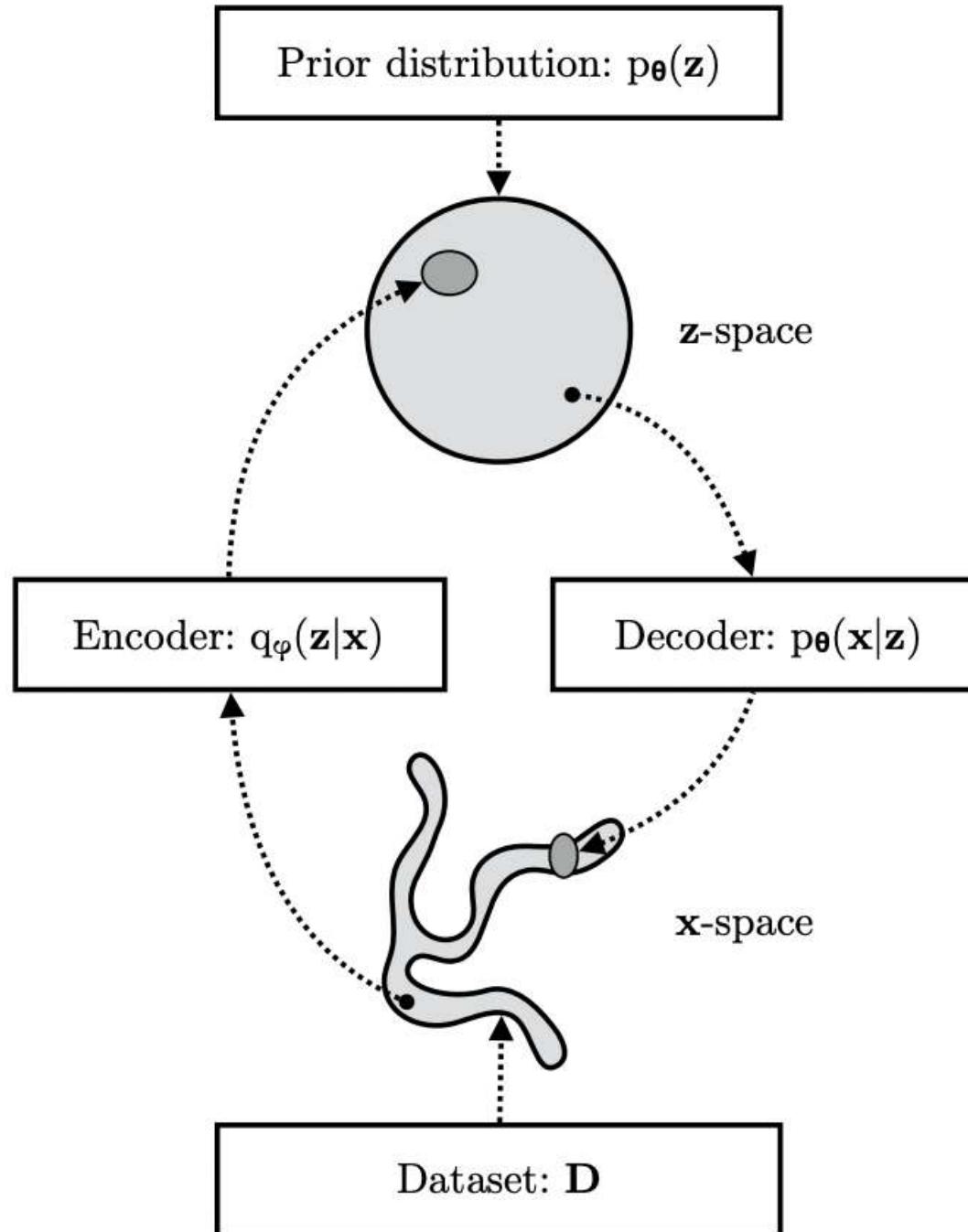
$$-\text{ELBO} = \langle \log P(\mathbf{x} | \mathbf{z}) \rangle_{\mathbf{z} \sim Q} + \text{KL}(Q(\mathbf{z}; \Theta) \| P(\mathbf{z})) ,$$

The equation is framed in a black box. Three red arrows point from labels below the box to specific parts of the equation: one arrow from the label ' $\sim \text{MSE}$ ' points to the term  $\langle \log P(\mathbf{x} | \mathbf{z}) \rangle_{\mathbf{z} \sim Q}$ ; another arrow from the label ' $\sim \text{Gaussian}$ ' points to the term  $\text{KL}(Q(\mathbf{z}; \Theta) \| P(\mathbf{z}))$ ; and a third arrow from the label 'REGULARIZATION TERM' points to the term  $\langle \log P(\mathbf{x} | \mathbf{z}) \rangle_{\mathbf{z} \sim Q}$ .

# WHAT WOULD BE THEN THE LOSS FUNCTION OF A VAE?

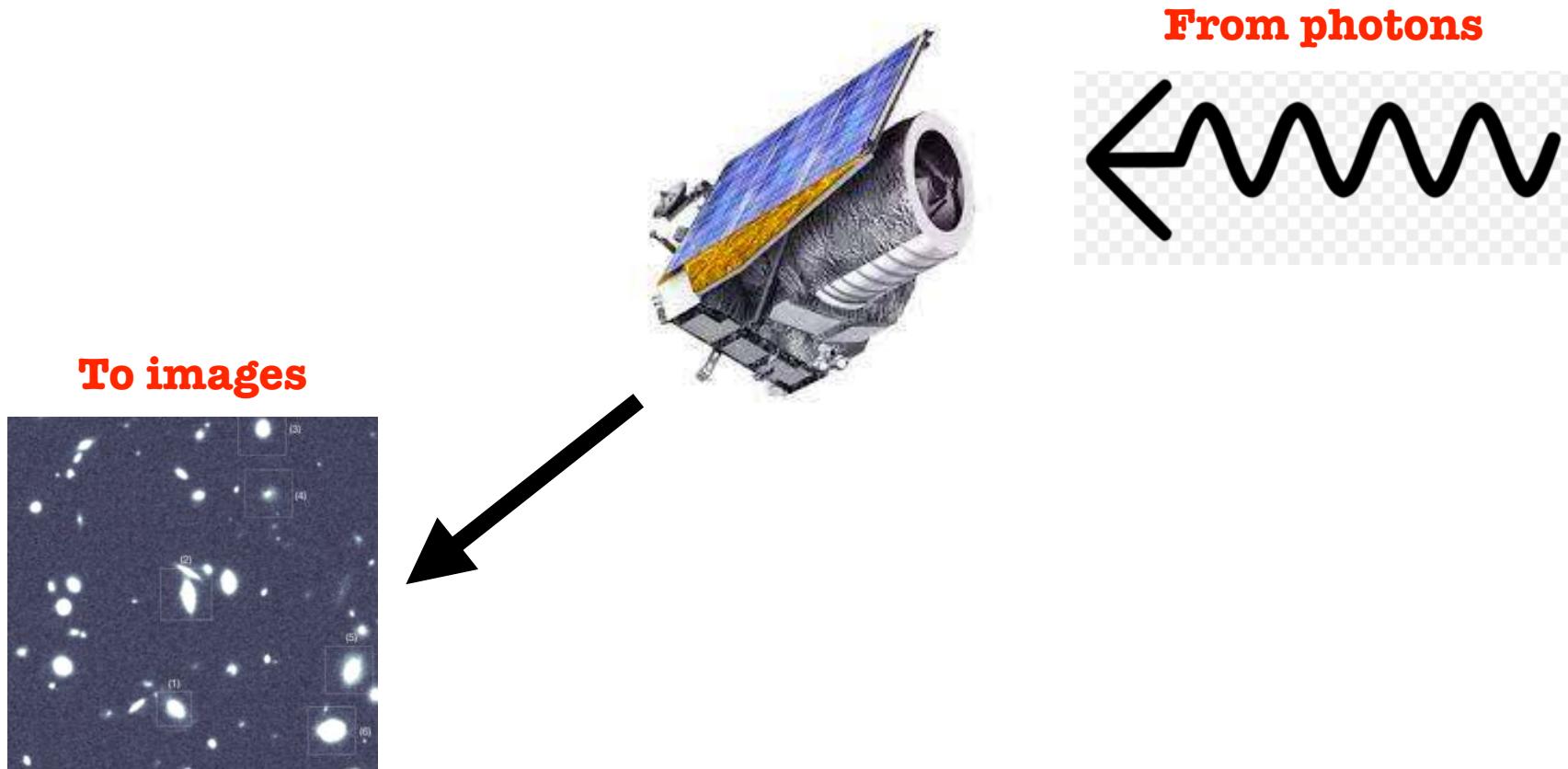
The key insight of VAE is that we are actually performing variational inference here, which then tells us what the loss function should be...





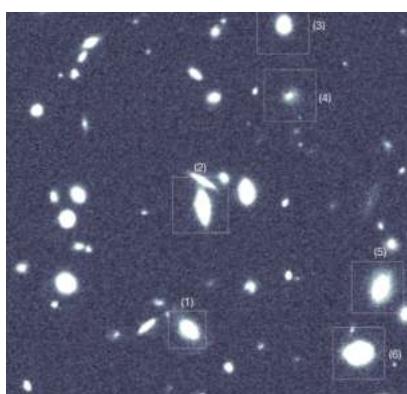
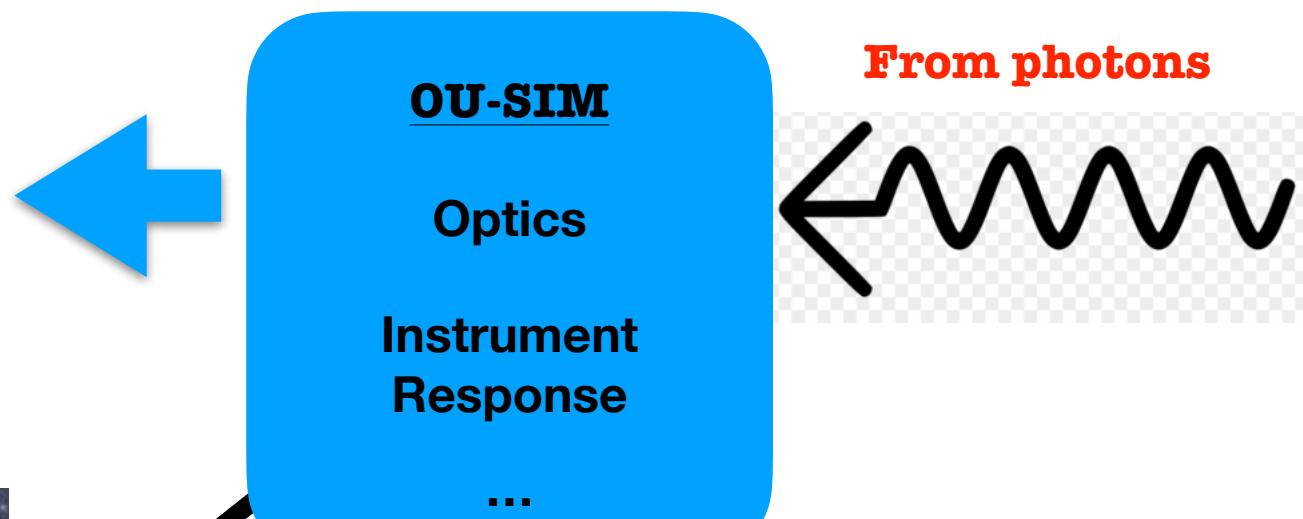
[Kingma and Welling, 2019]

# **Preparing a space mission requires very detailed simulations**



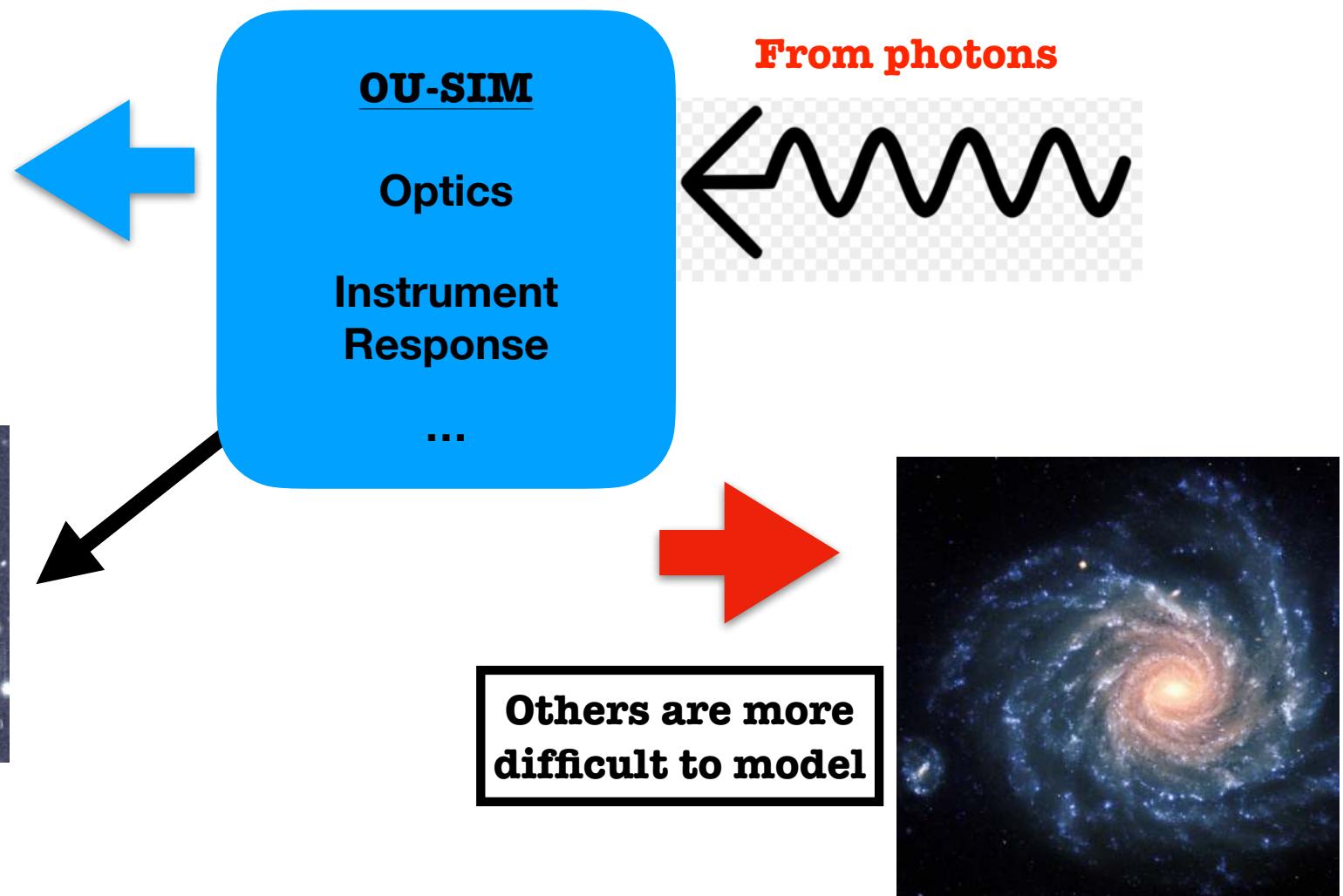
# Preparing a space mission requires very detailed simulations

There is a physical model for many of these effects



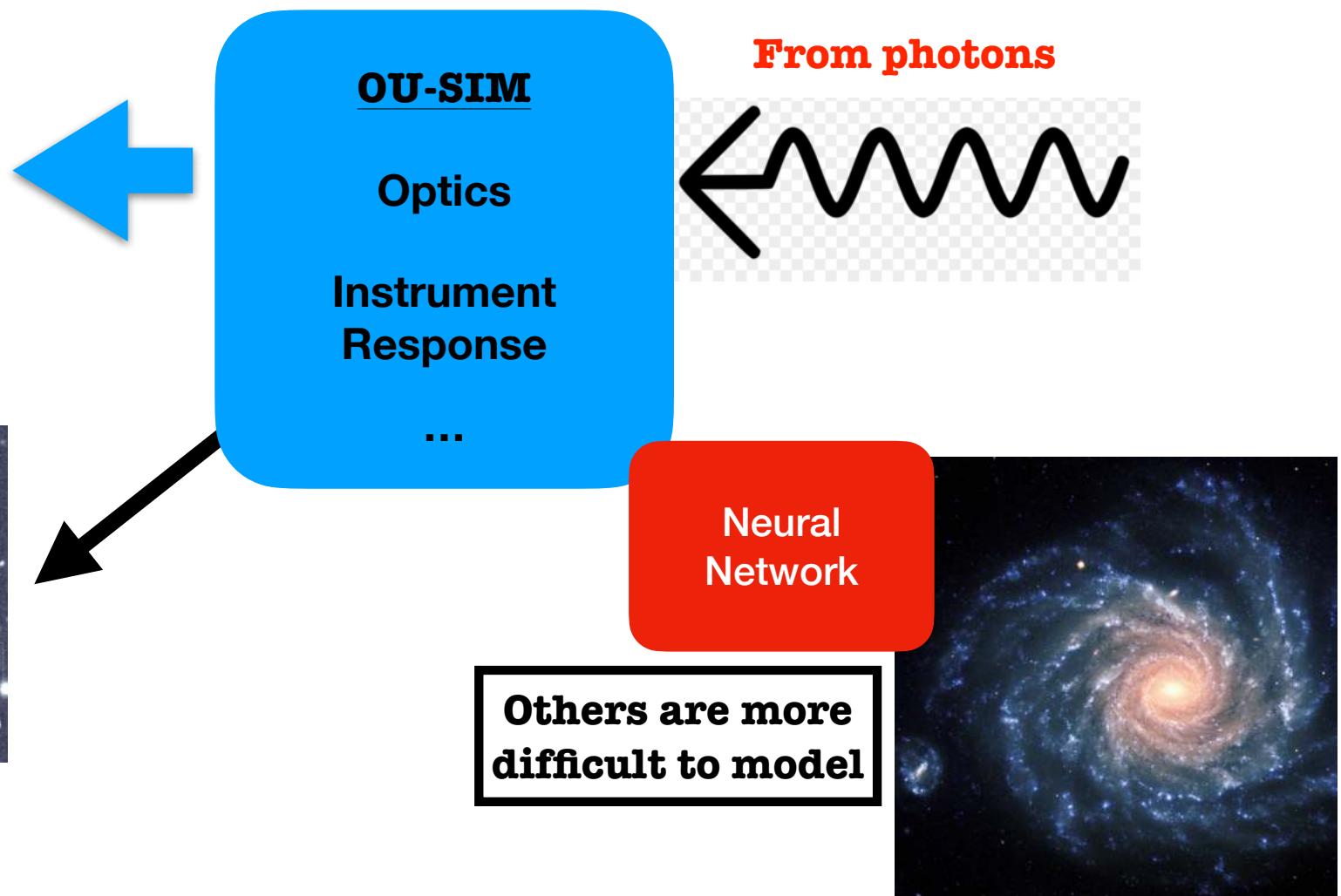
# **Preparing a space mission requires very detailed simulations**

**There is a physical model for many of these effects**

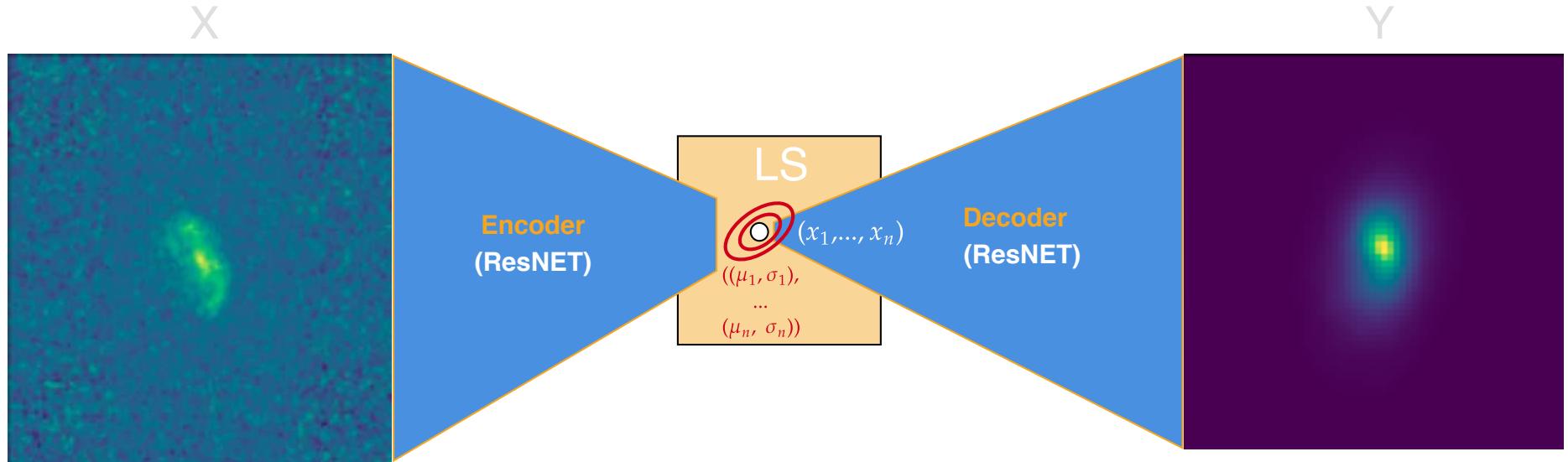


# **Preparing a space mission requires very detailed simulations**

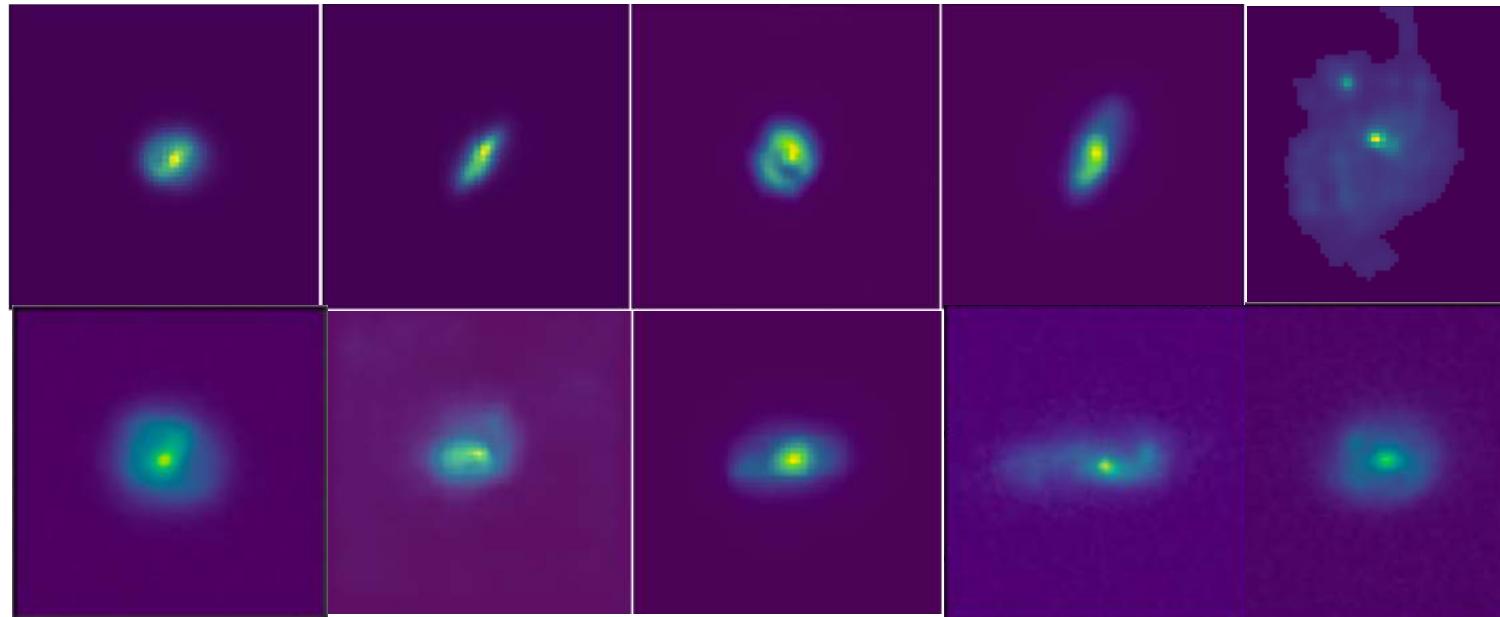
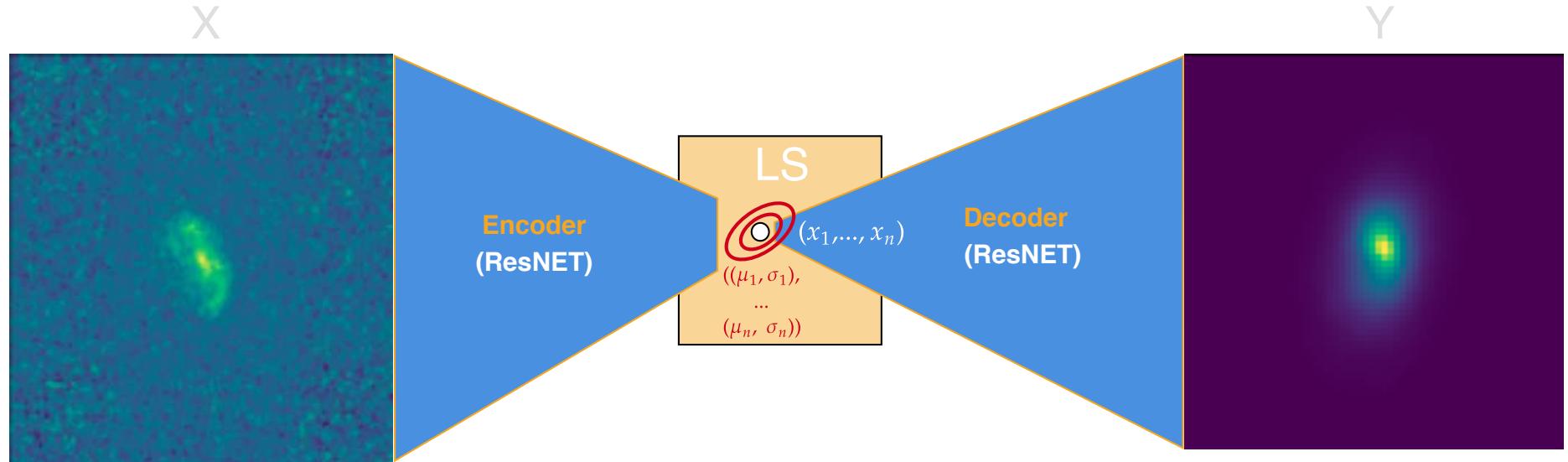
**There is a physical model for many of these effects**

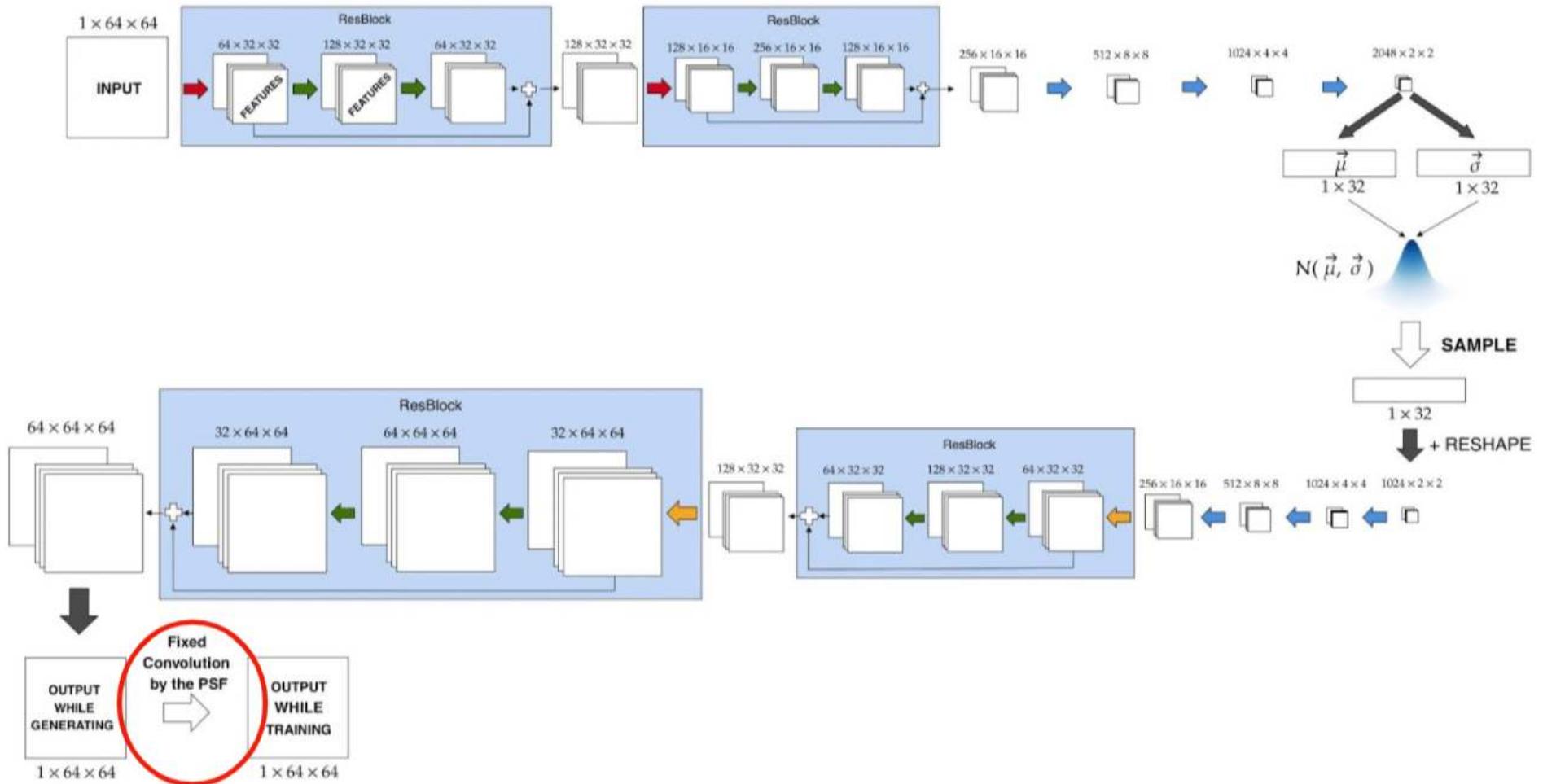
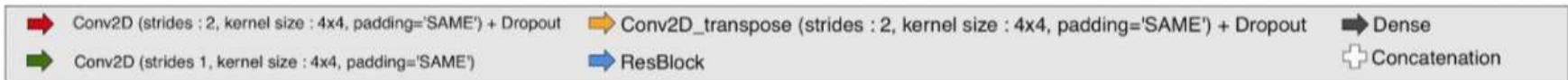


# Include a Generative Model for galaxy generation in the Euclid Simulation Pipeline to model process for which we do not have a physical model



# Include a Generative Model for galaxy generation in the Euclid Simulation Pipeline to model process for which we do not have a physical model





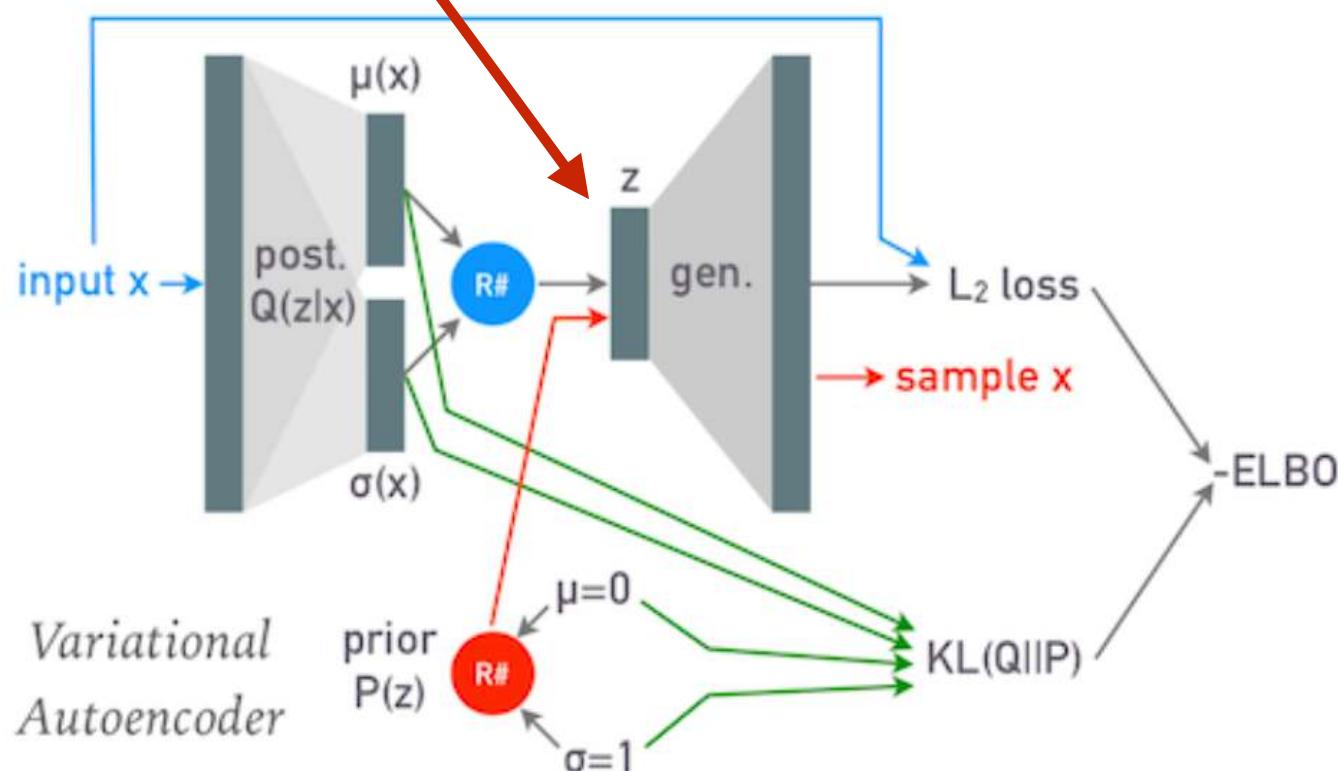
**Physical information is included as well within the architecture**

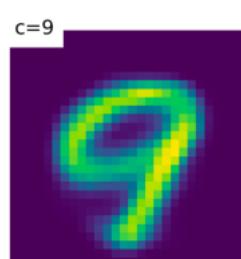
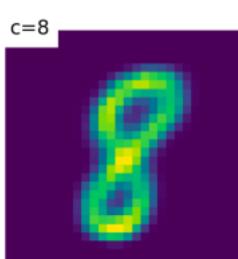
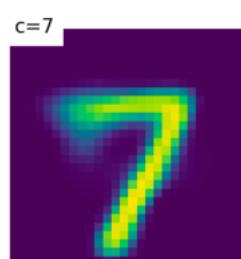
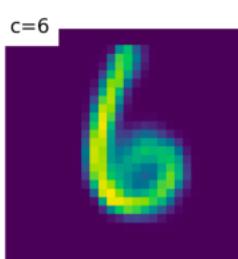
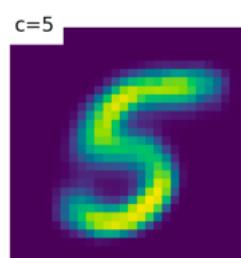
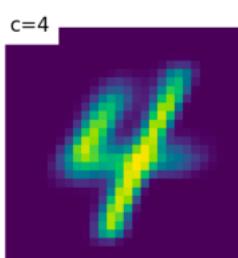
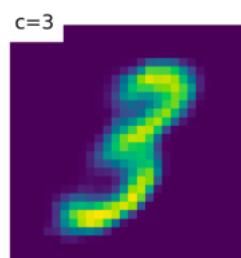
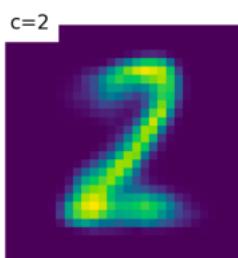
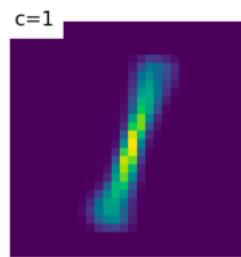
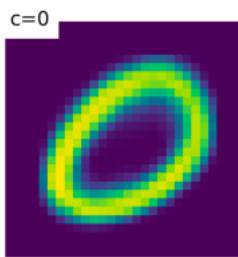
VAE's can also be conditioned, to generate a given class

concatenate labels



$$P(X|Y)$$



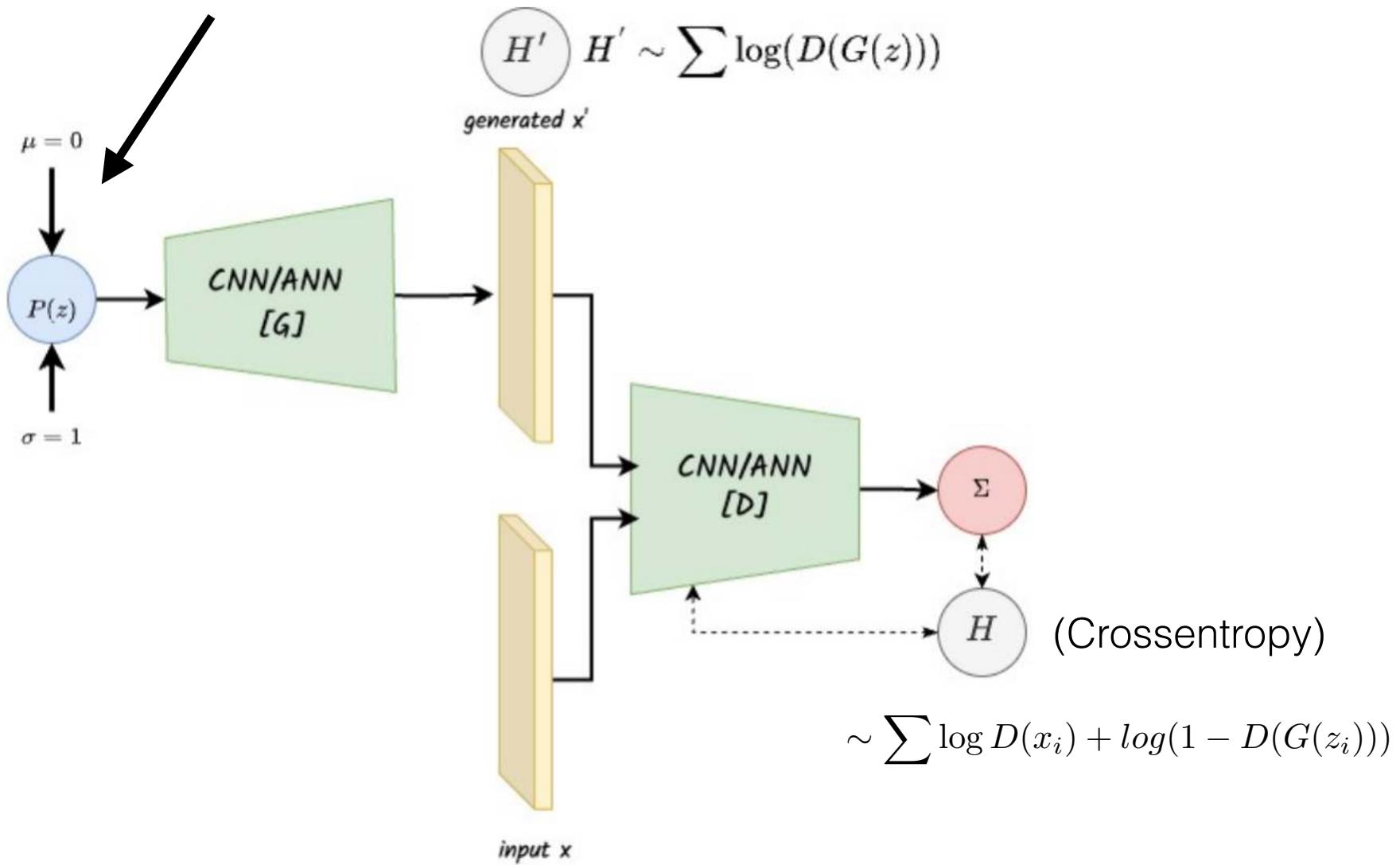


# Generative Adversarial Networks

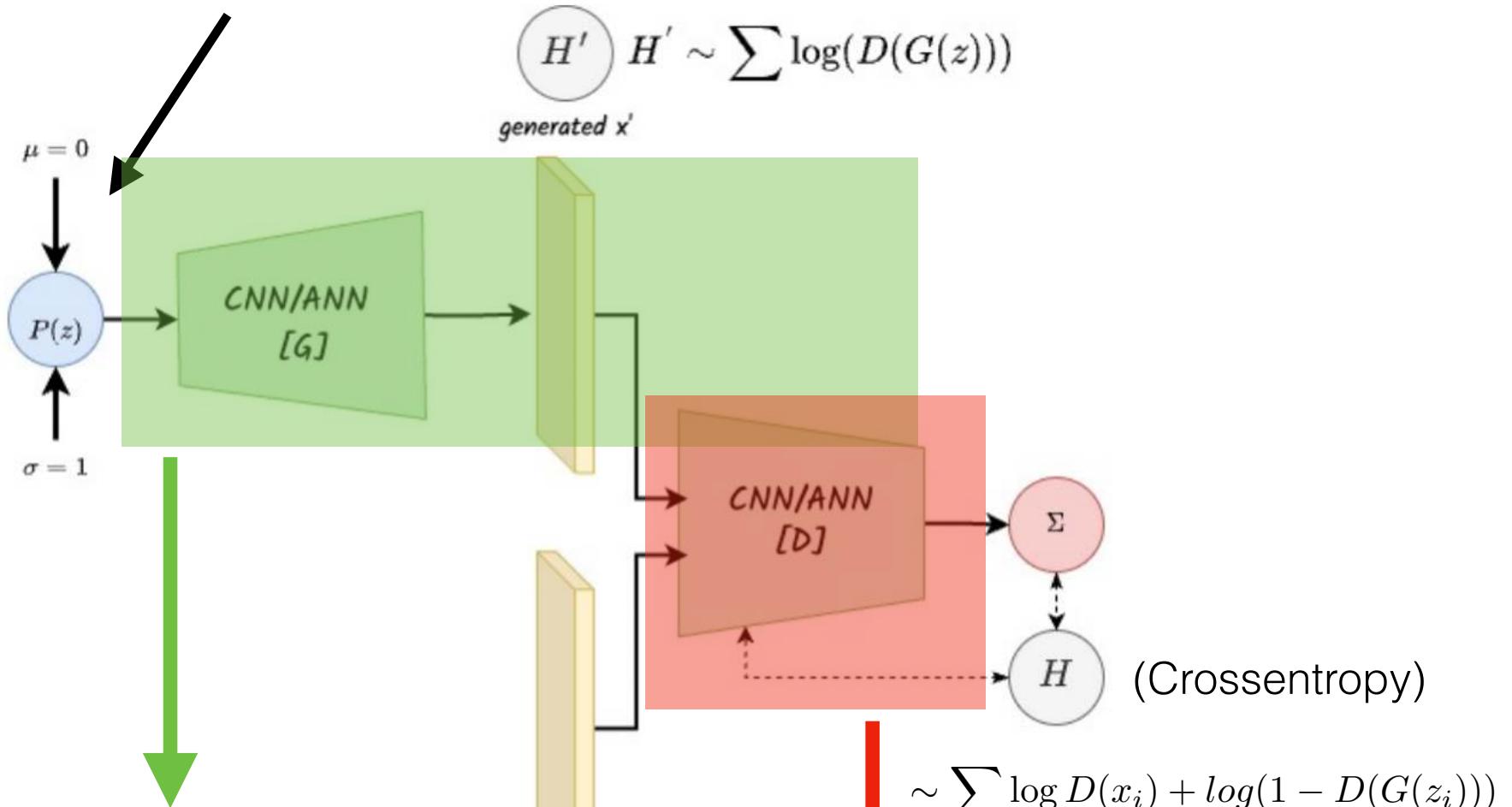
As for VAEs,  
the goal of generative Adversarial Networks (GANs) is to estimate  
 $p(z|x)$

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

# The latent variable is here



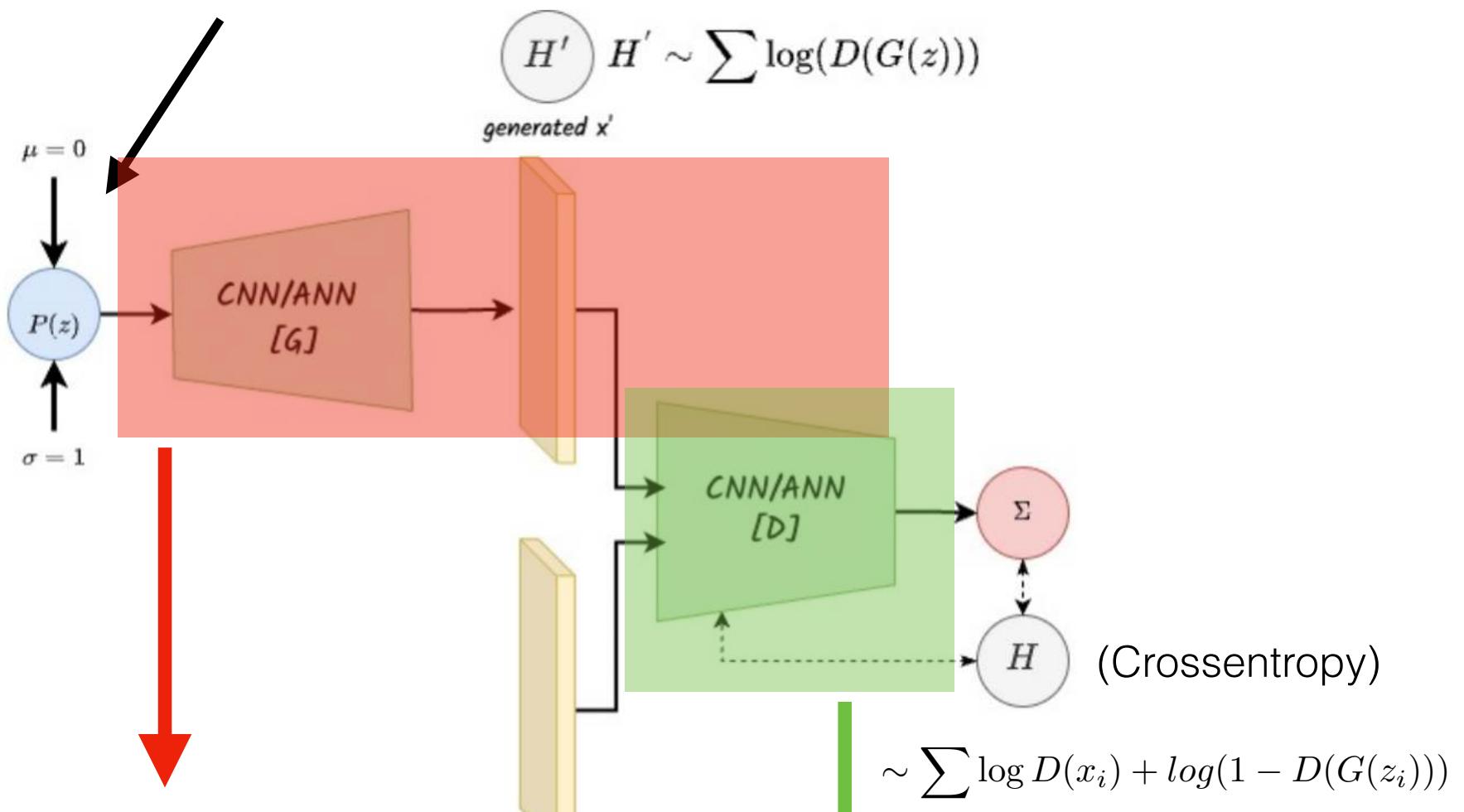
# The latent variable is here



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

Frozen

# The latent variable is here

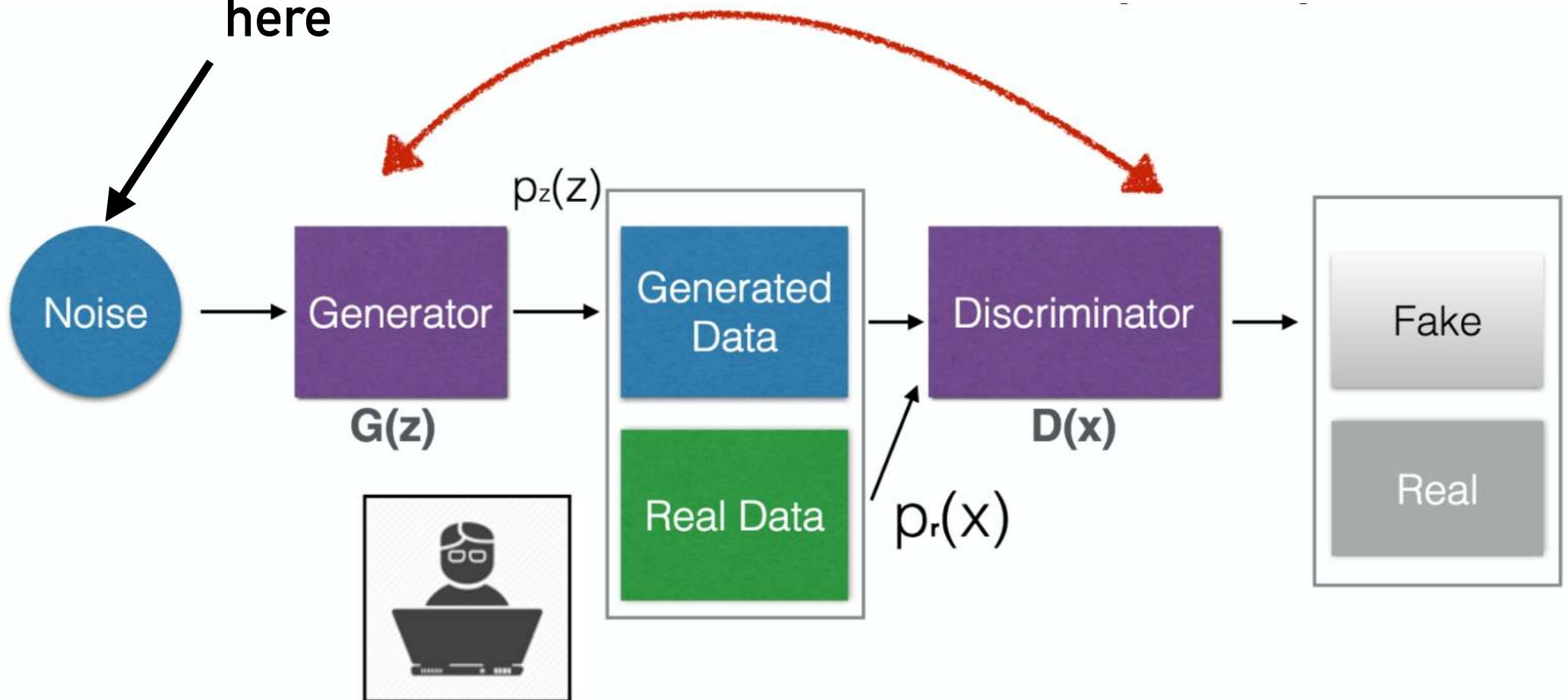


Every N iterations the discriminator  
is trained to force to distinguish between  
real and fake

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+14)

The latent variable is  
here

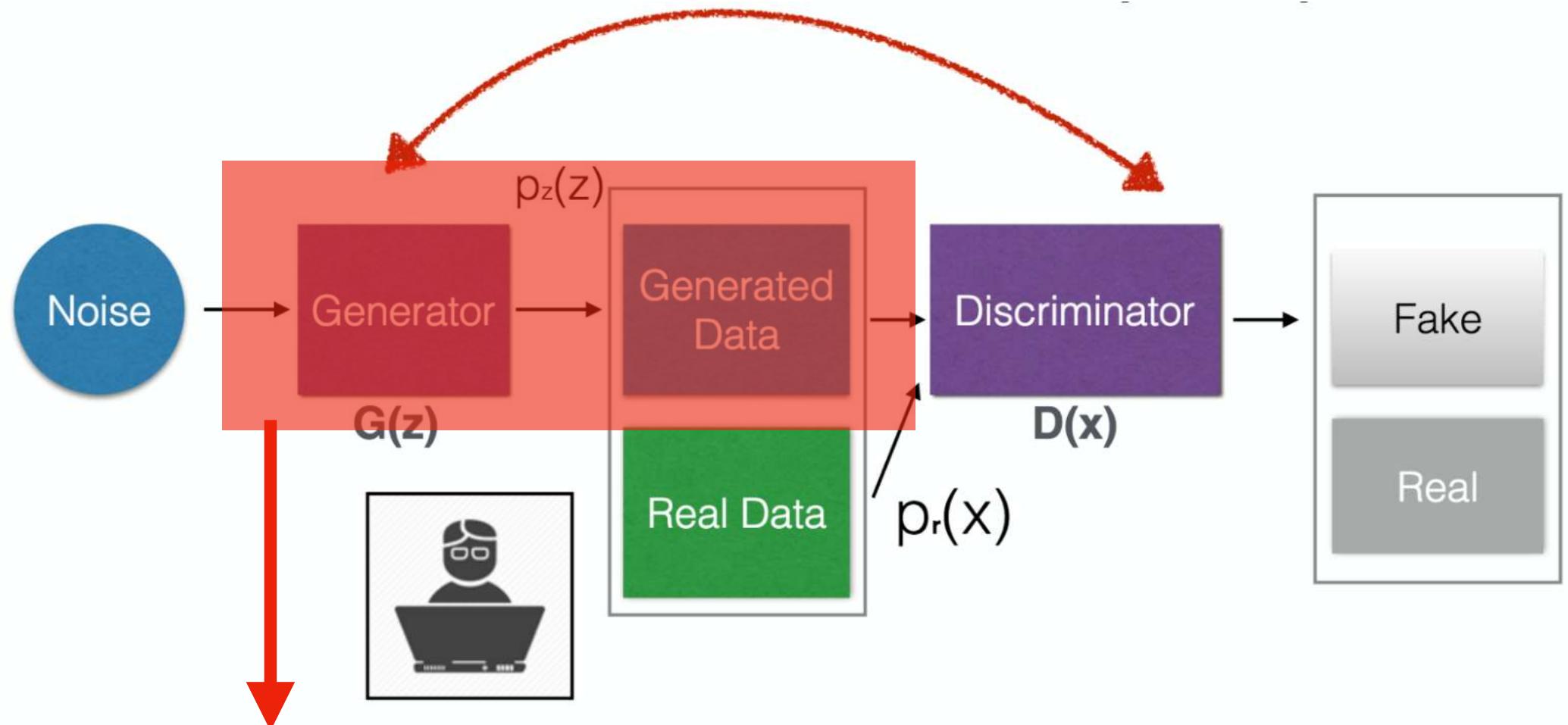


TWO COMPETING NETWORKS

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+)

## TWO COMPETING NETWORKS

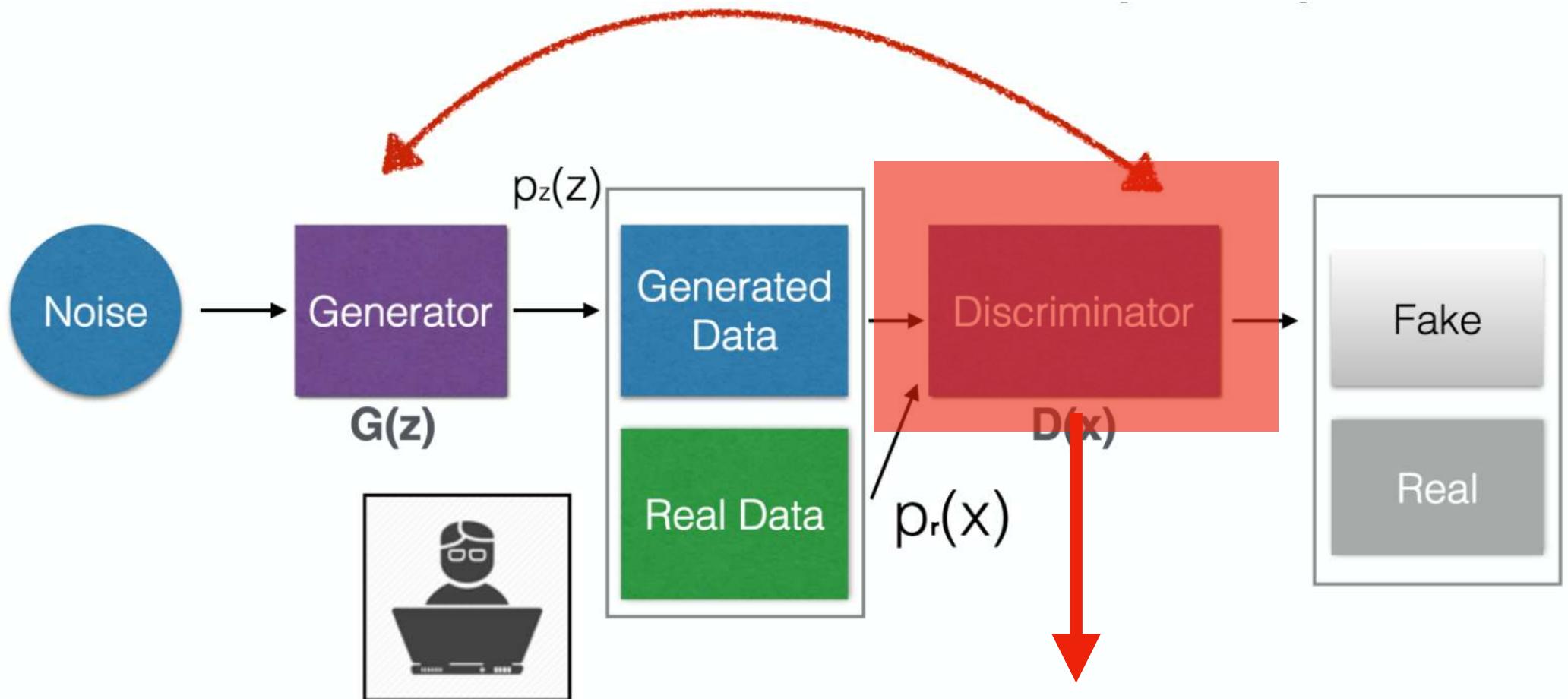


Every N iterations the generator  
is trained to force the discriminator  
to classify as real

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+)

## TWO COMPETING NETWORKS

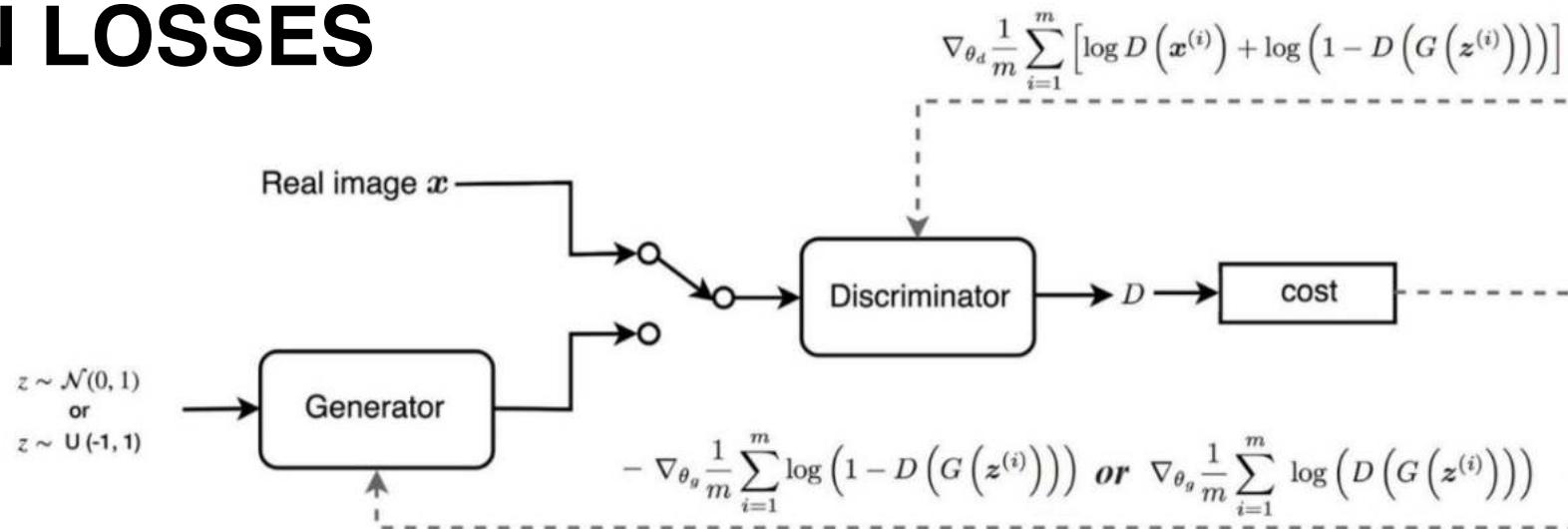


Every N iterations the discriminator  
is trained to force to distinguish between  
real and fake

## IN PRACTICE

### DISCRIMINATOR LOSS (CROSS-ENTROPY)

## GAN LOSSES



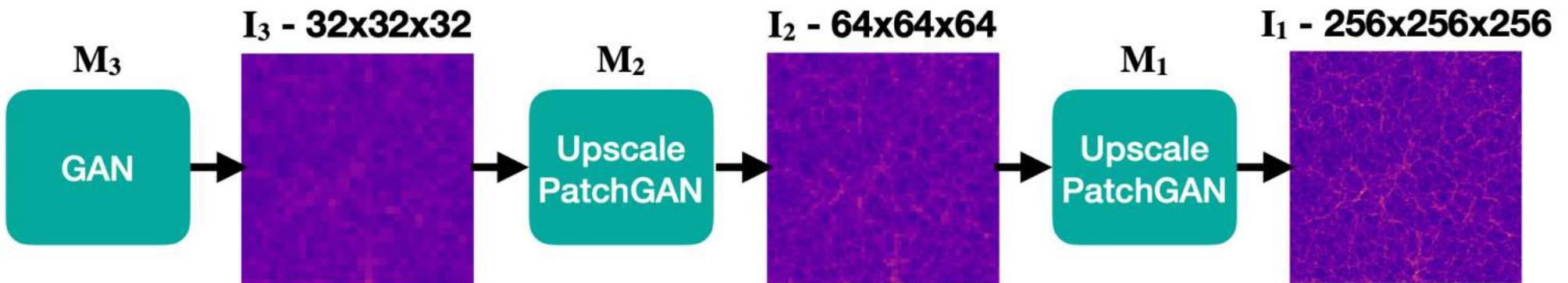
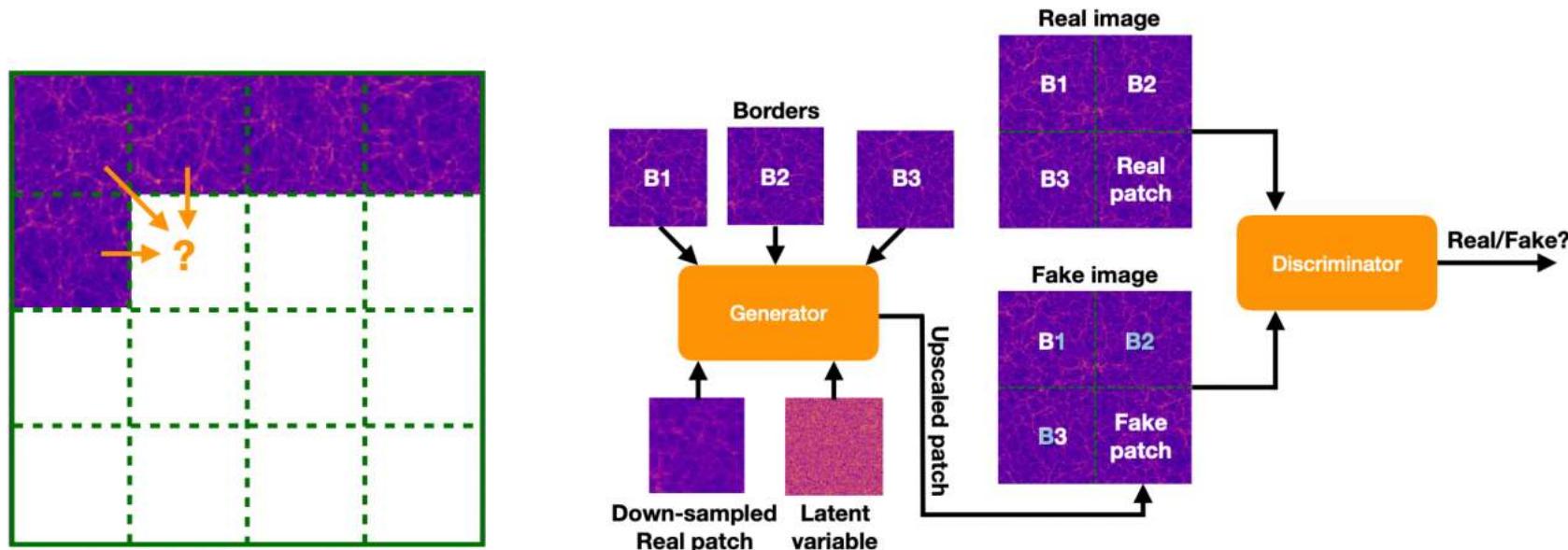
### GENERATOR LOSS

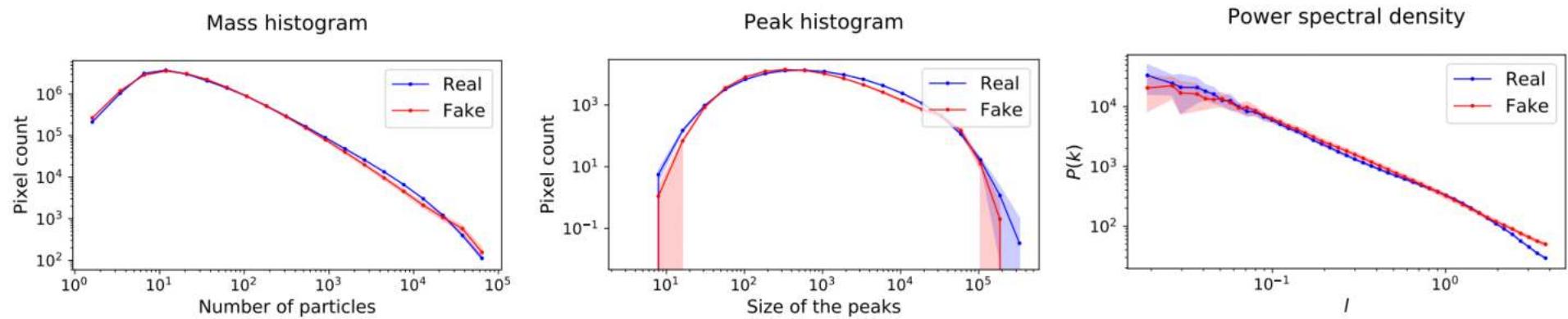
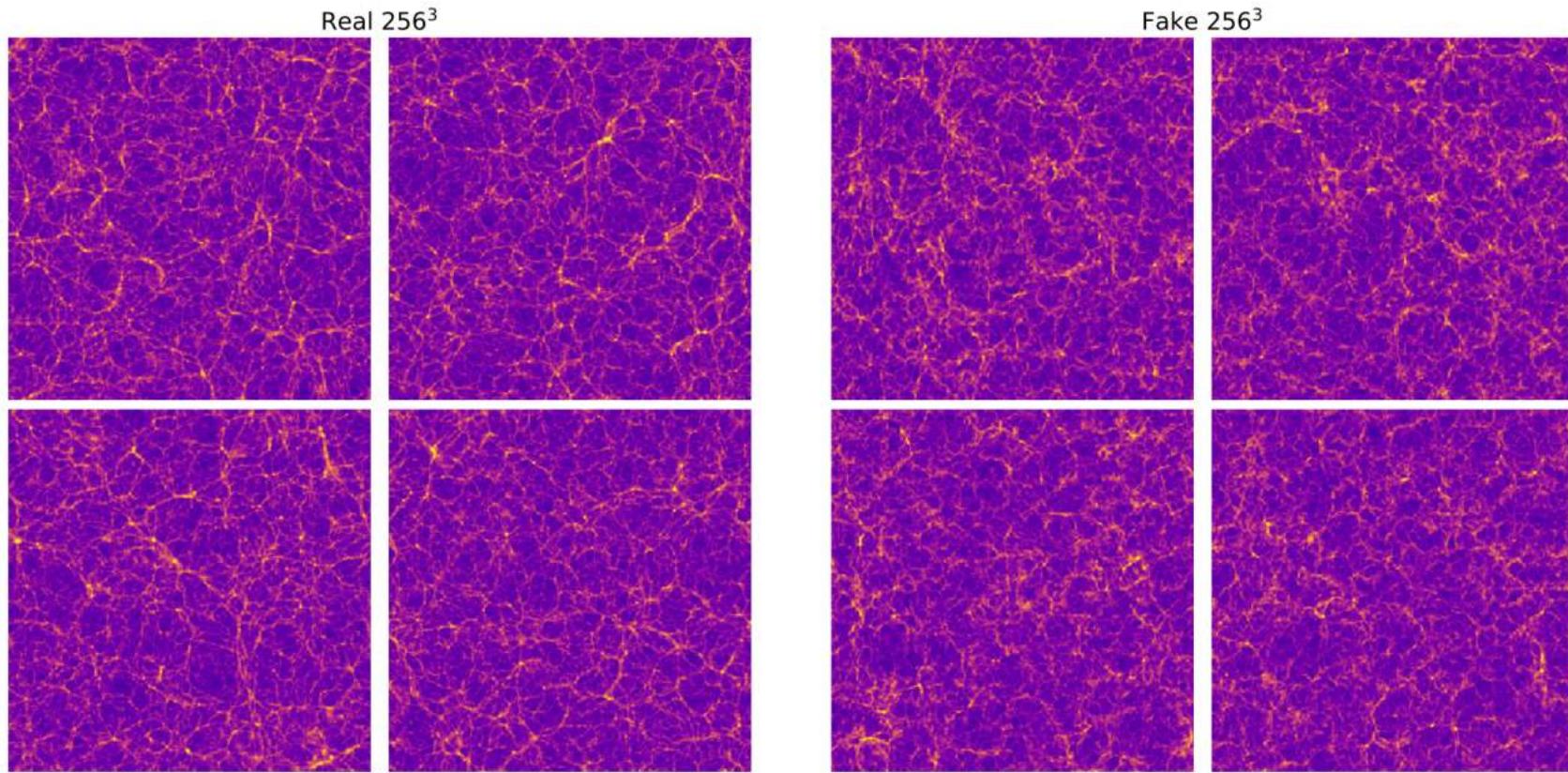
# GANs have remained state-of-the-art for data generation (until diffusion...)

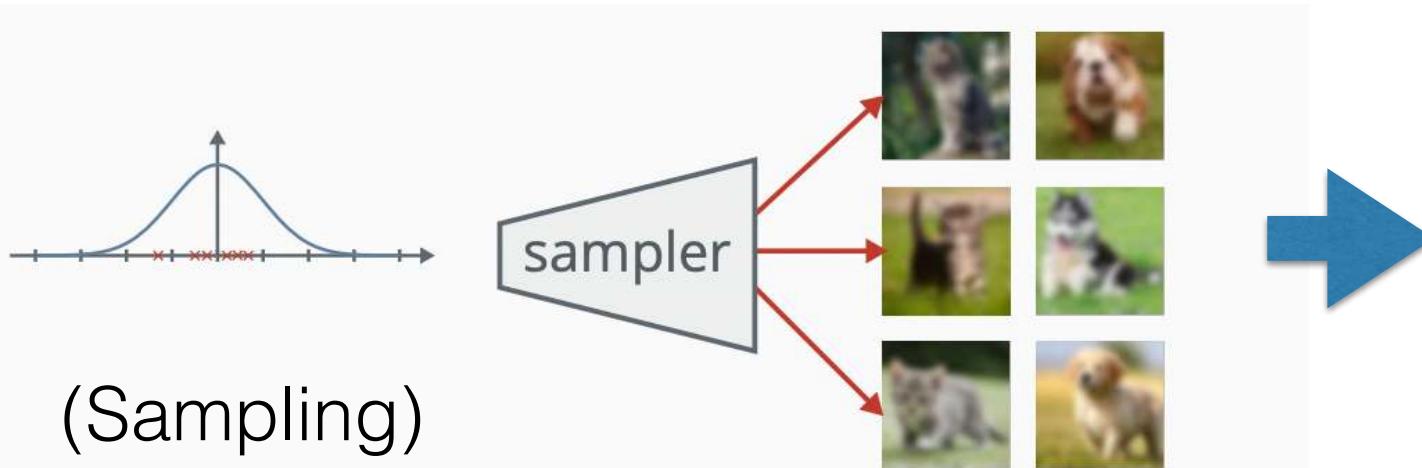


Karras+19

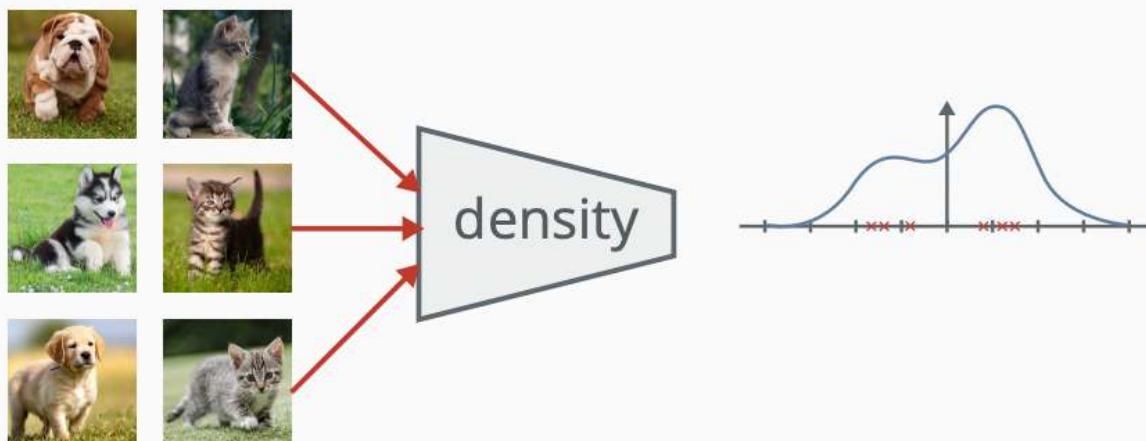
# N-body emulation by Deep Generative Modelling







easy with  
VAEs and GANs



~ VAEs  
difficult with  
GANs

(Density estimation)

# WHAT IS AN ANOMALY OR OUTLIER?

$p(X)$



*your data follows some  
probability distribution  $p$*

# WHAT IS AN ANOMALY OR OUTLIER?

$$p(X)$$



*your data follows some probability distribution  $p$*

Then an object will be anomalous if:

$$p(x_i) < \epsilon$$

# WHAT IS AN ANOMALY OR OUTLIER?

$p(X)$



*your data follows some  
probability distribution  $p$*

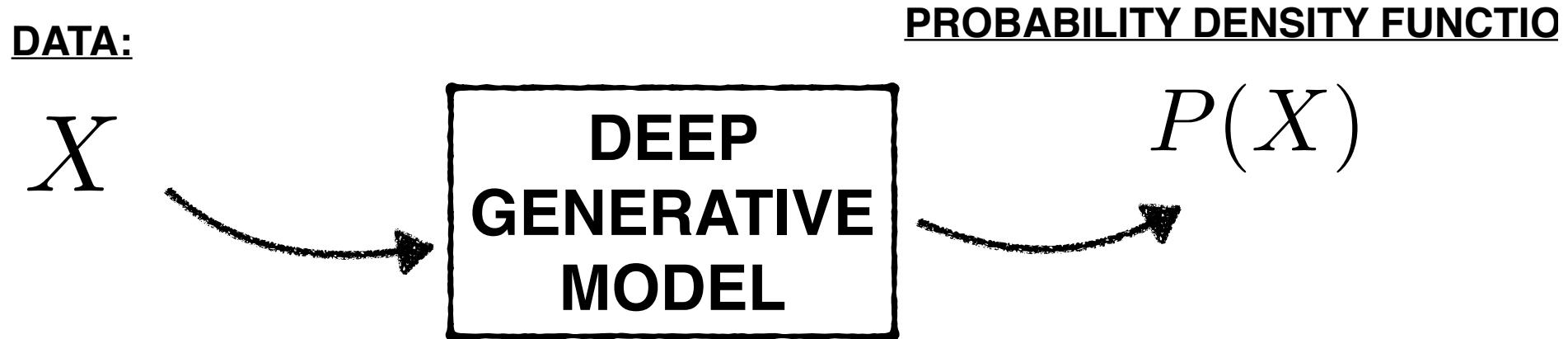


Then an object will be anomalous if:

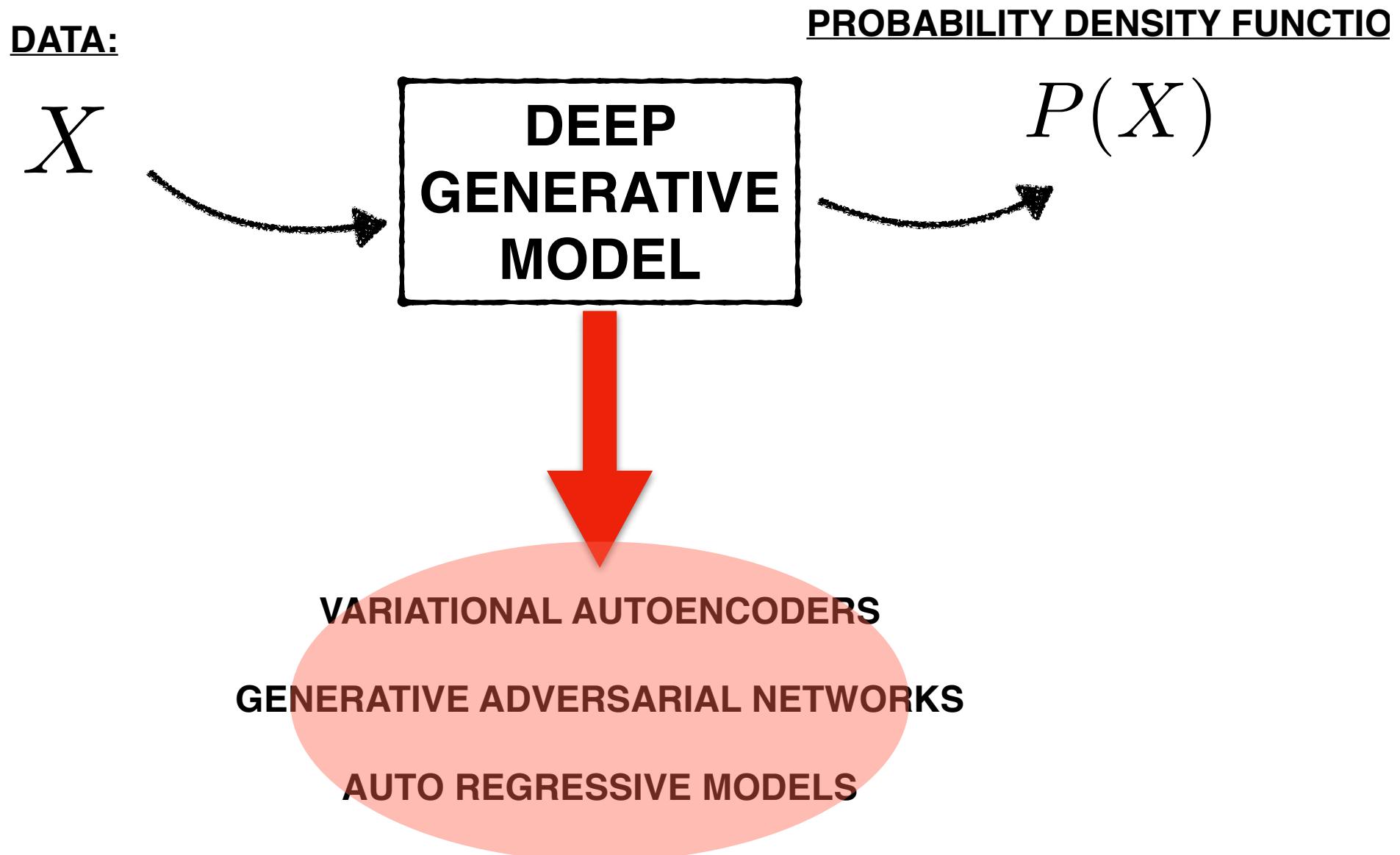
$$p(x_i) < \epsilon$$

**HOW DO WE COMPUTE THE PROBABILITY  
DISTRIBUTION  $p$ ?**

## GENERATIVE MODELS DO PRECISELY THAT:



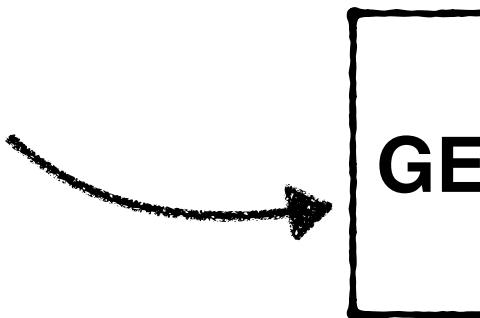
## GENERATIVE MODELS DO PRECISELY THAT:



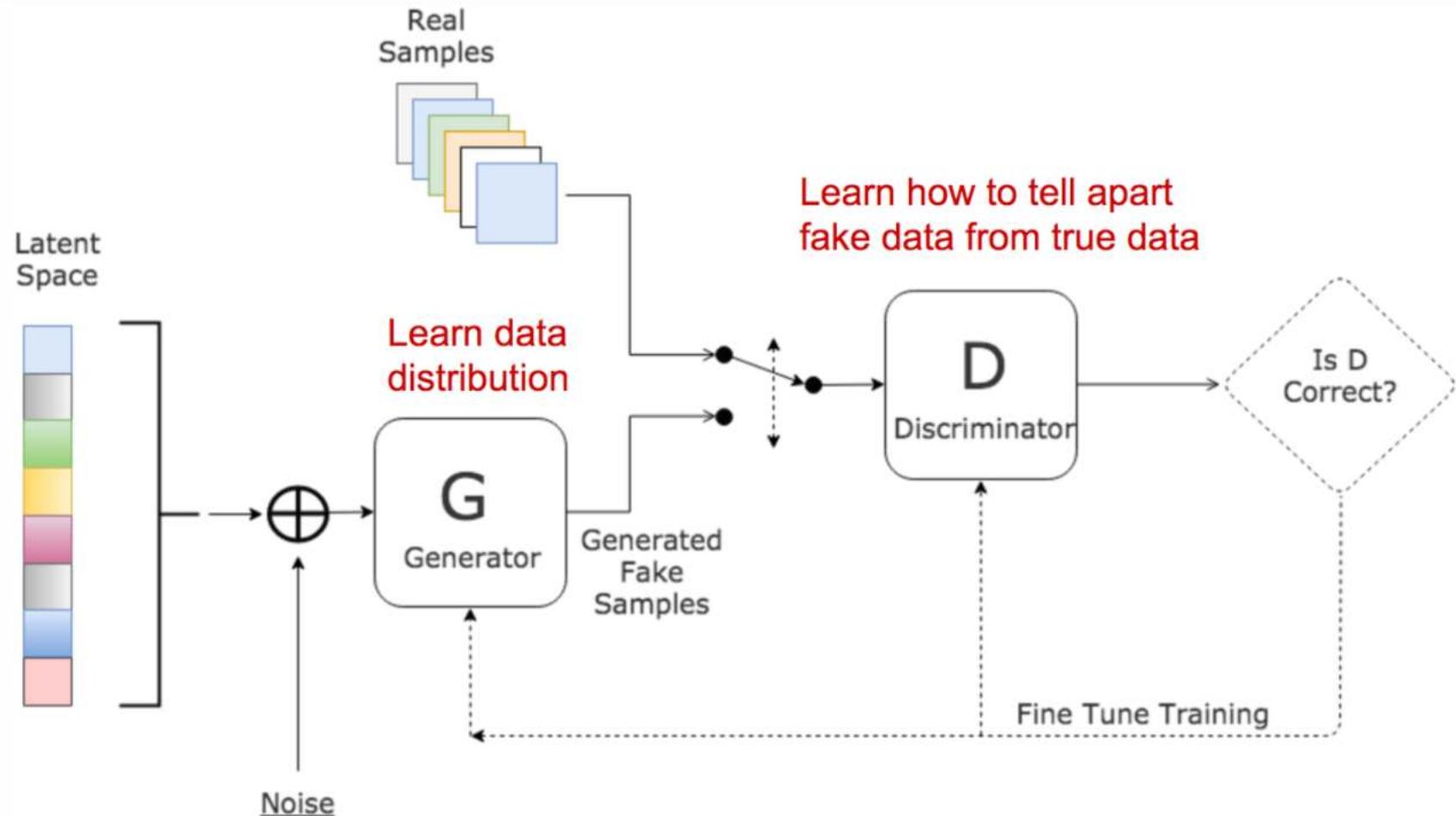
# GENERATIVE MODELS DO PRECISELY THAT:

DATA:

$X$



# USE GENERATIVE MODELING TO LEARN $P(X)$ [NORMAL DATA]



(Image source: [www.kdnuggets.com/2017/01/generative-...-learning.html](http://www.kdnuggets.com/2017/01/generative-...-learning.html))

**Hyper Suprime Cam Survey  
(HSC)**  
**(Subaru telescope 1400 sq.  
deg, r~26, 436 million**

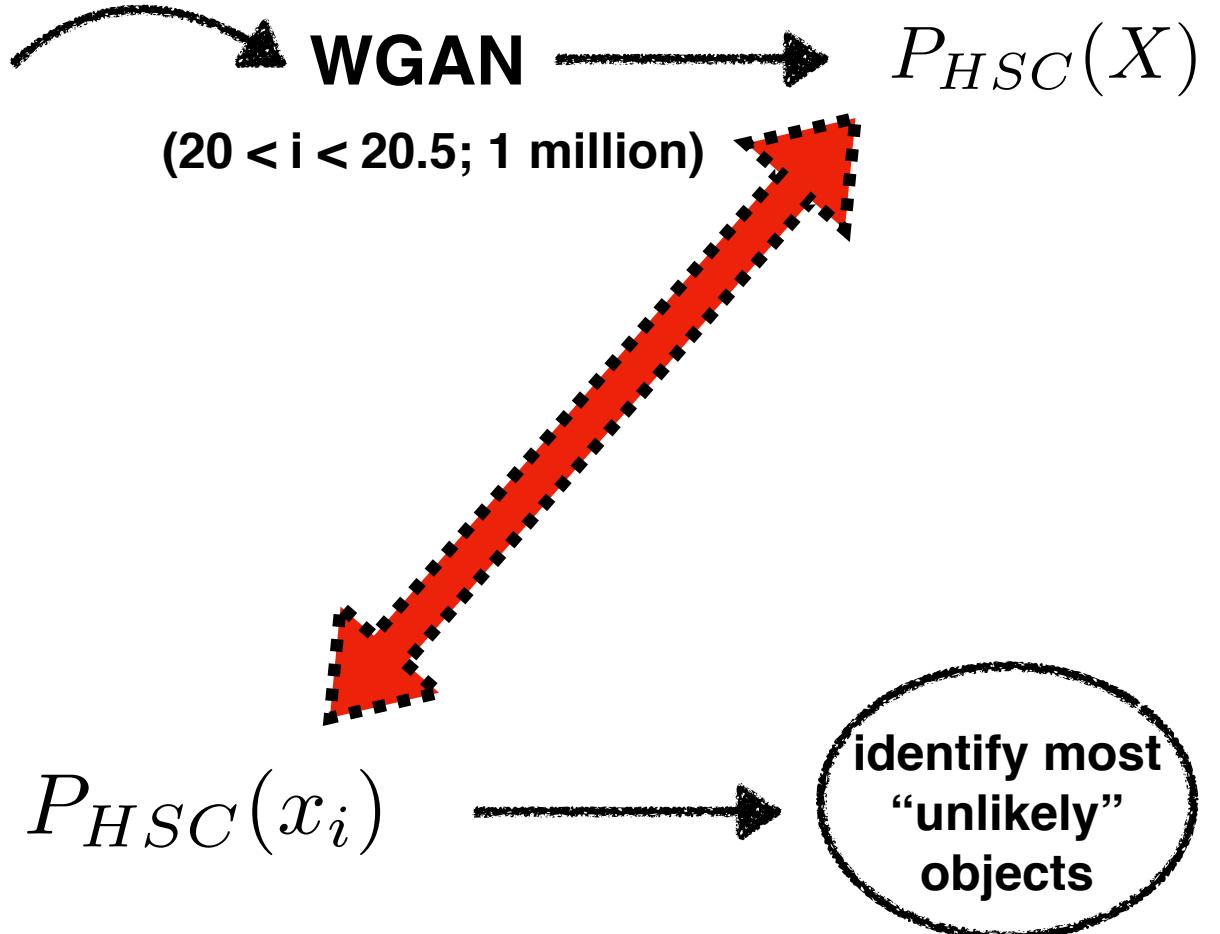


WGAN  $P_{HSC}(X)$   
( $20 < i < 20.5$ ; 1 million)

**Hyper Suprime Cam Survey  
(HSC)**  
**(Subaru telescope 1400 sq.  
deg, r~26)**

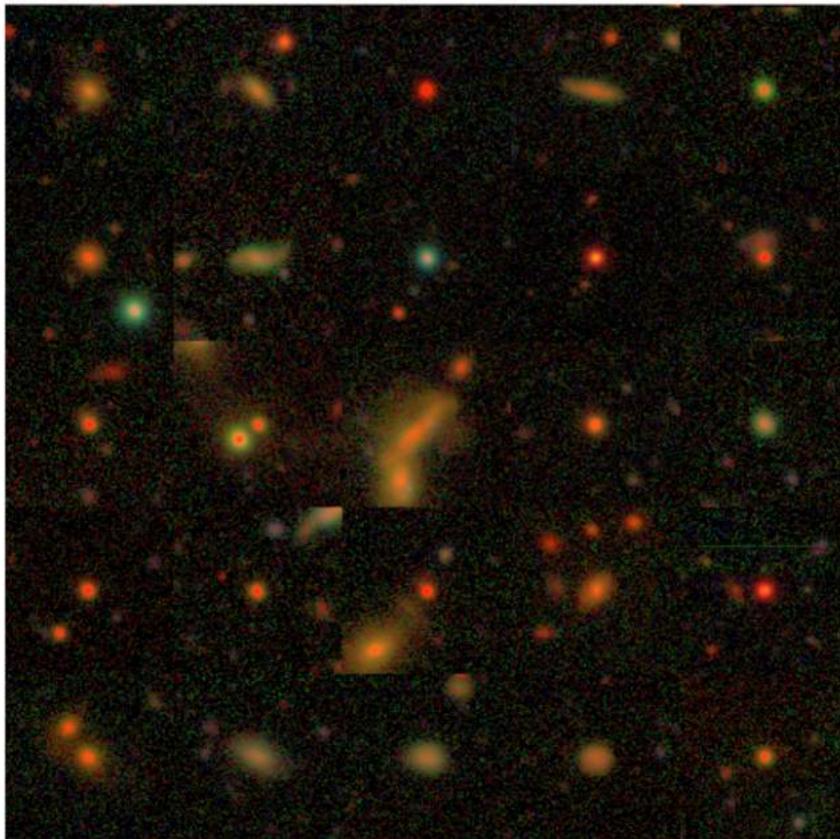


$x_i$  →  
**(individual detection)**

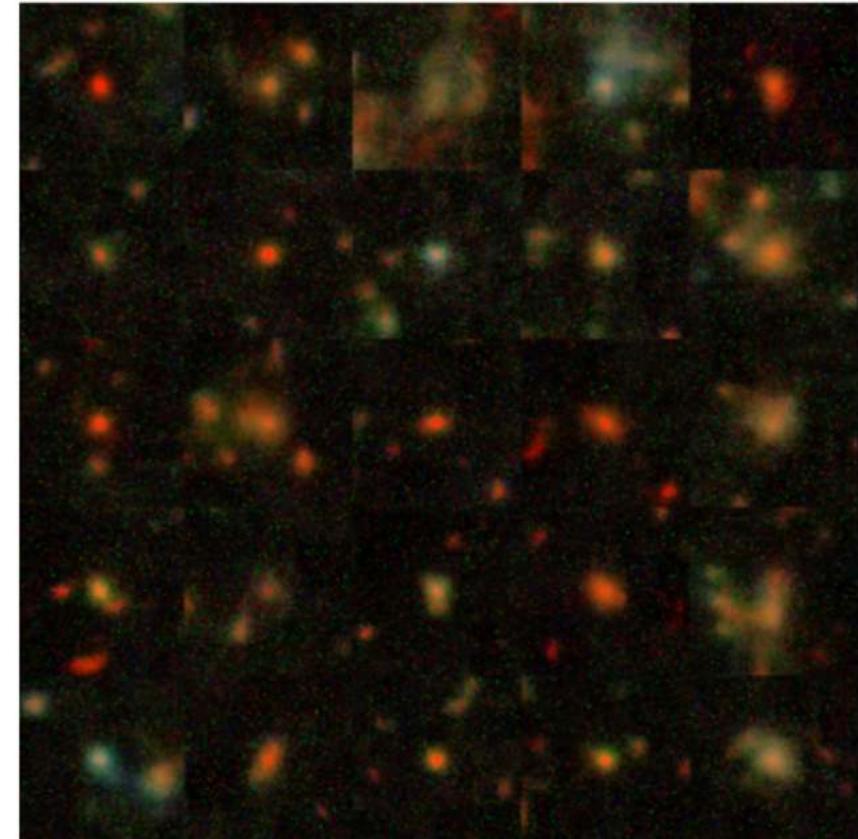


**Storey-Fisher, MHC+21**  
**Margalef-Bentabol, MHC+20**

**REAL**



**WGAN GENERATED**



# COMPUTE ANOMALY SCORE BASED ON WGAN RECONSTRUCTION ERROR

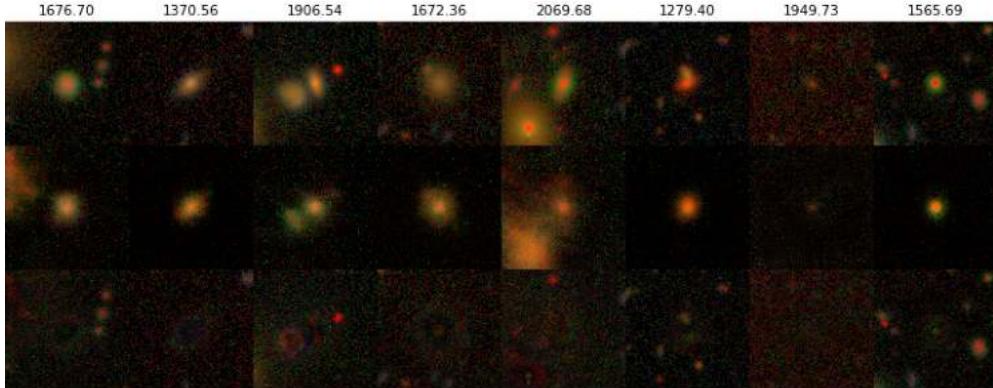
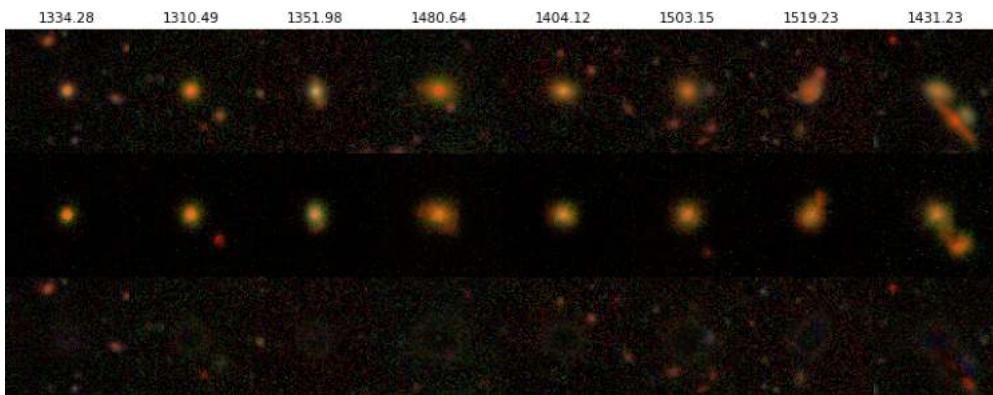
ANOMALY SCORE:

$$AS = \lambda G + (1 - \lambda)C$$

Hyperparameter

RMS btw input  
and generated

RMS btw critic features  
from input and generated



Schlegl+17

Storey-Fisher+21

Margalef-Bentabol, MHC+20

# COMPUTE ANOMALY SCORE BASED ON WGAN RECONSTRUCTION ERROR

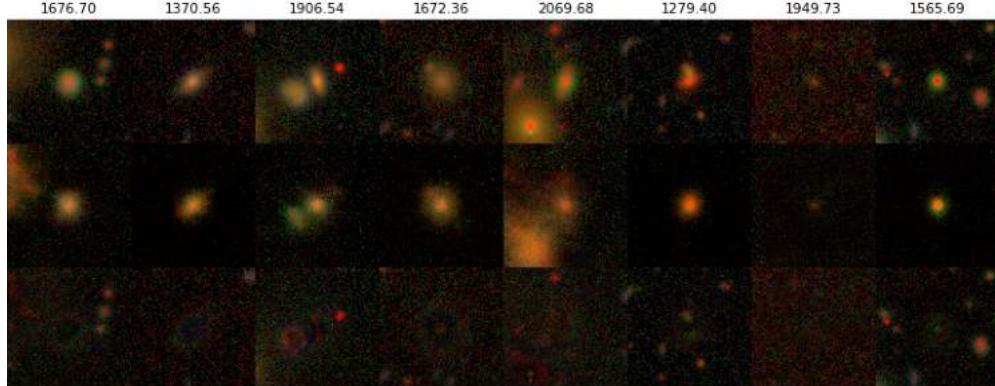
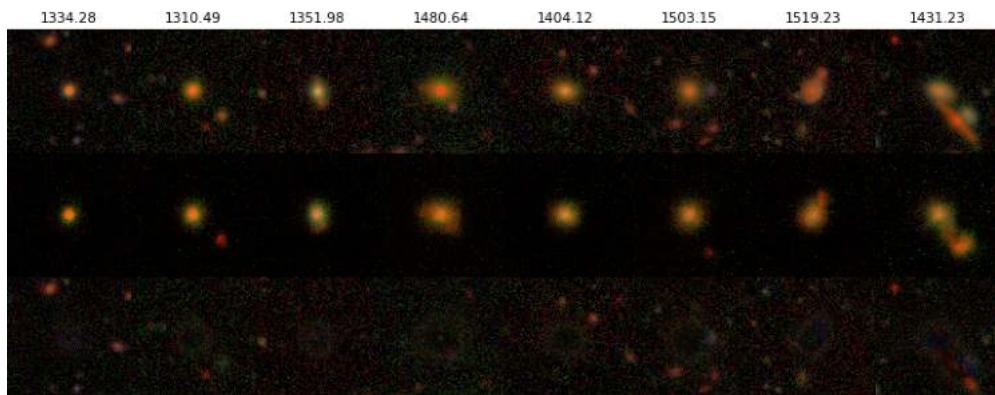
ANOMALY SCORE:

$$AS = \lambda G + (1 - \lambda)C$$

Hyperparameter

RMS btw input  
and generated

RMS btw critic features  
from input and generated



BEST RECONSTRUCTION

RESIDUAL

REAL

BEST RECONSTRUCTION

RESIDUAL

Schlegl+17

Storey-Fisher+20 (in prep)

Margalef-Bentabol, MHC+20

# COMPUTE ANOMALY SCORE BASED ON WGAN RECONSTRUCTION ERROR

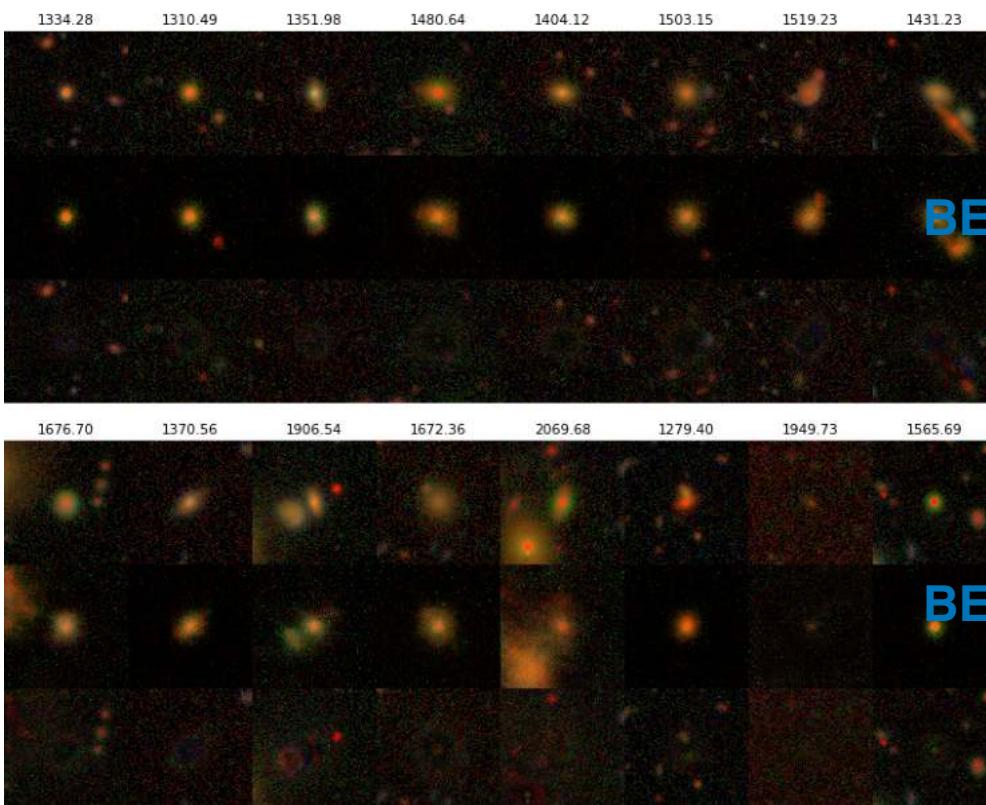
ANOMALY SCORE:

$$AS = \lambda G + (1 - \lambda)C$$

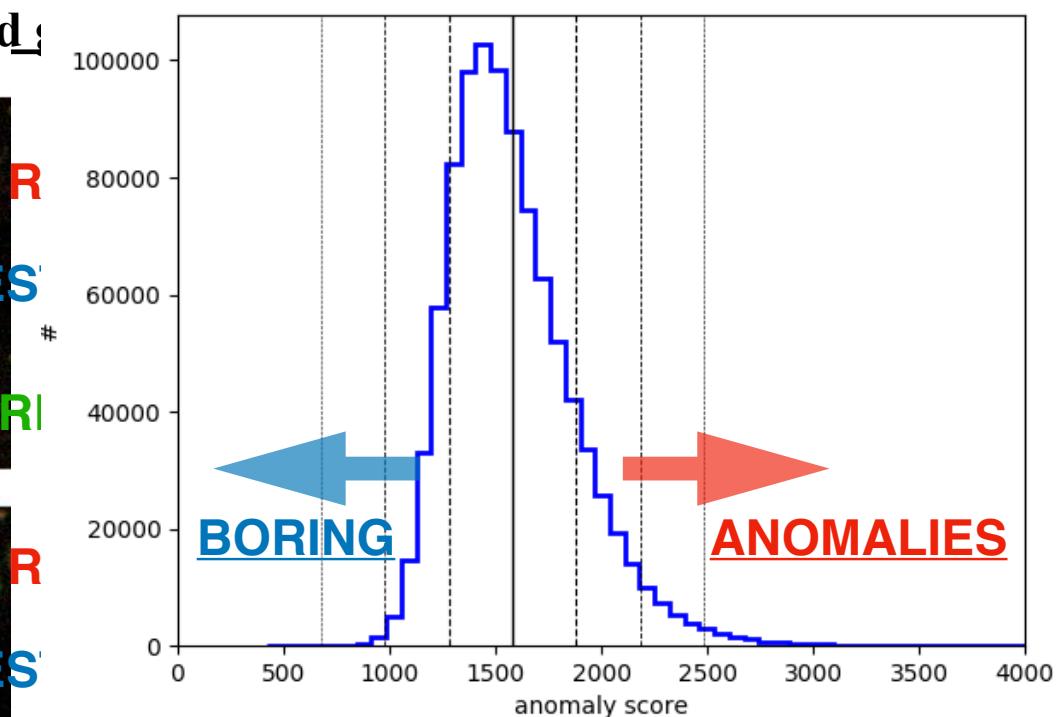
Hyperparameter

RMS  
and

RMS btw critic features



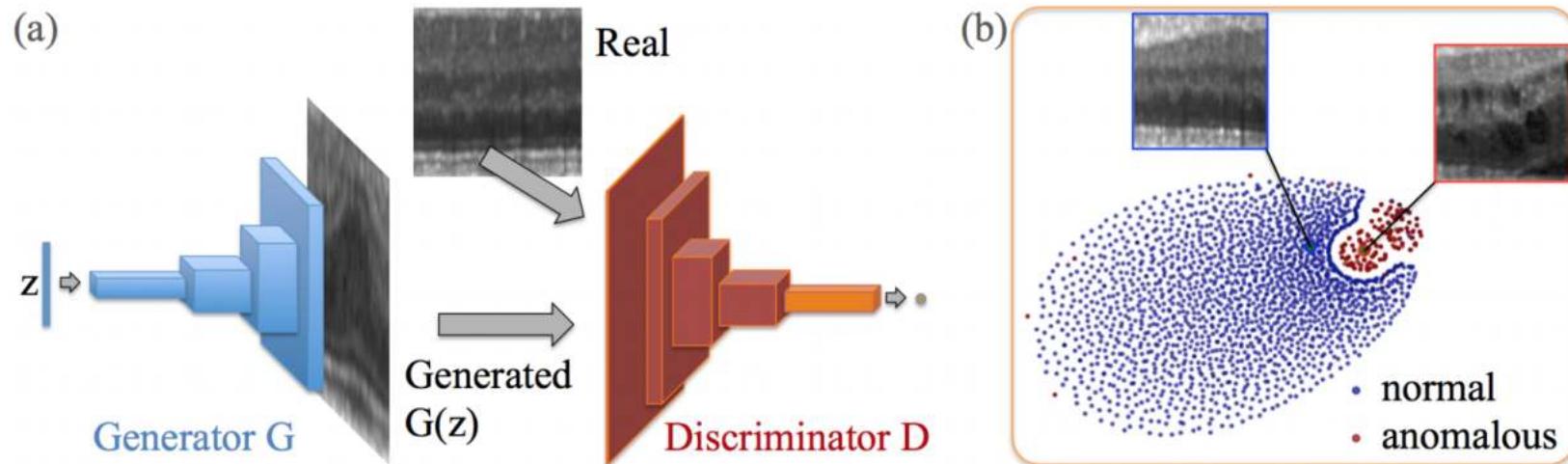
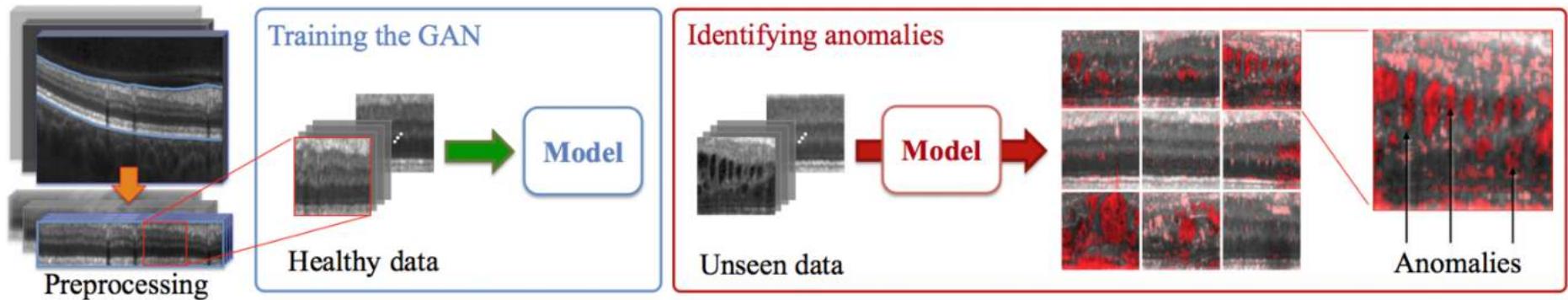
**RESIDUAL**



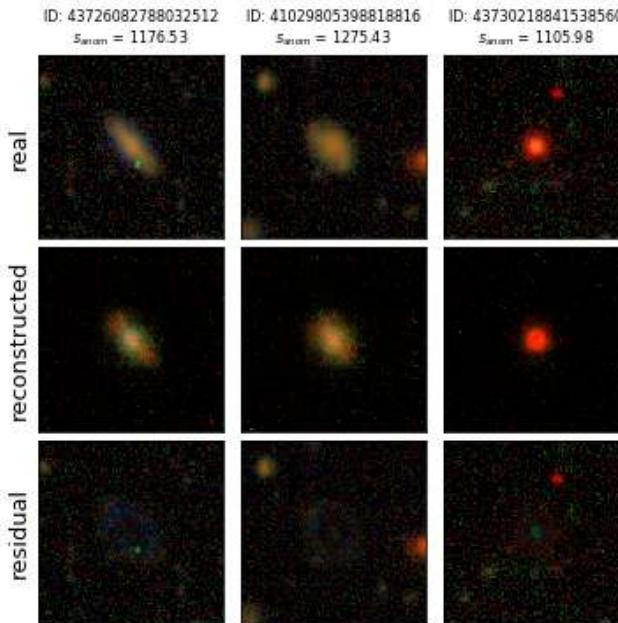
Storey-Fisher+21

Margalef-Bentabol, MHC+20

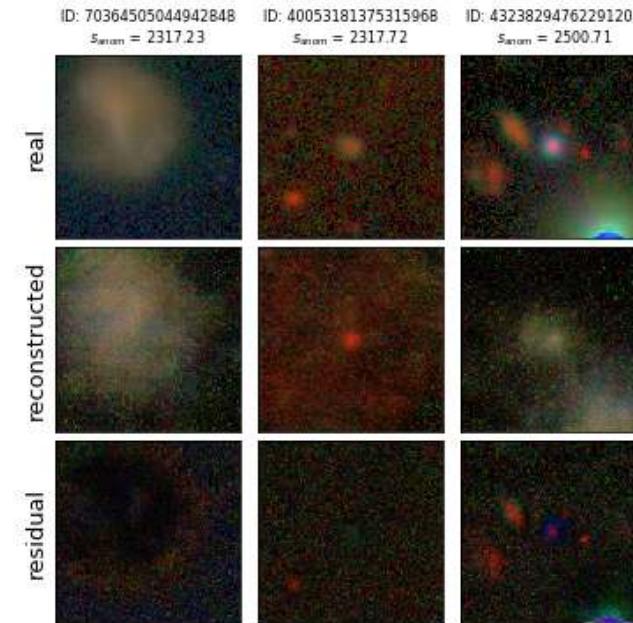
# This is again inspired by biomedical applications



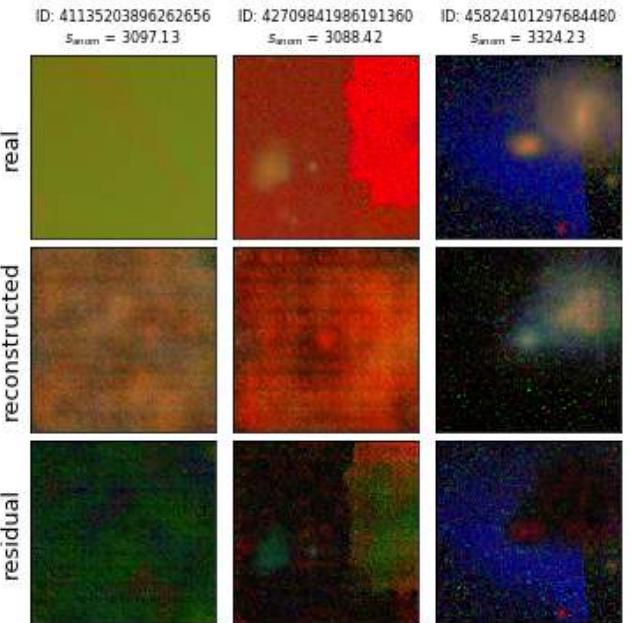
## AVERAGE



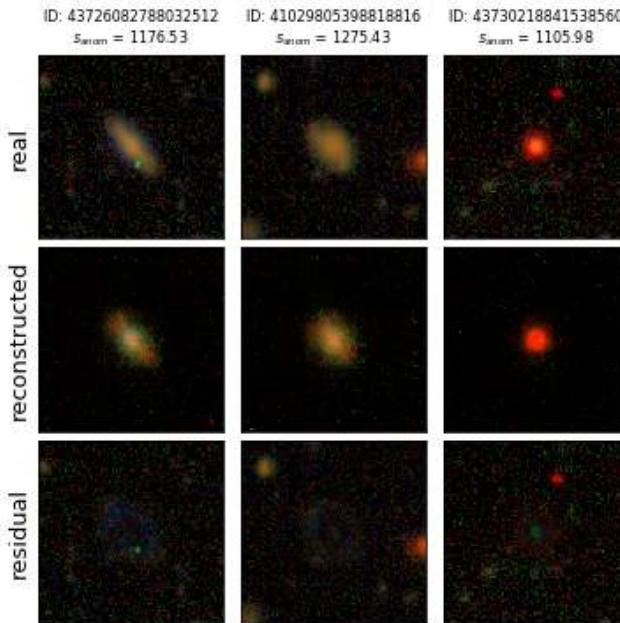
## 1-SIGMA



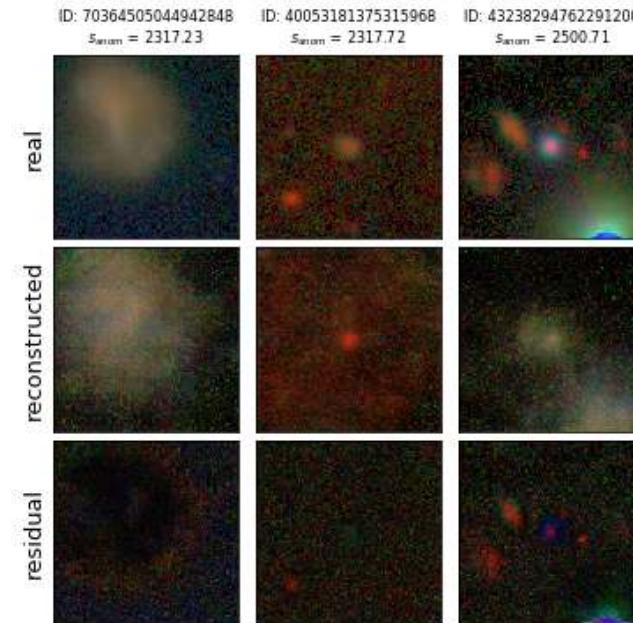
## >3-SIGMA



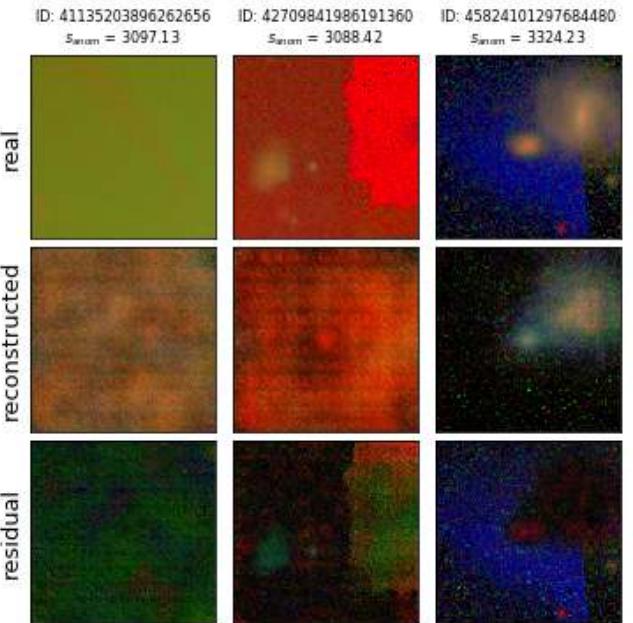
## AVERAGE



## 1-SIGMA

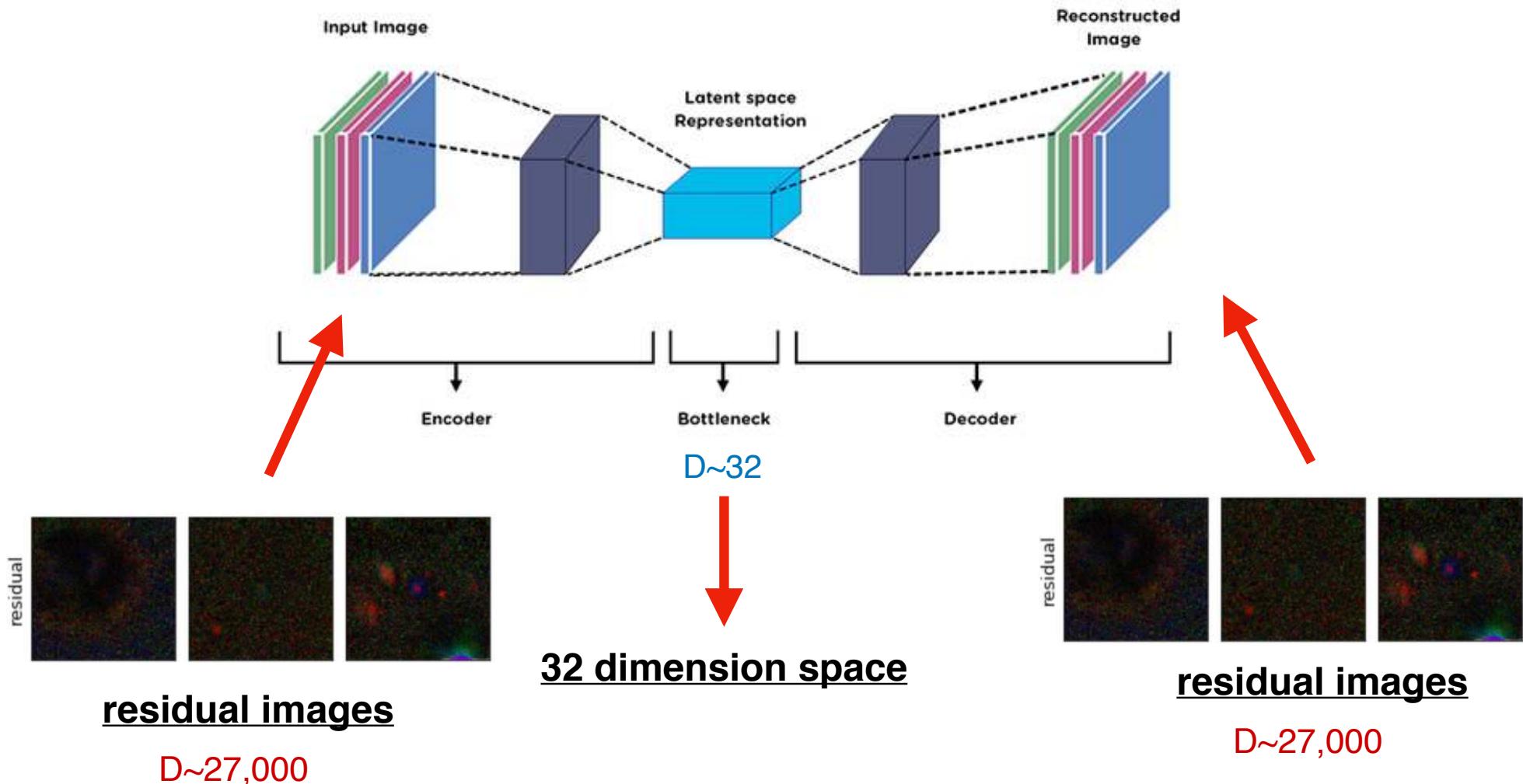


## >3-SIGMA



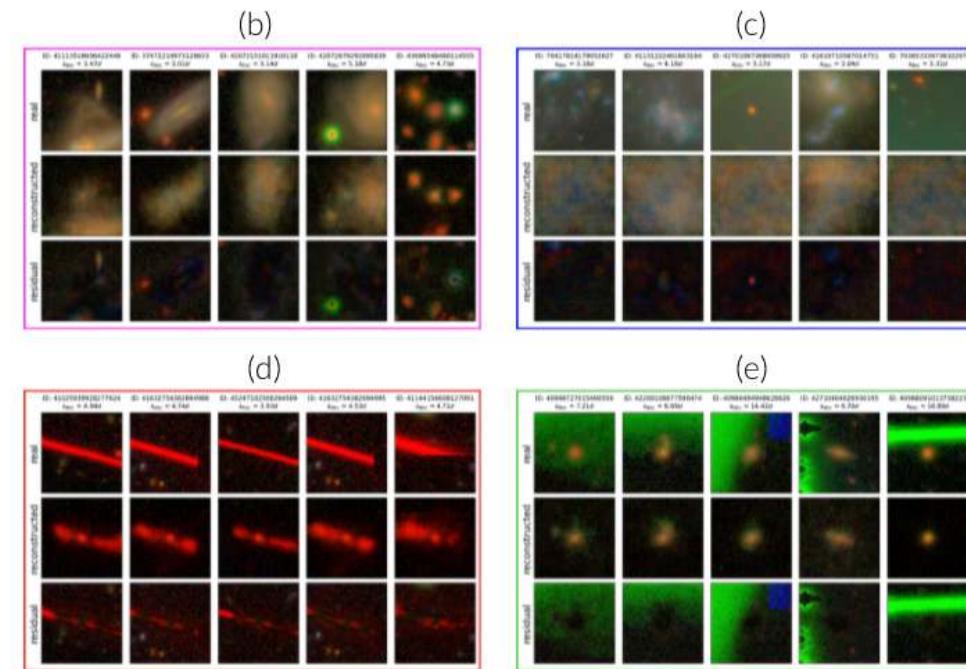
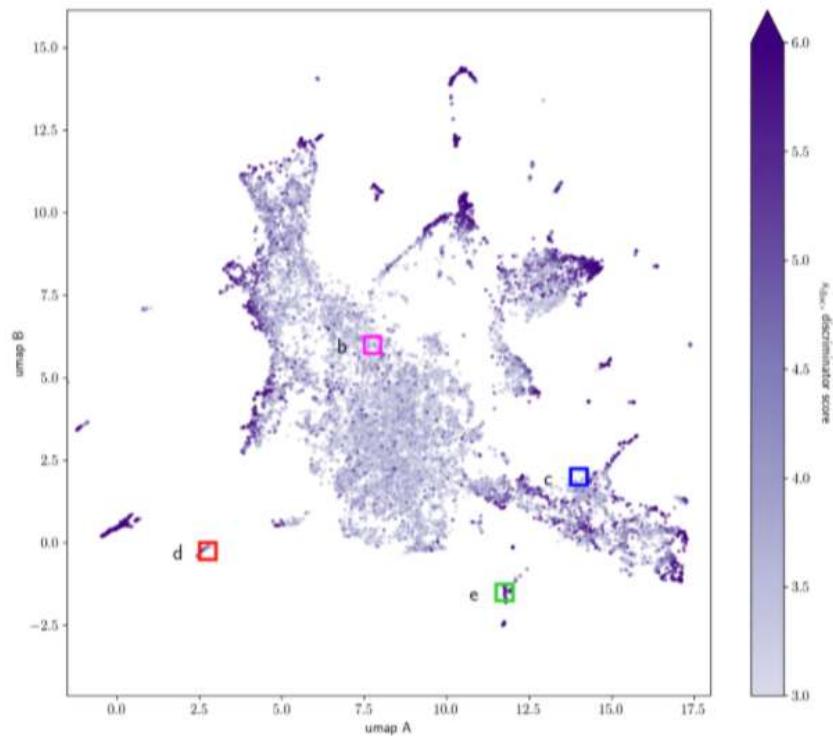
**HOW DO WE FILTER OUT “INTERESTING”  
ANOMALIES?**

# HUNTING “INTERESTING” OUTLIERS

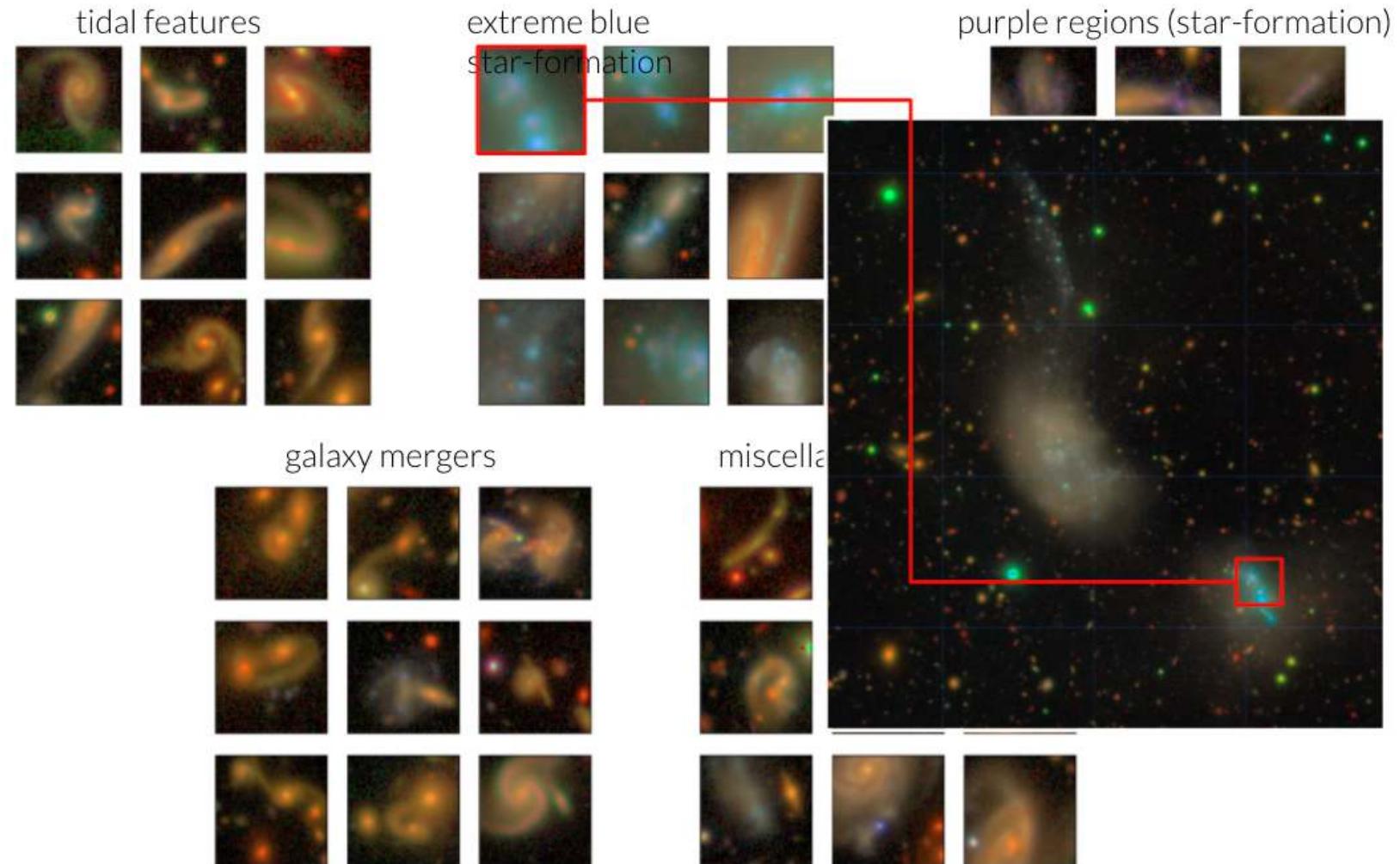


AUTO-ENCODERS TO REDUCE THE DIMENSIONALITY  
OF RESIDUAL IMAGES

# UMAP of high-anomaly sample ( $s_{\text{disc}} > 3\sigma$ )



# Detected interesting anomalies



### 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

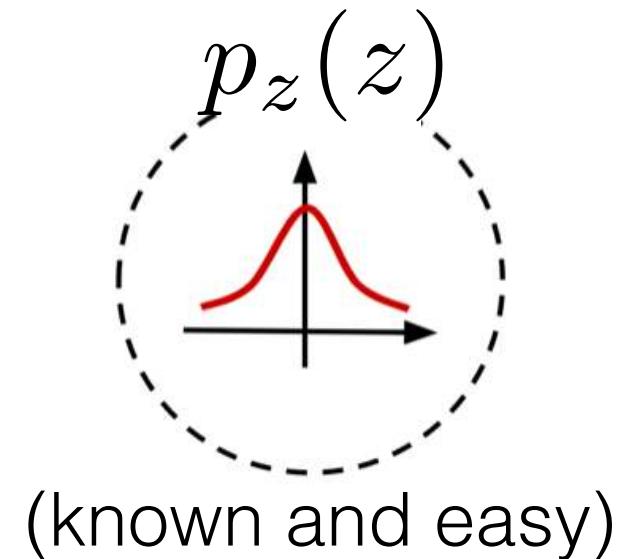
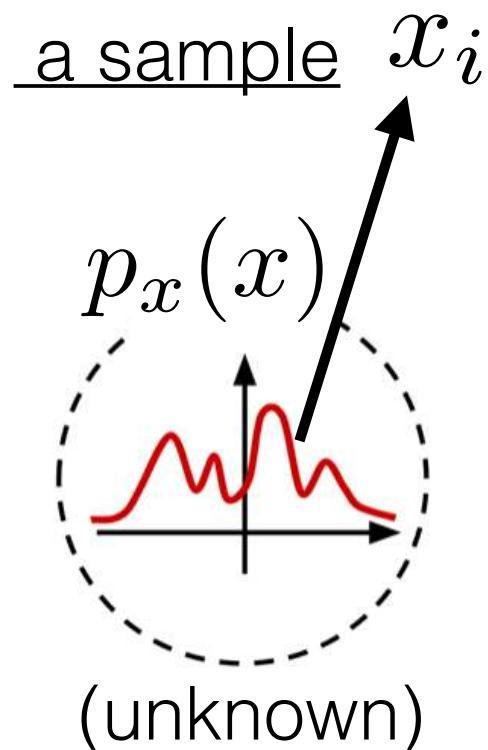
Invertible differentiable function

Determinant of the Jacobian of  $f(x)$

tractable pdf  
(typically Gaussian)

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

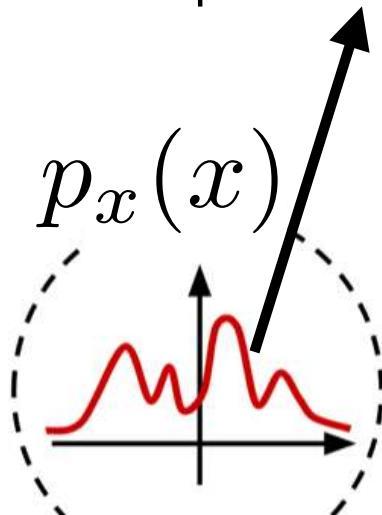
Can represent  $p_X(x)$  in terms of  $f(x)$  and  $p_Z(x)$   
(notice there is no dimensionality reduction as opposed to VAEs )



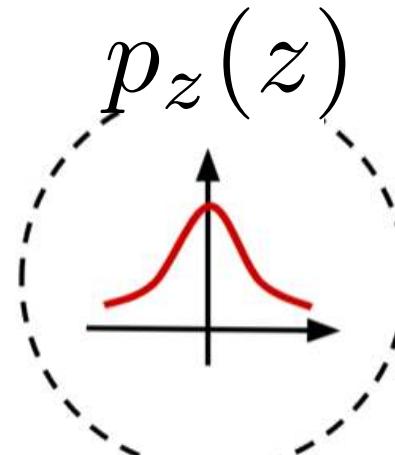
$$p_X(x) = p_Z(f(x)) |\det J_{f(x)}|$$

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  $x_i \longrightarrow f_w(x_i)$



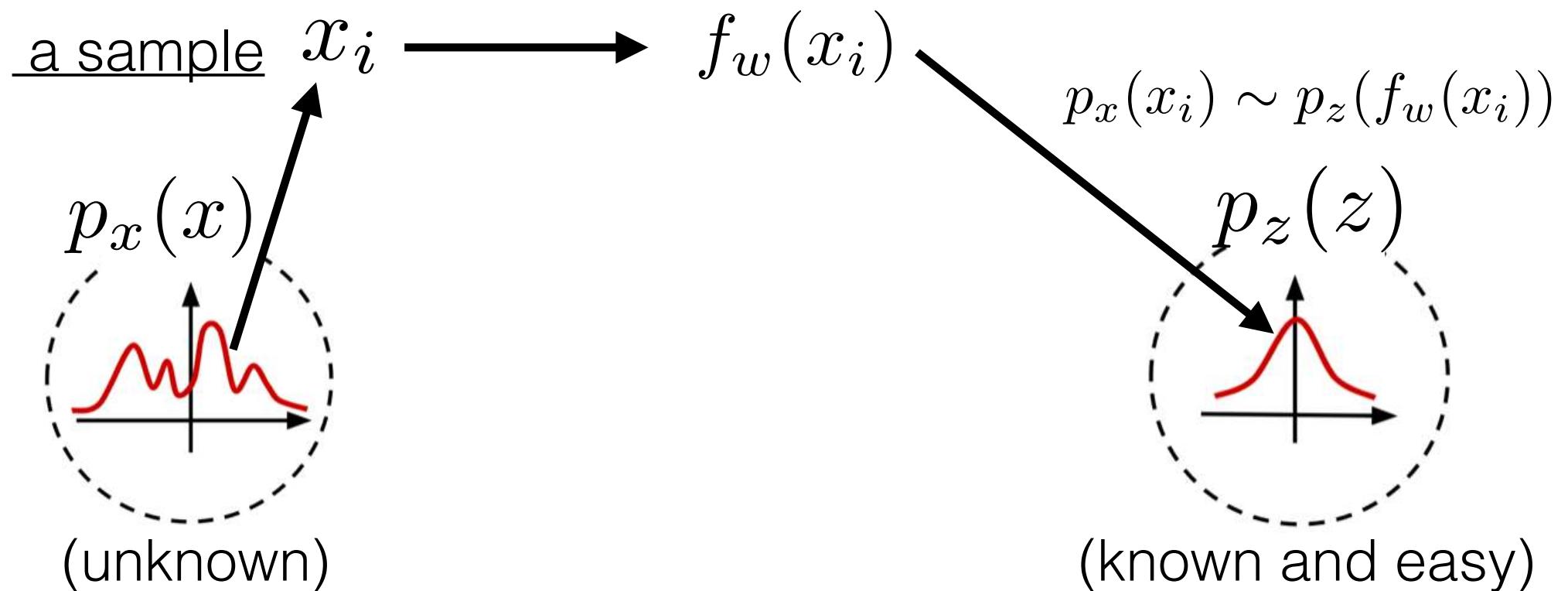
(unknown)



(known and easy)

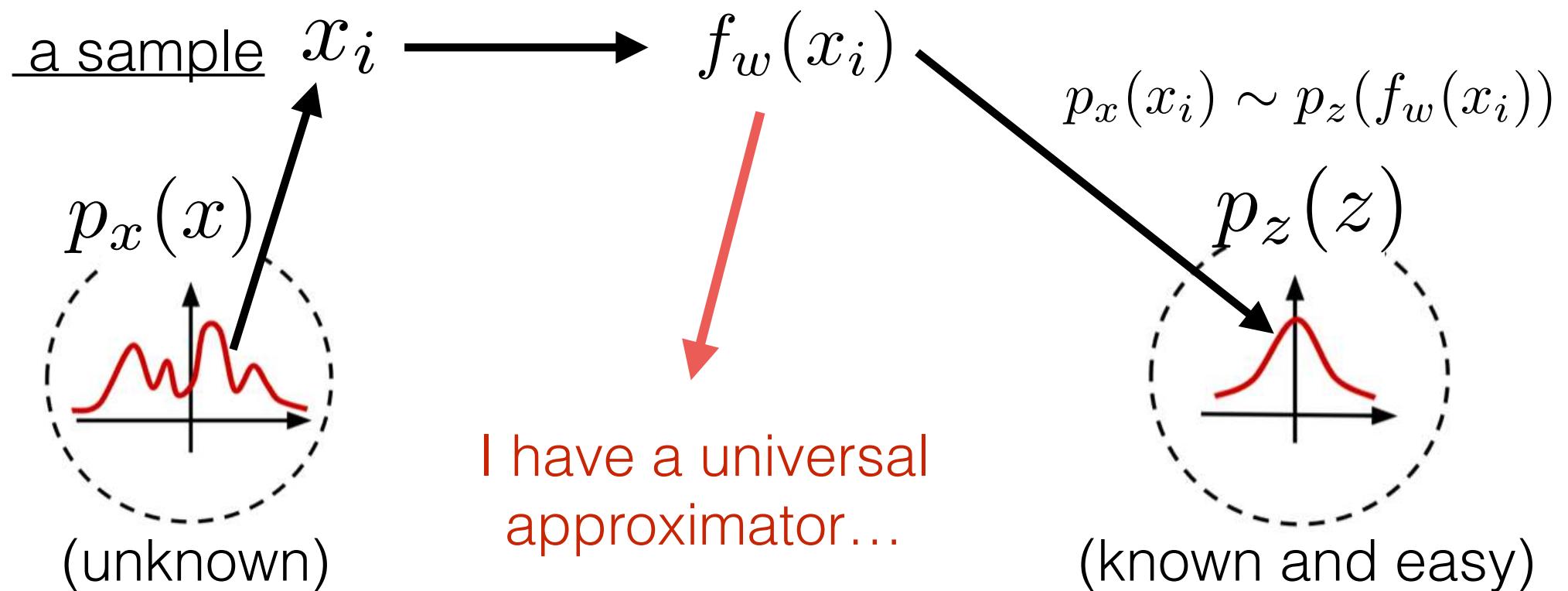
$$p_X(x) = p_Z(f(x)) |\det J_{f(x)}|$$

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



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

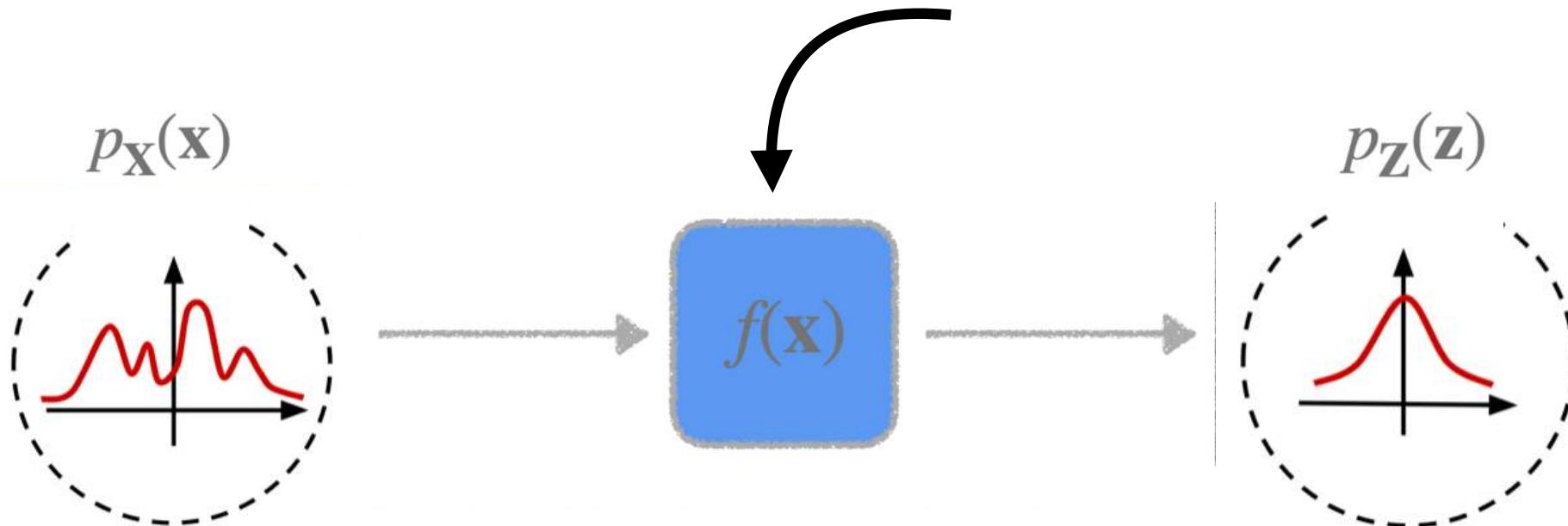
Can represent  $p_X(x)$  in terms of  $f(x)$  and  $p_Z(x)$   
(notice there is no dimensionality reduction as opposed to VAEs )



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

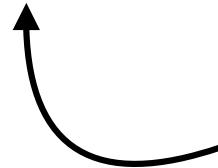
Can represent  $p_X(x)$  in terms of  $f(x)$  and  $p_Z(x)$   
(notice there is no dimensionality reduction as opposed to VAEs )

Learn this transformation with a NN  
(remember it's differentiable)



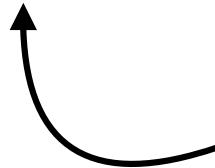
Then  $f(x)$  becomes:

$$f(x; W)$$



NN learnable weights

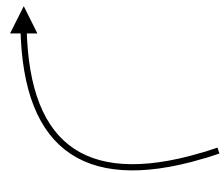
Then  $f(x)$  becomes:  $f(x; W)$



NN learnable weights

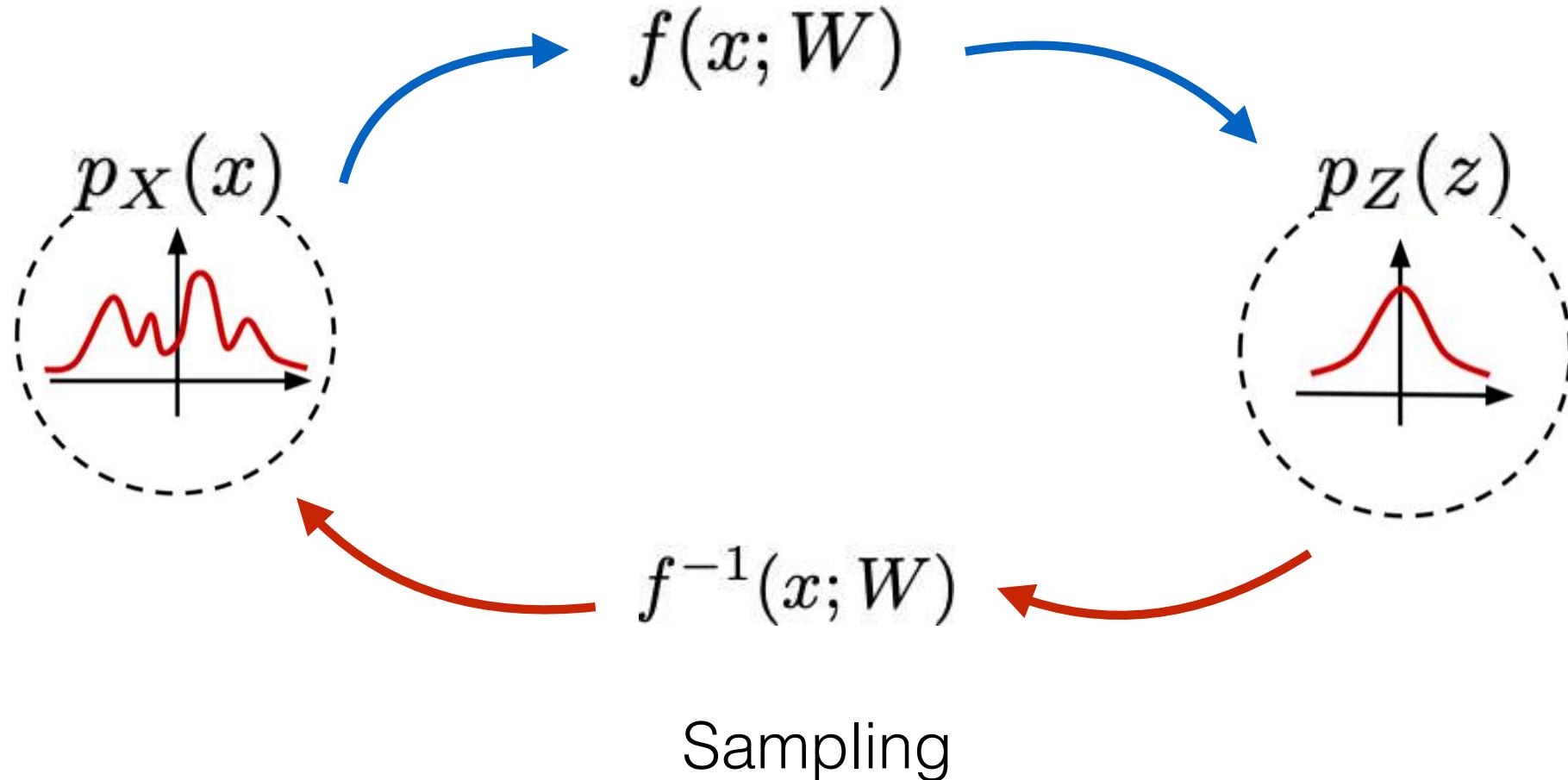
And the loss function: log likelihood

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

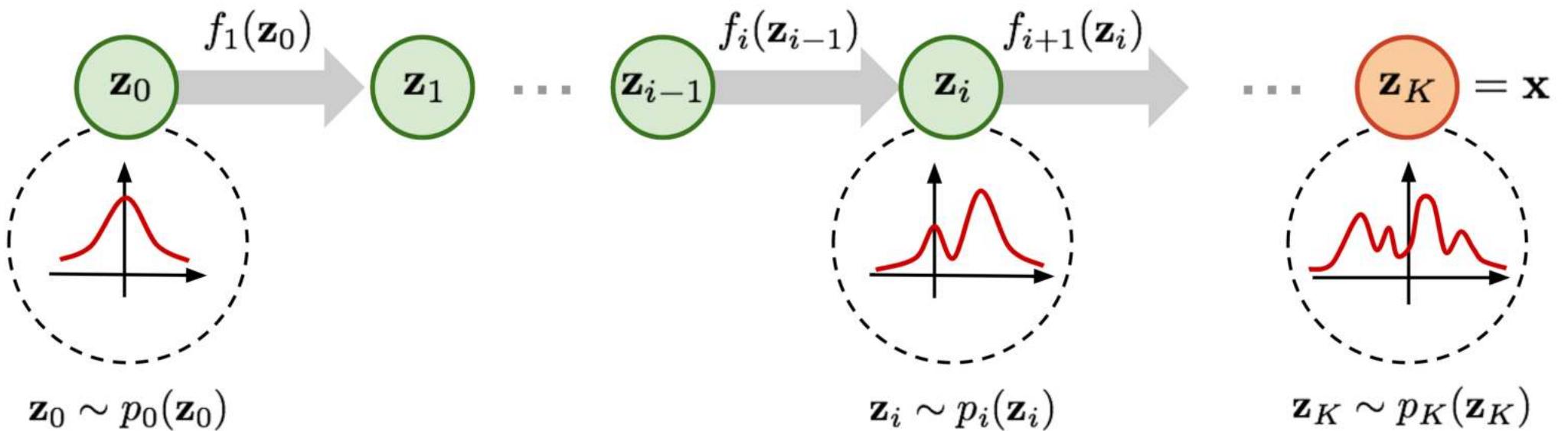


This is a gaussian

## Likelihood evaluation



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

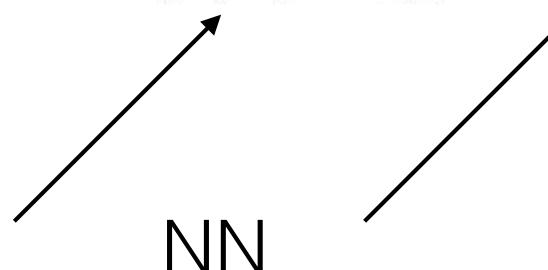
# What functions satisfy these conditions?

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

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

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

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



- Easily invertible:

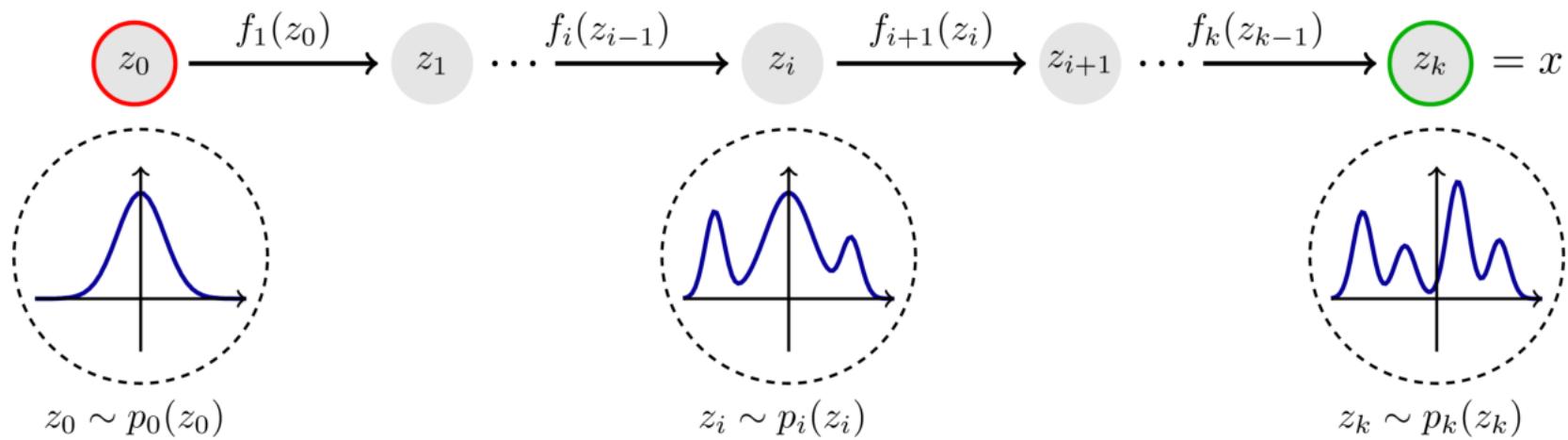
$$\begin{cases} \mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} &= \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} &= (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot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

- Jacobian:

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

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

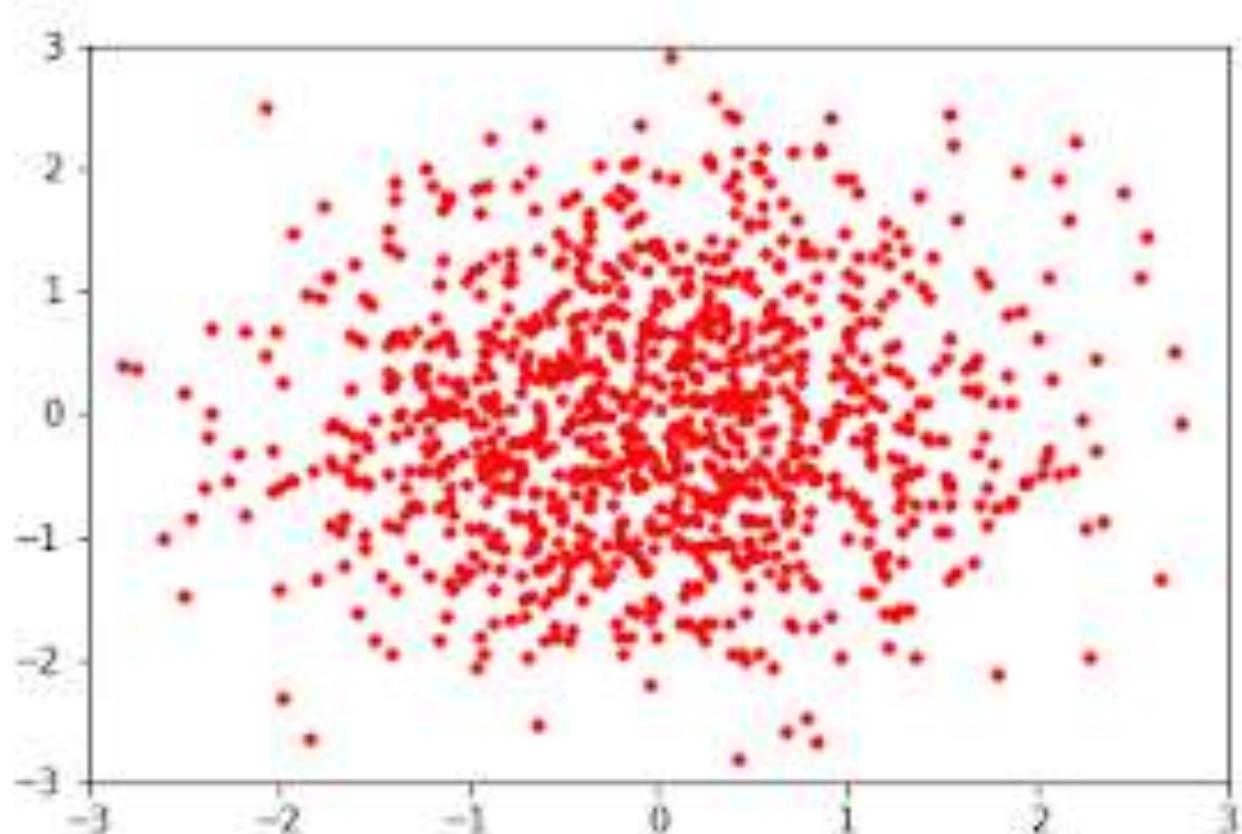
**normalizing flows** are generative models that are easy to evaluate and flexibly expressive



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

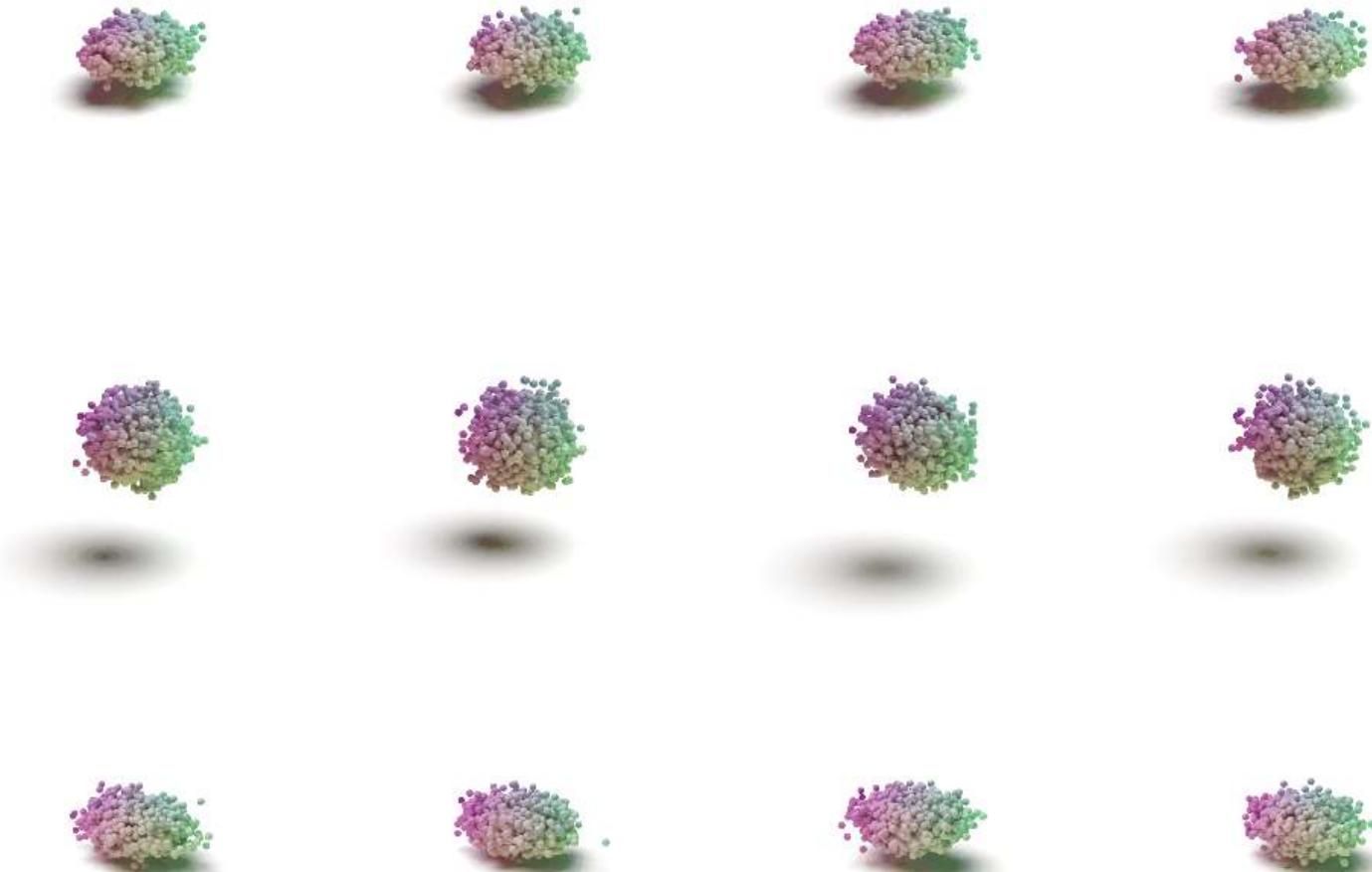
$f = f_1 \circ f_2 \dots \circ f_{k-1} \circ f_k$  is also invertible and differentiable

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



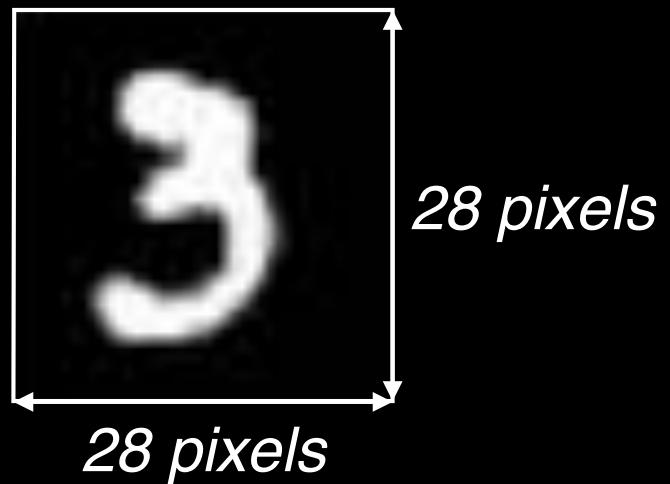
*credit: Eric Jang*

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



3 4 2 1 9 5 6 2 / 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*  
 $p(\text{pixels}) \approx q_\phi(\text{pixels})$

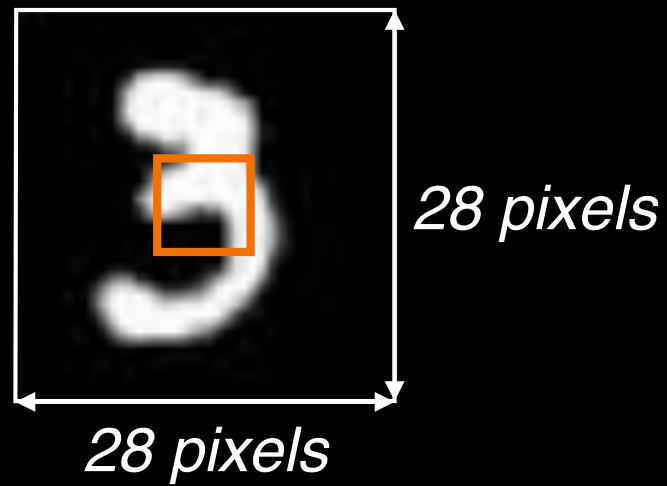


sample  $q_\phi(\text{pixels})$

sample  $q_\phi$ (pixels)

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

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



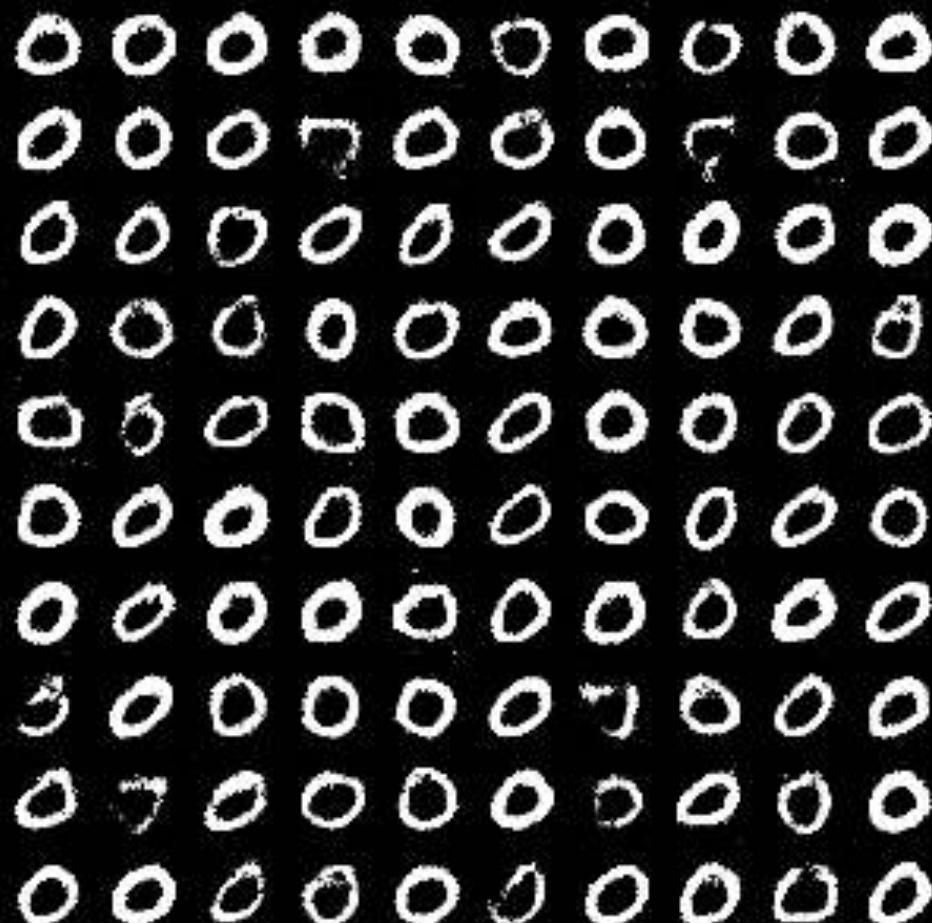
$$p(\text{pixels} \mid \text{central pixels}) \approx q_{\phi}(\text{pixels} \mid \text{central pixels})$$

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



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

$$p(\text{pixels} \mid \text{central pixels}) \sim \frac{q_{\theta}}{X}$$

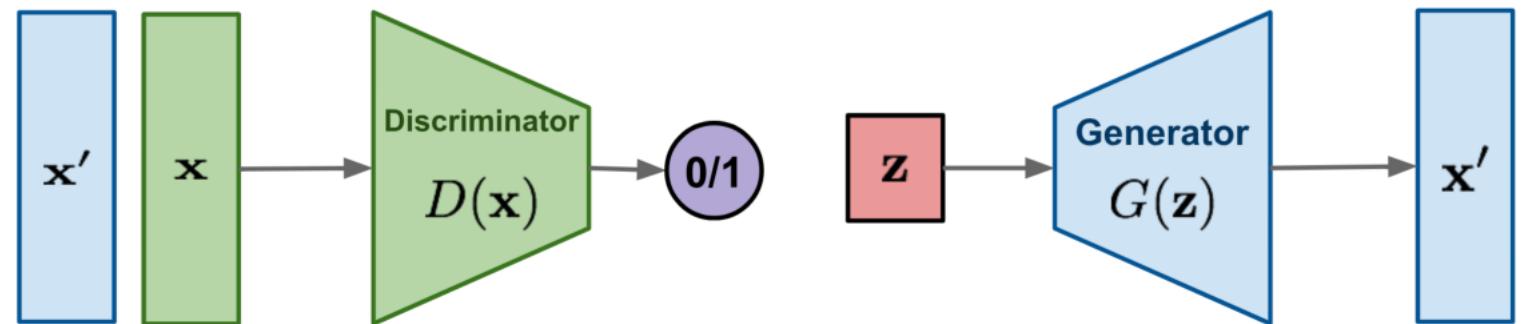
# GANs AND VAEs ARE VERY POWERFUL BUT DO NOT PROVIDE AN EXPLICIT LIKELIHOOD

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

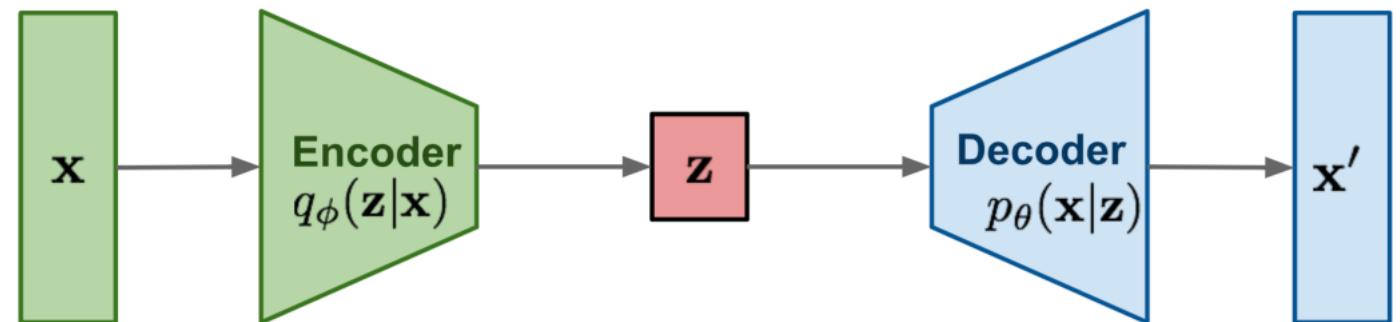
Grathwohl+18

\* there is a third member of the family now: diffusion models which are as accurate as GANs and provide exact likelihood. (gradient)

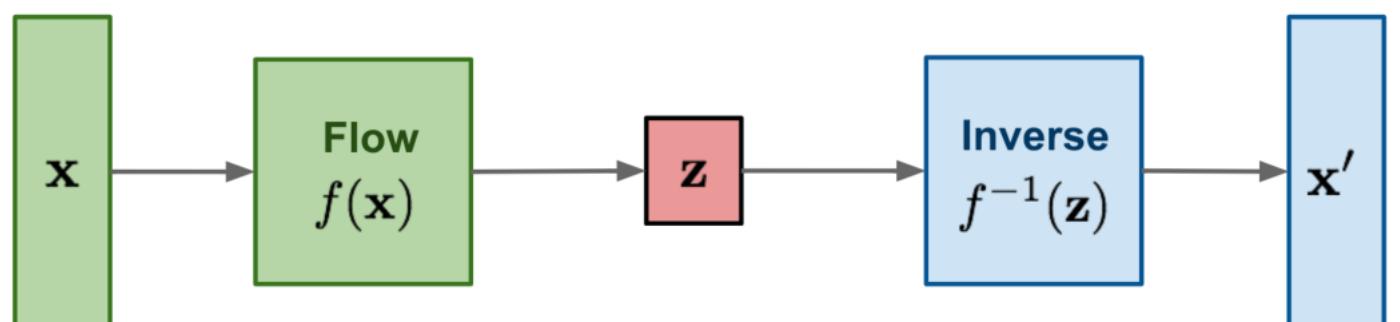
**GAN:** minimax the classification error loss.



**VAE:** maximize ELBO.

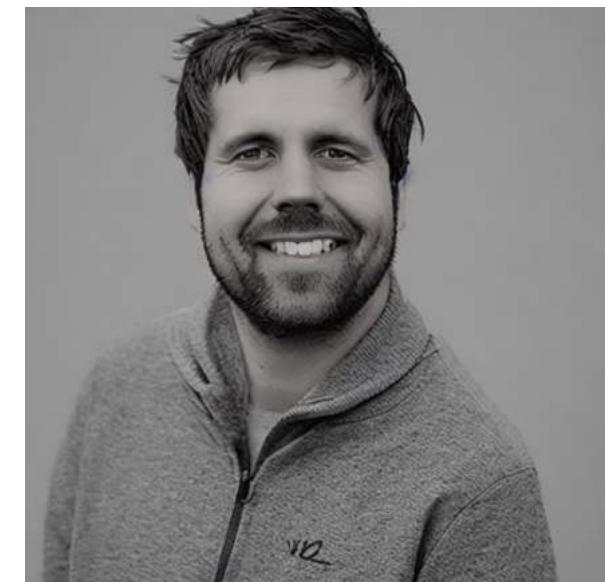
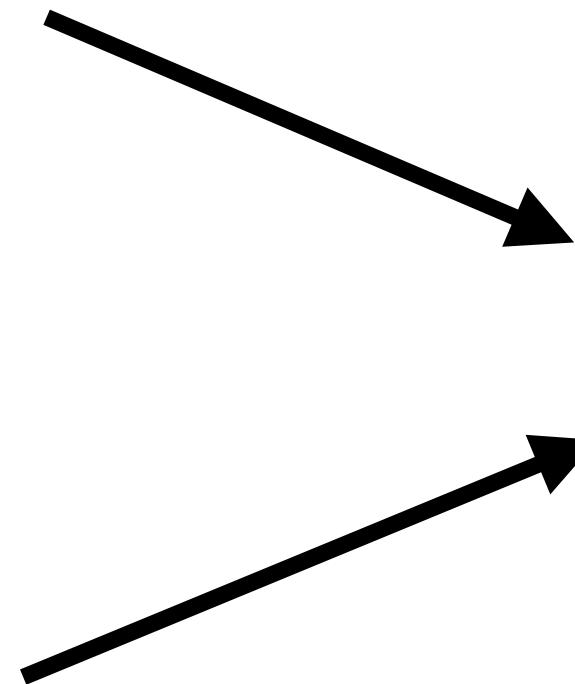


**Flow-based generative models:** minimize the negative log-likelihood



(+ score based diffusion models)





[https://huggingface.co/spaces/  
lambdalabs/image-mixer-demo](https://huggingface.co/spaces/lambdalabs/image-mixer-demo)

## “standard” bayesian inference

goal of inference

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

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

↑  
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)**

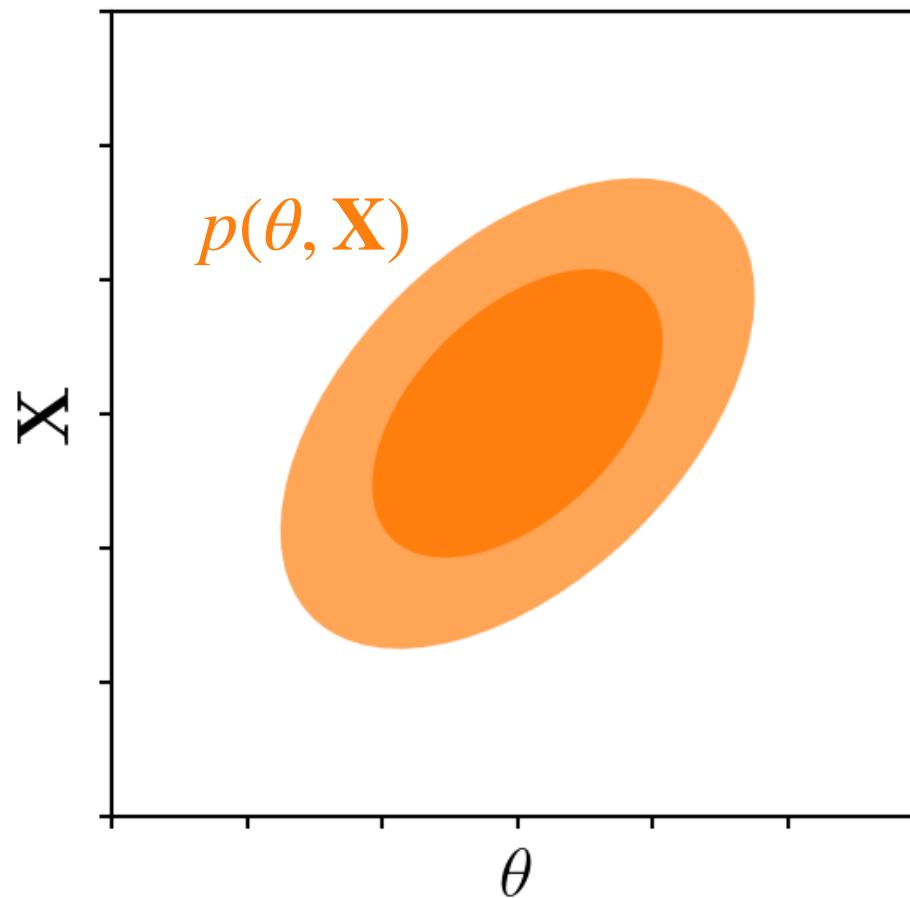
Gaussian likelihood is often **an incorrect assumption (not enough data)**

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

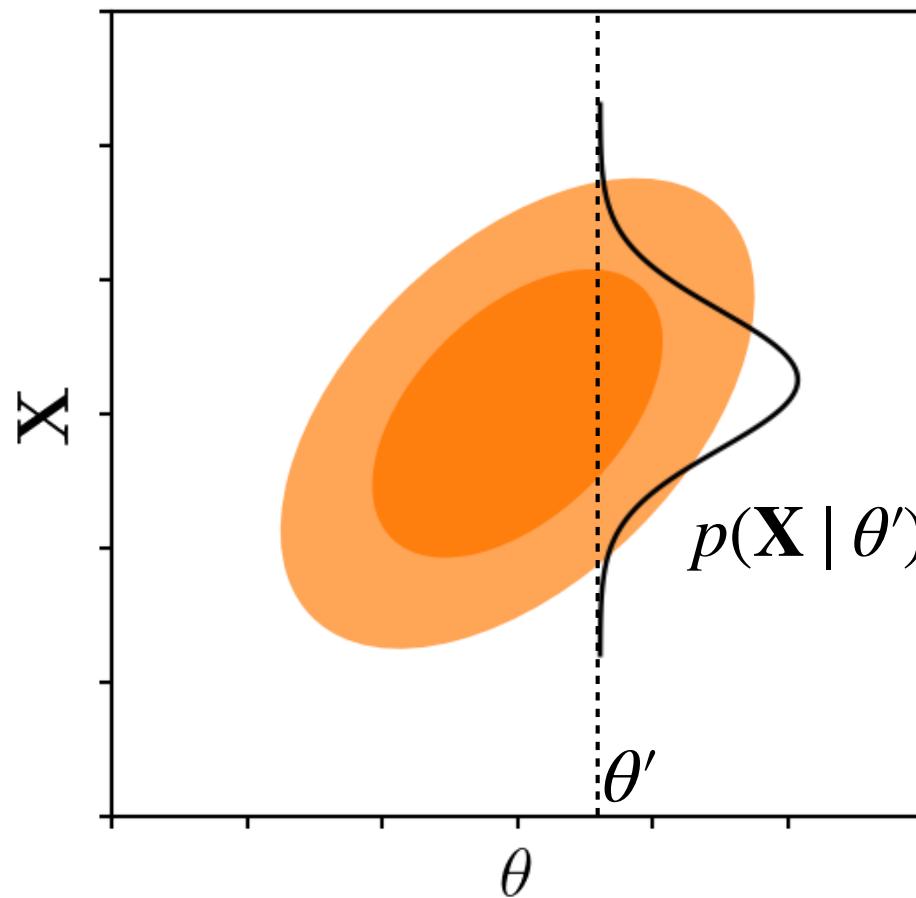
Sometimes the likelihood is intractable. e.g. Baryonic properties of DM haloes?

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

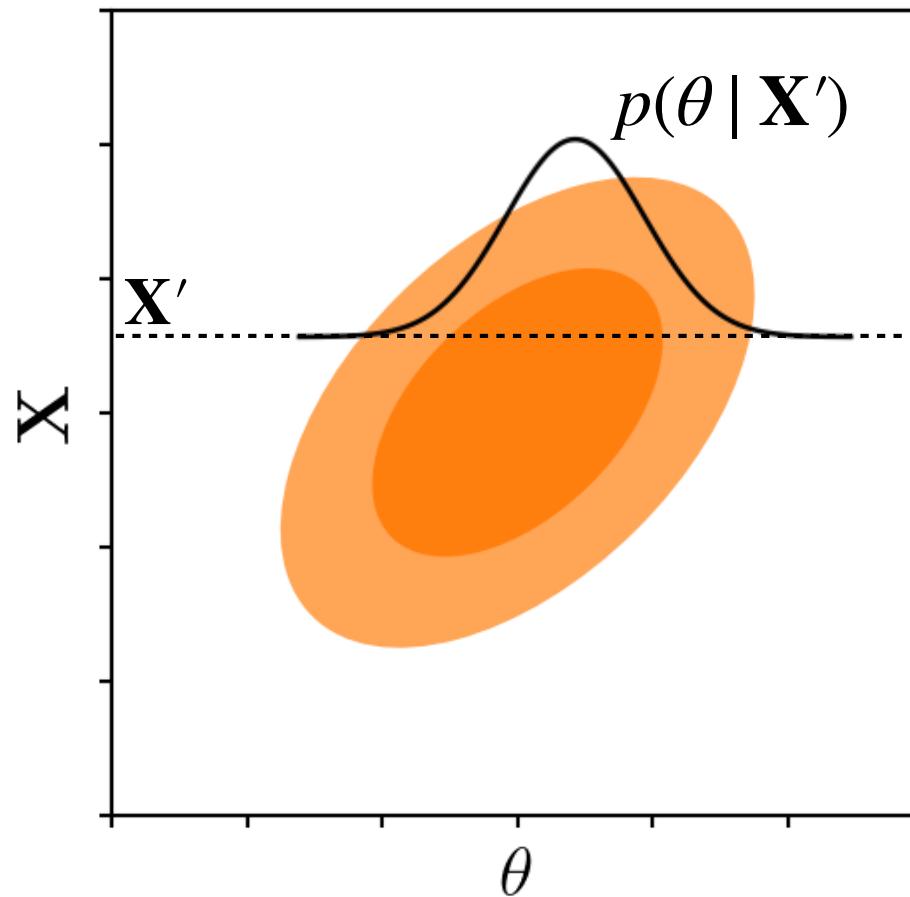
consider an experiment with *1D observation  $\mathbf{X}$*  and *1D parameter  $\theta$*



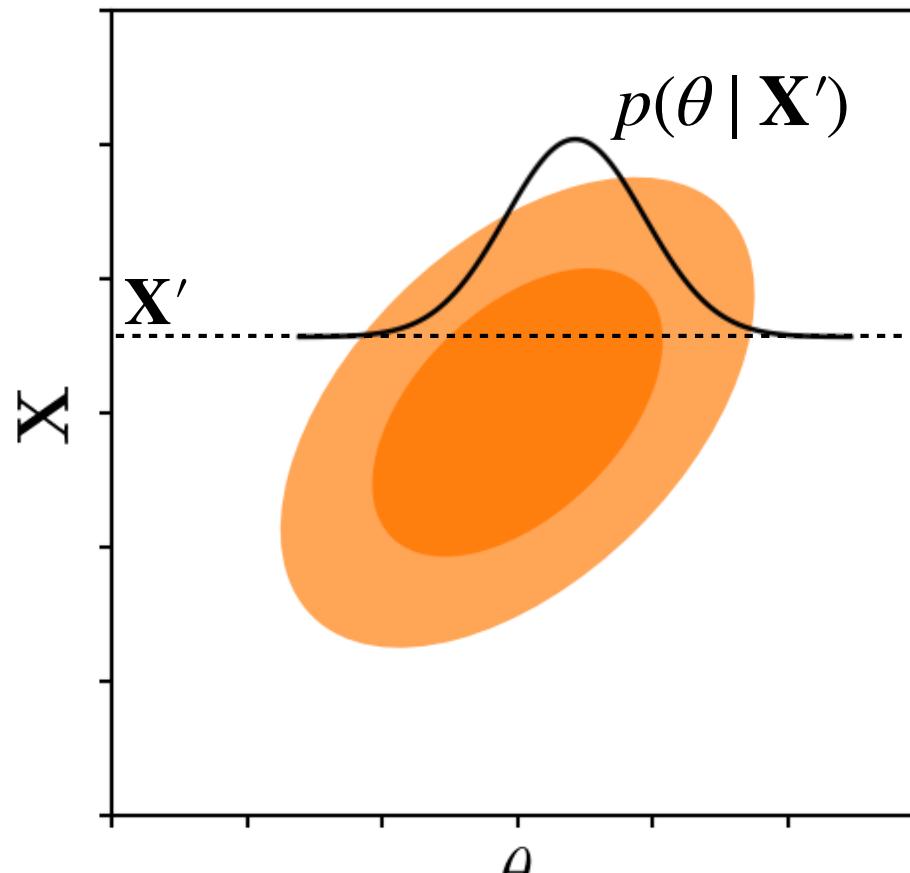
$p(\theta, \mathbf{X})$  includes both the **likelihood**  $p(\mathbf{X} | \theta')$  and posterior  $p(\theta | \mathbf{X}')$



$p(\theta, \mathbf{X})$  includes both the likelihood  $p(\mathbf{X} | \theta')$  and **posterior**  $p(\theta | \mathbf{X}')$



$p(\theta, \mathbf{X})$  includes both the likelihood  $p(\mathbf{X} | \theta')$  and posterior  $p(\theta | \mathbf{X}')$

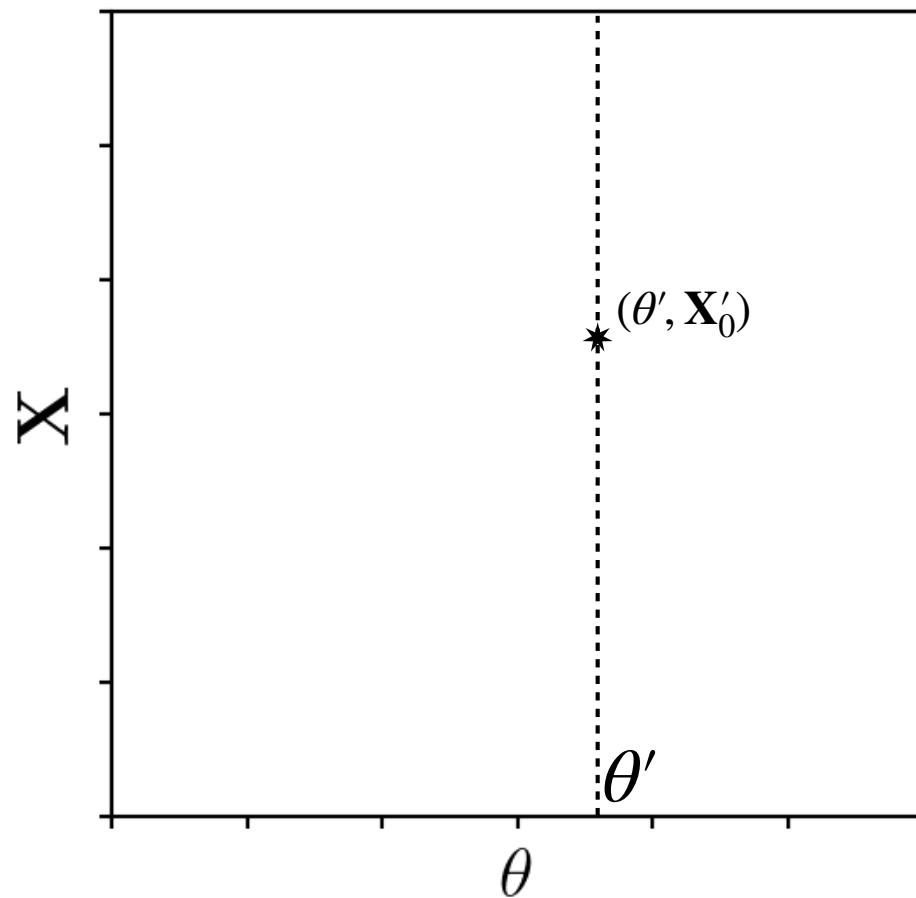


$$p(\theta | X) = \frac{p(X | \theta) p(\theta)}{p(X)}$$

$p(\theta, X) = p(X | \theta) p(\theta) = p(\theta | X) p(X) \longrightarrow$  Bayes rule

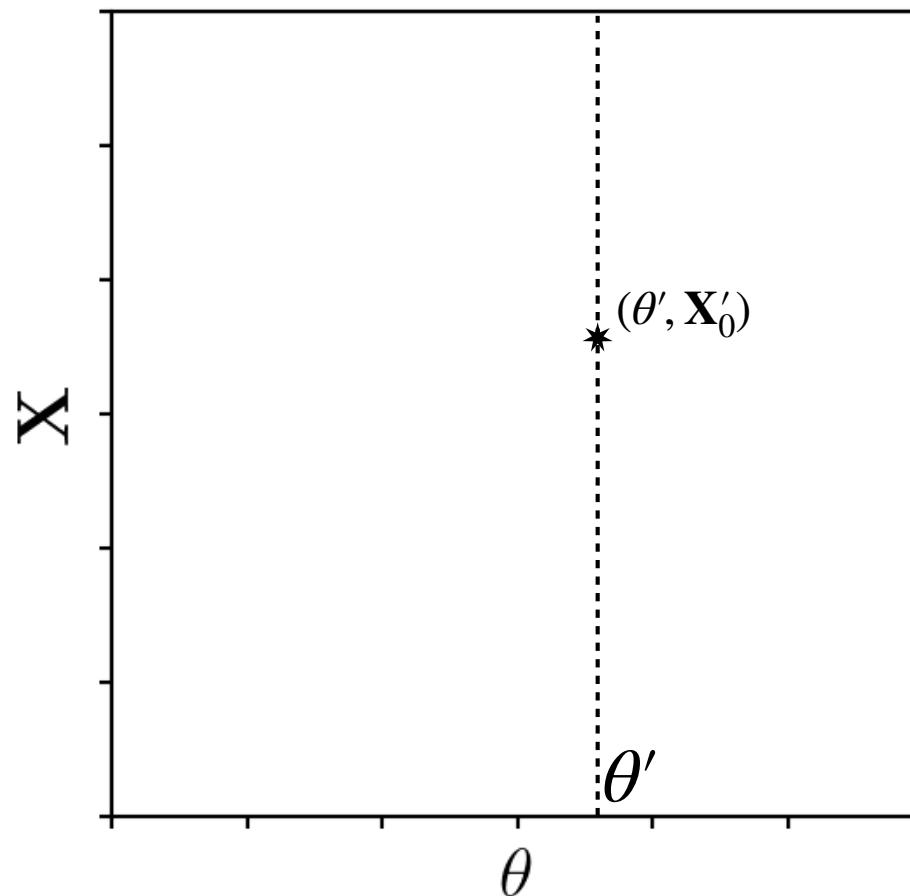
forward model  $F$  that can simulate observations:

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



forward model  $F$  that can simulate observations:

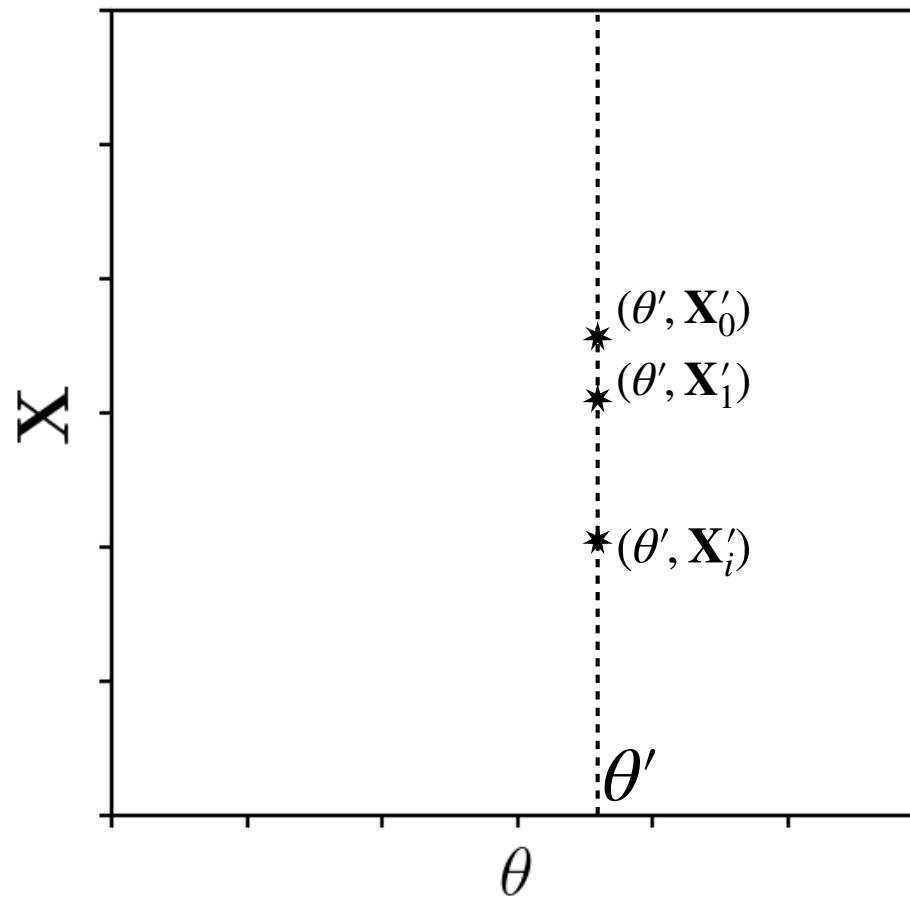
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



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

forward model  $F$  that can simulate observations:

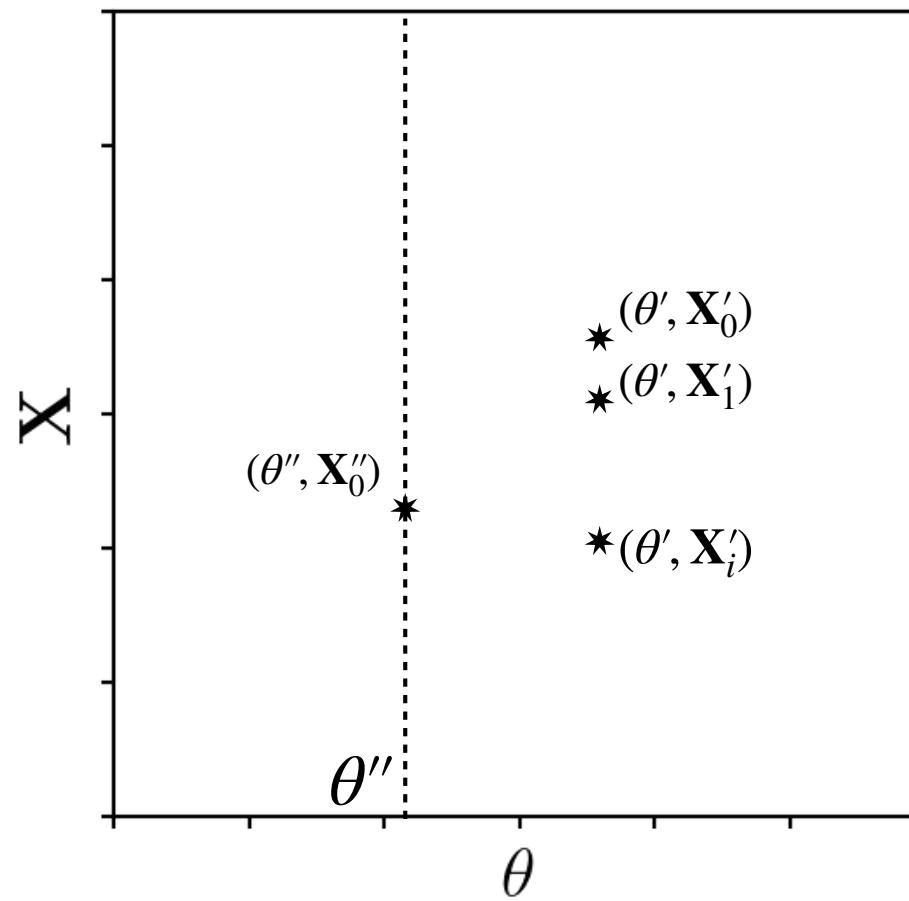
$$\mathbf{X}'_0, \mathbf{X}'_1 \dots \sim F(\theta')$$



$F$  must include noise model

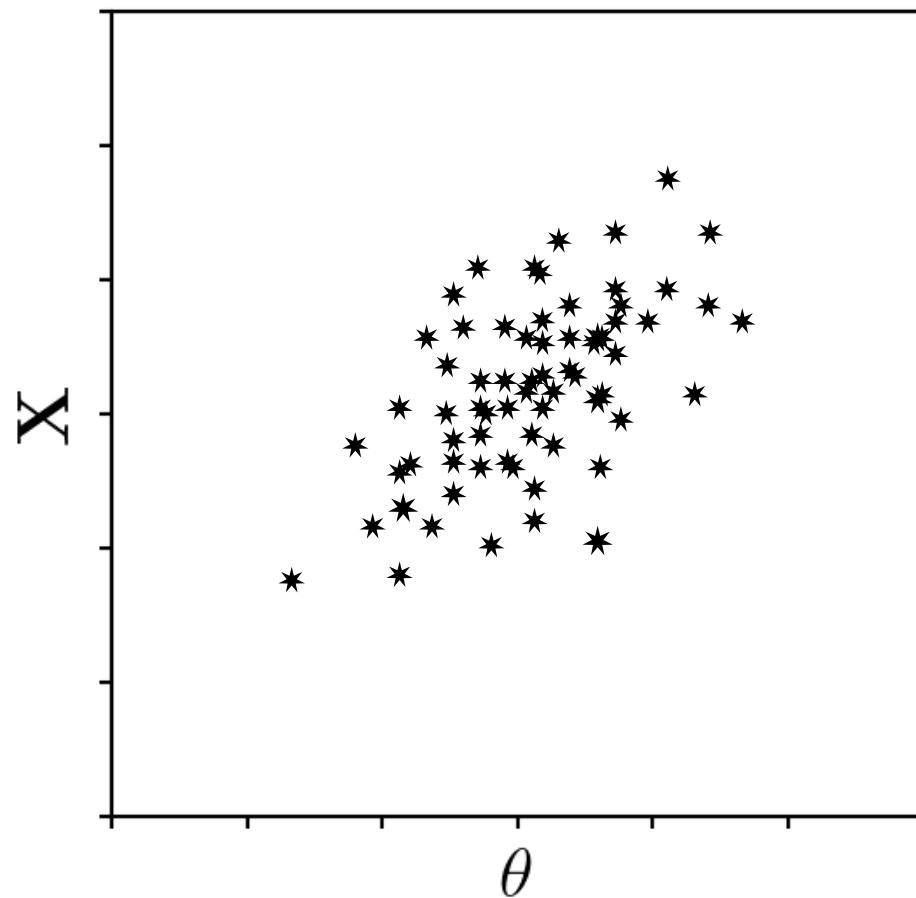
forward model  $F$  that can simulate observations:

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



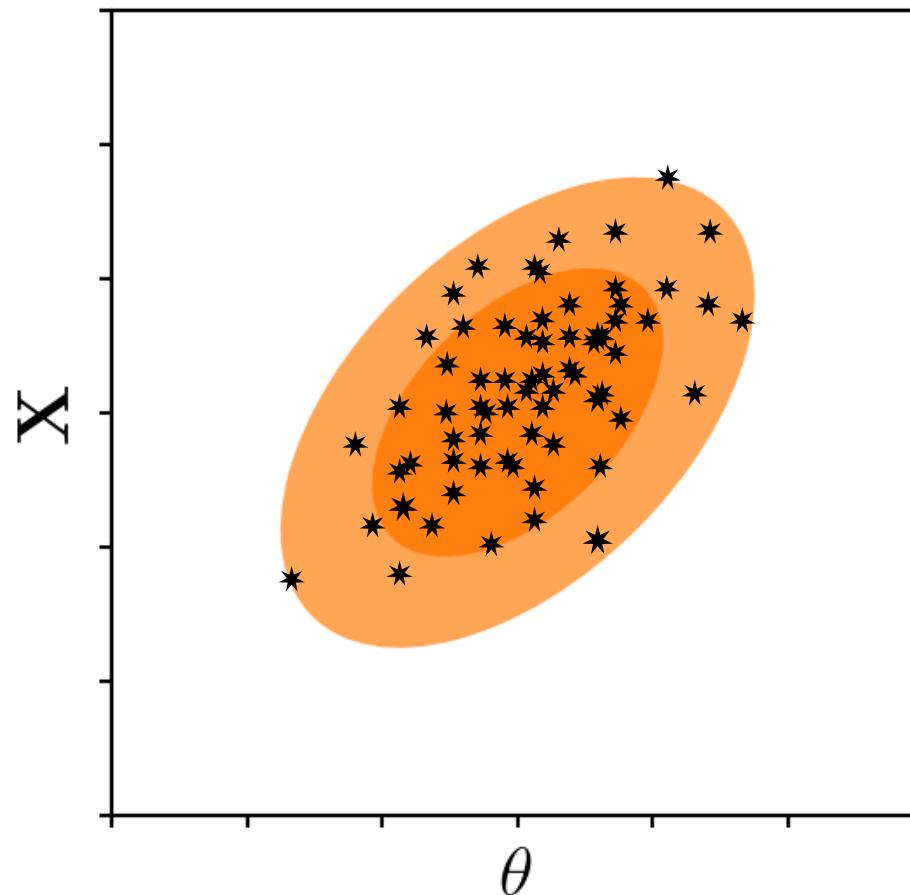
forward model  $F$  that can simulate observations:

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



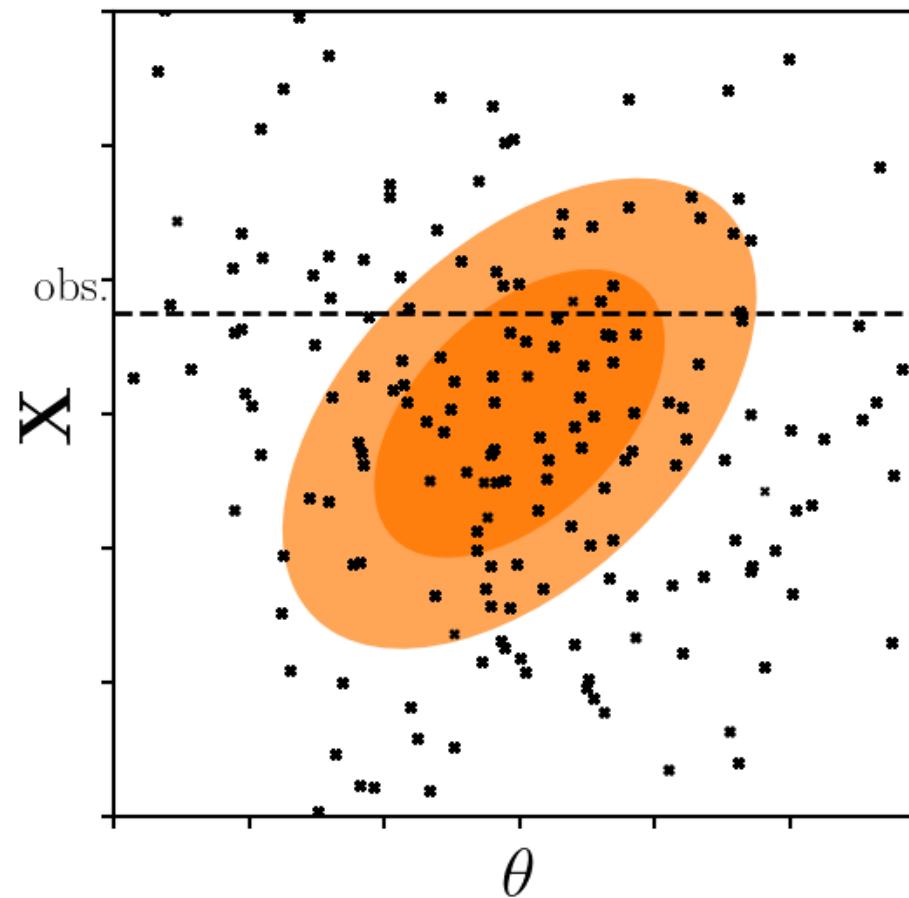
forward model  $F$  that can simulate observations:

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

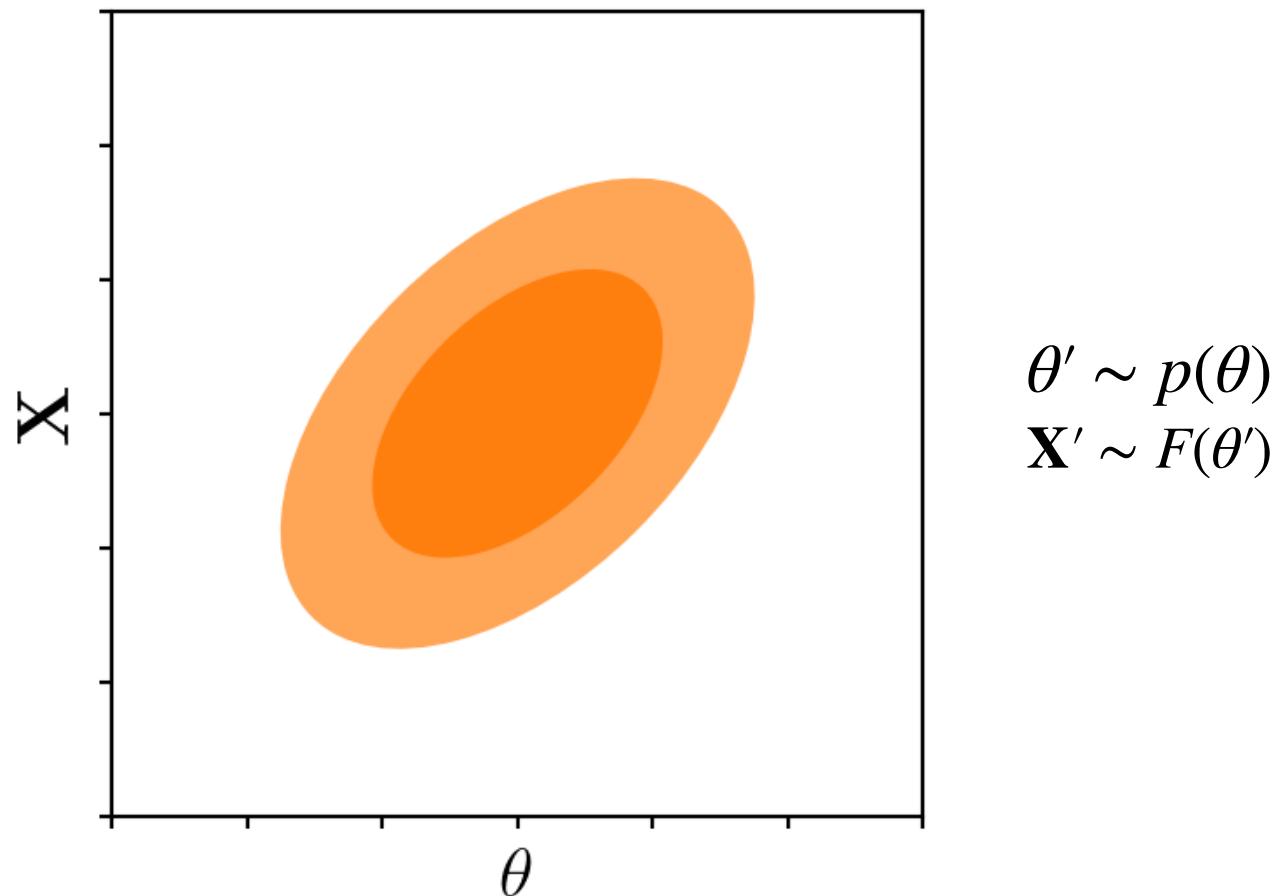


$F$  allows us to sample  $(\theta', X') \sim p(\theta, X)$  ← Unknown

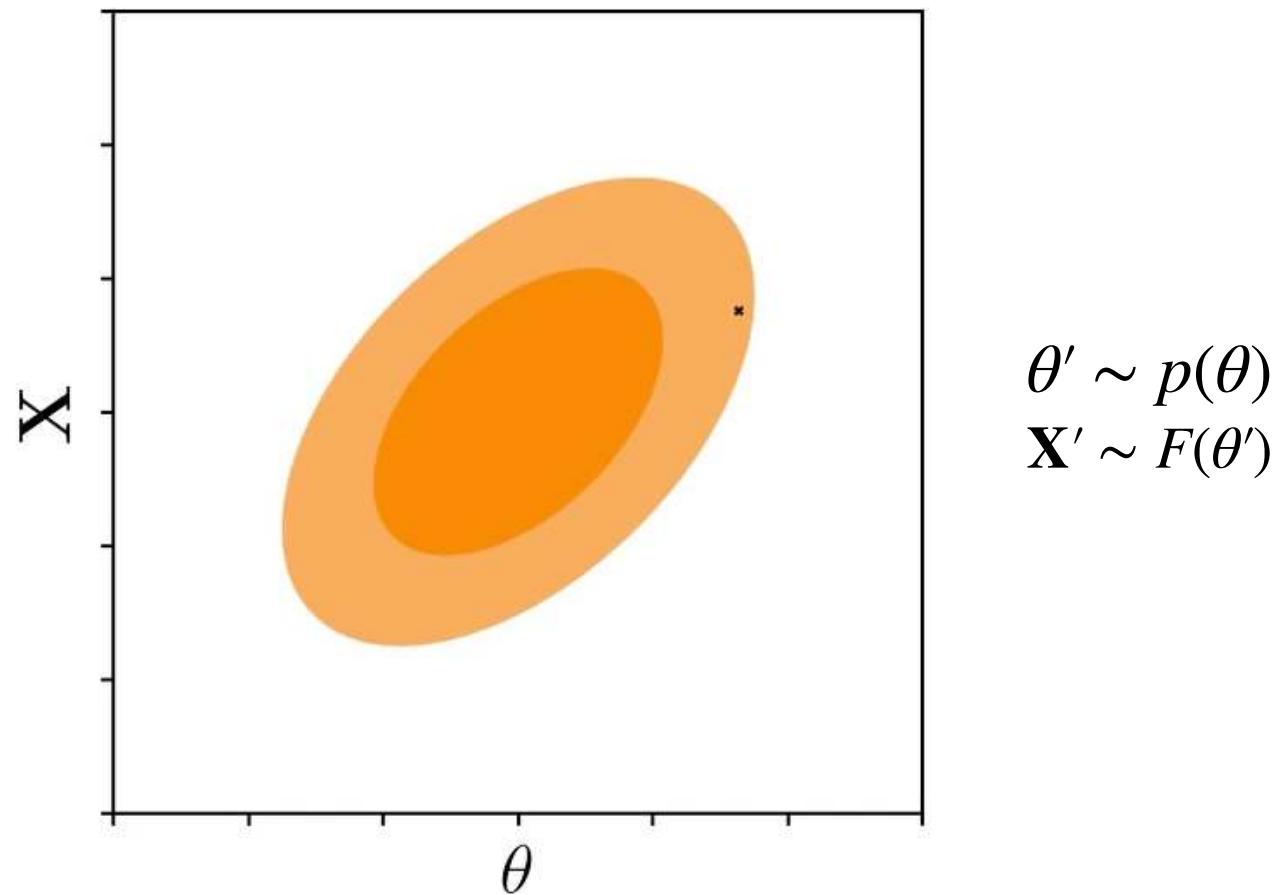
with  $F$  one way to infer the posterior is through *brute force* —  
**approximate bayesian computation (ABC)**



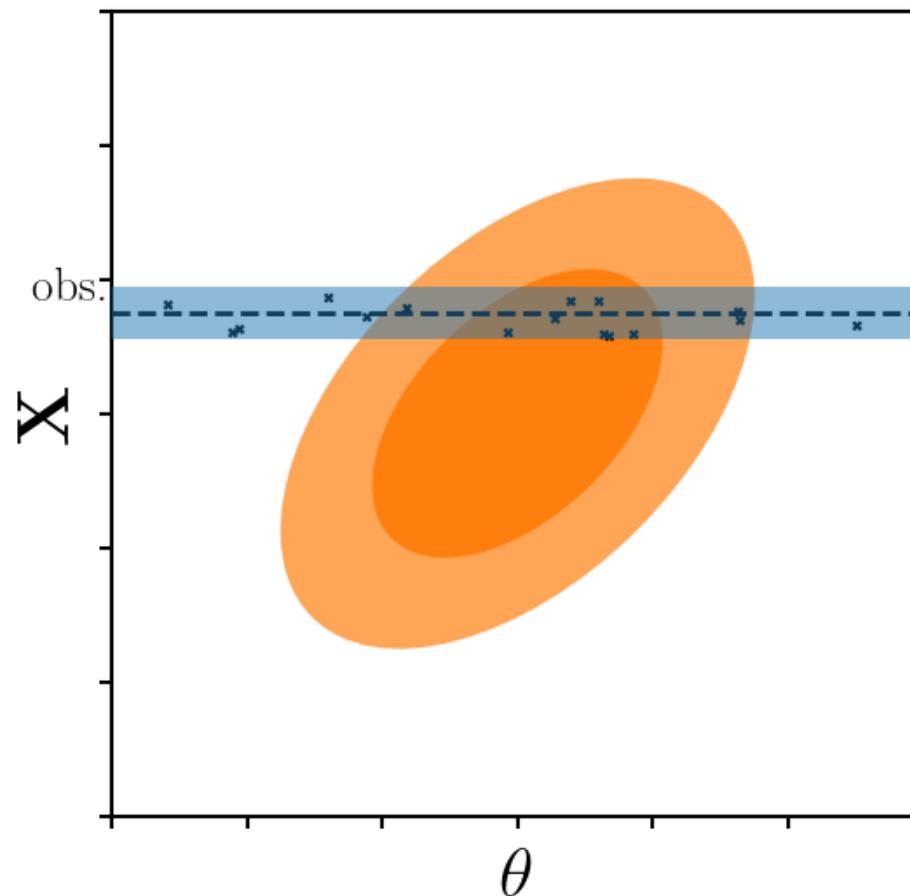
with  $F$  one way to infer the posterior is through *brute force* —  
**approximate bayesian computation (ABC)**



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

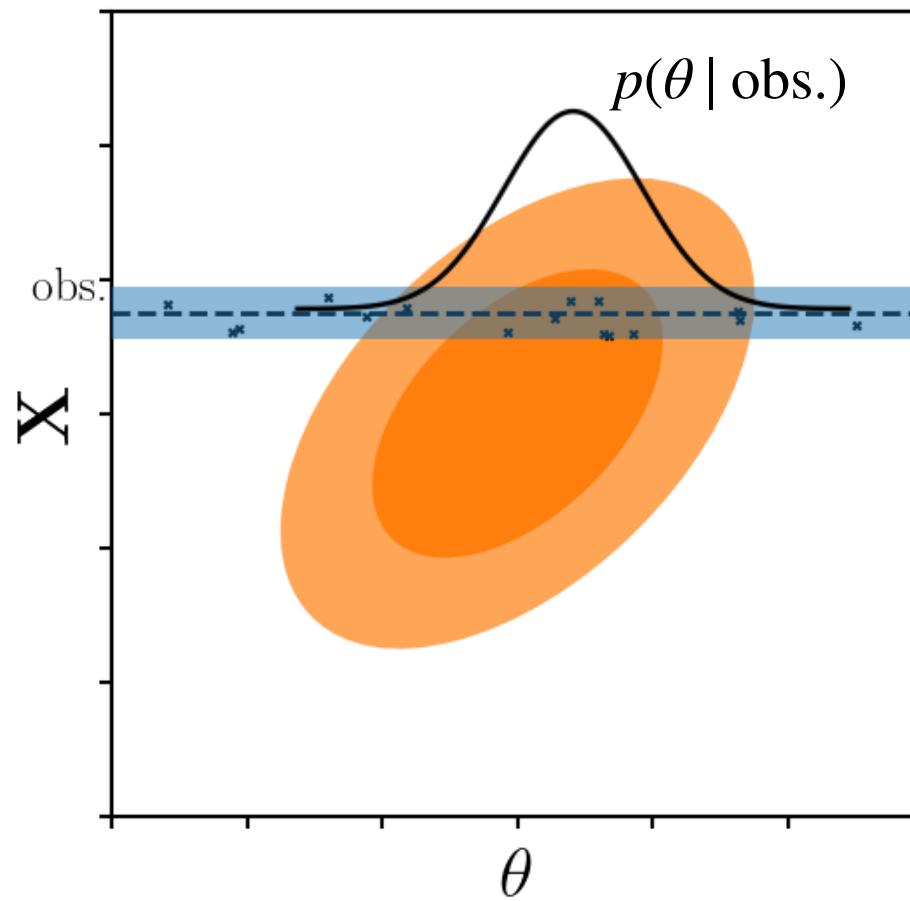


with  $F$  one way to infer the posterior is through *brute force* —  
**approximate bayesian computation (ABC)**



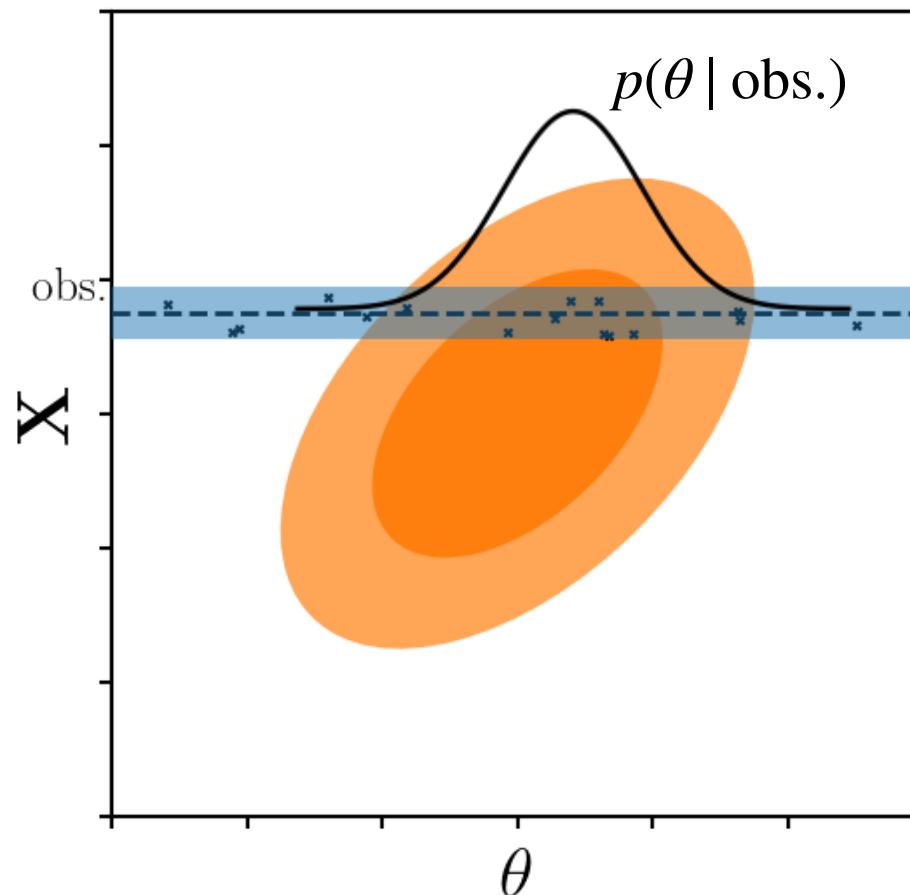
only keep the simulations “close” to the observation

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



only keep the simulations “close” to the observation

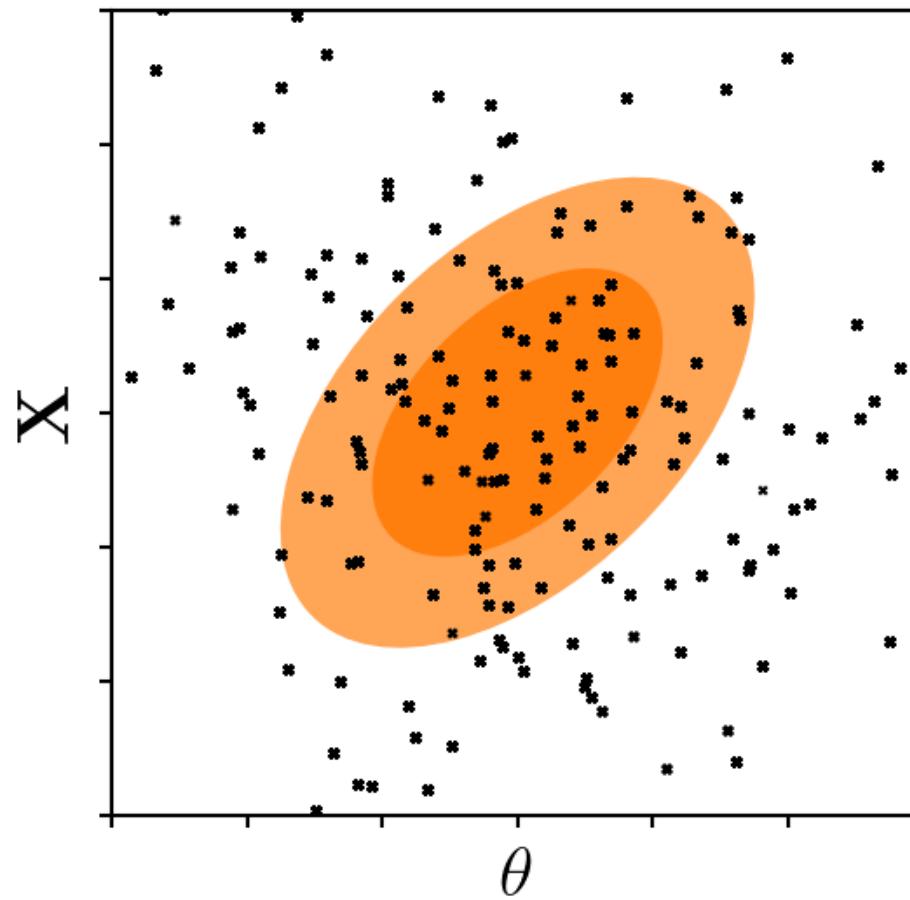
with  $F$  one way to infer the posterior is through *brute force* —  
**approximate bayesian computation** (ABC)



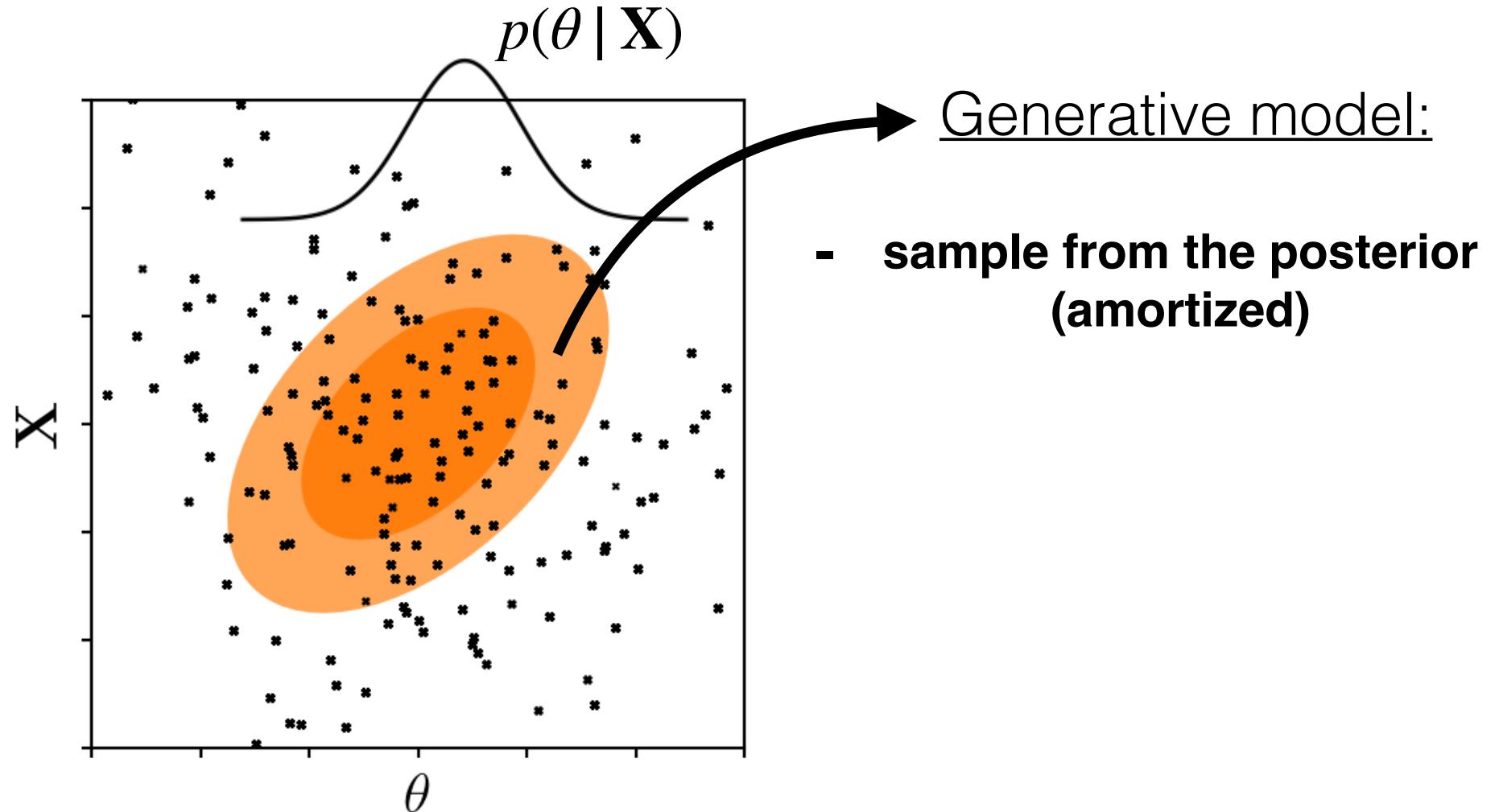
only keep the simulations “close” to the observation

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

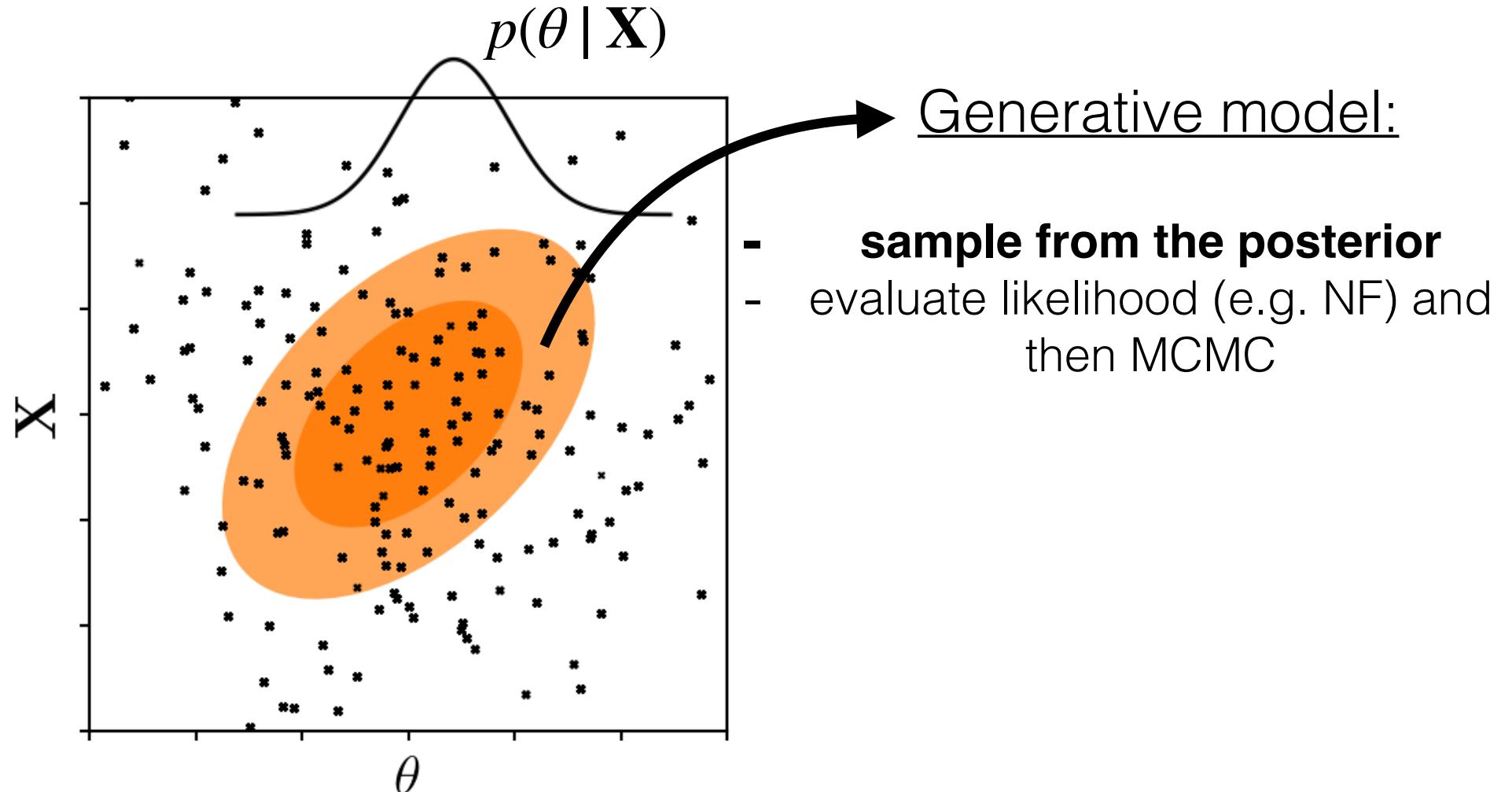
simulation-based inference is a **density estimation problem!**



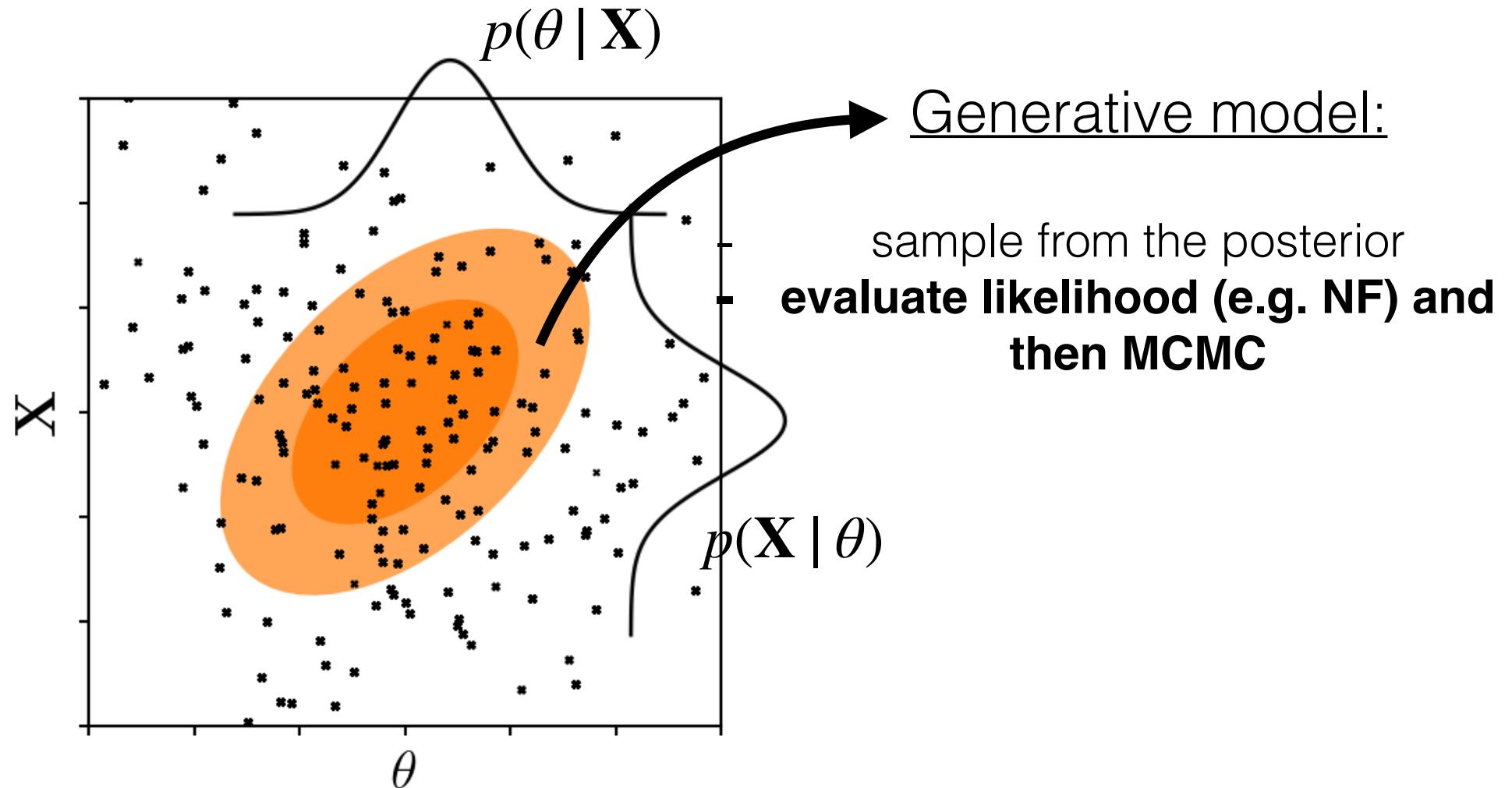
simulation-based inference is a **density estimation problem!**



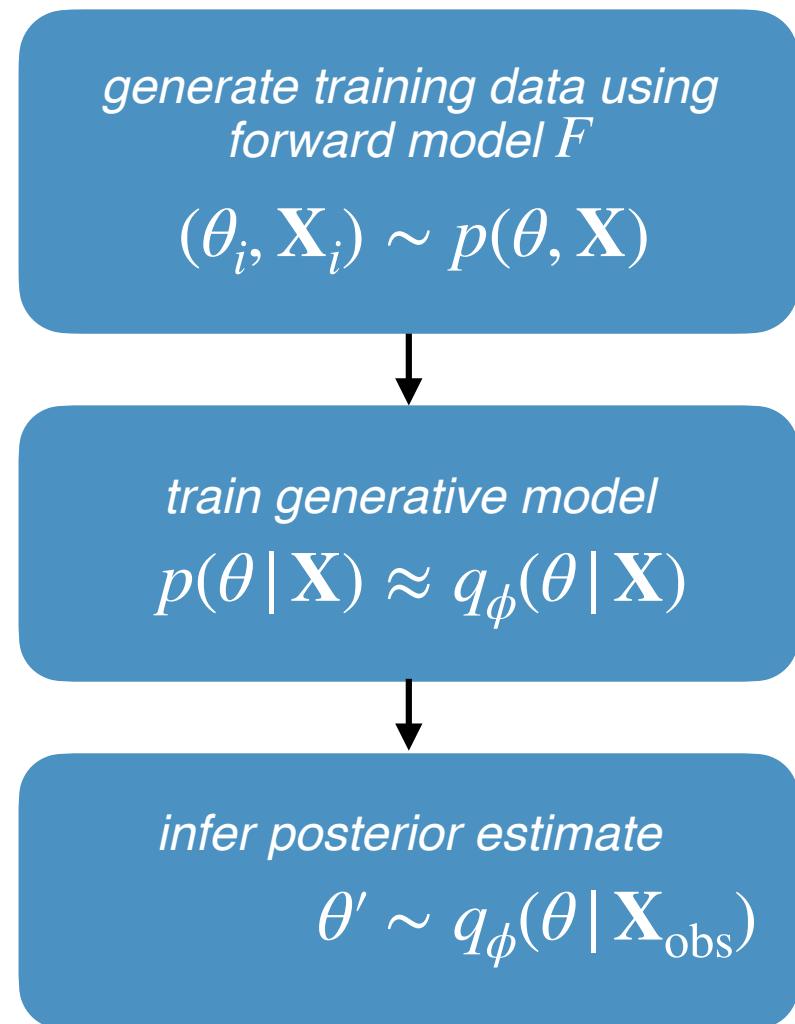
simulation-based inference is a **density estimation problem!**



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



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



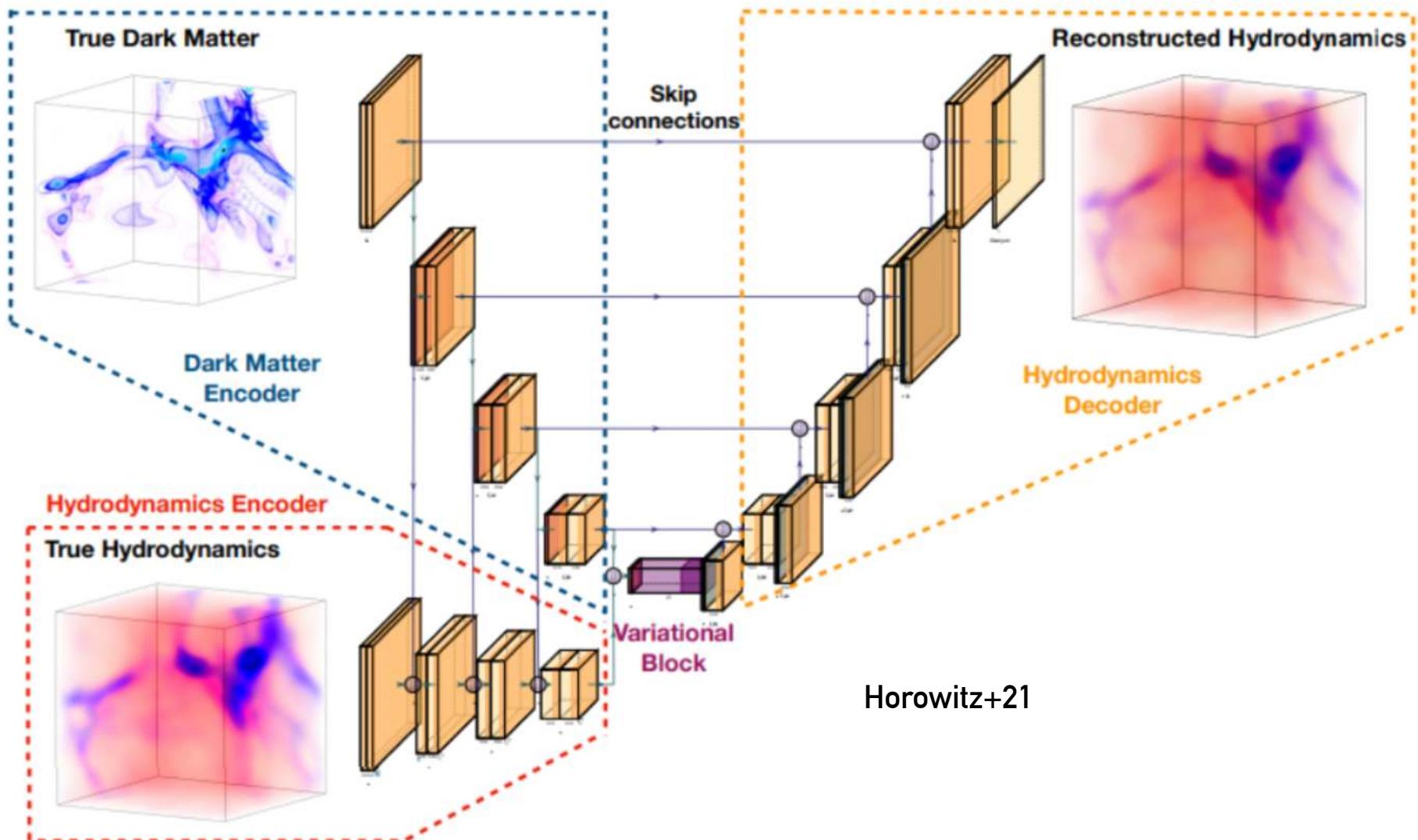
no guarantee it converges  
to the true posterior



*infer posterior estimate*

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

# Painting Baryons





simulation-based inference for Bayesian **galaxy SED modeling**

galaxy SED as a **composite stellar population**: star formation history, metallicity history, stellar population synthesis, dust

$$f_{\text{CSP}}(\lambda, t) = \int_{t'=0}^{t'=t} \text{SFR}(t') f_{\text{SSP}}(t', Z(t')) e^{-\tau_{\text{dust}}(t')} dt'$$

many galaxy SED models available: e.g. ***PROVABGS***

$\log M_*$

$\tau_{\text{BC}}$

*birth-cloud dust optical depth*

$\tau_{\text{ISM}}$

*diffuse cloud dust optical depth*

$\delta$

*index of Calzetti+(2000) attenuation curve*

$\beta_1^{\text{SFH}}, \beta_2^{\text{SFH}}, \beta_3^{\text{SFH}}, \beta_4^{\text{SFH}}, f_{\text{burst}}, t_{\text{burst}}$

$\gamma_1^{\text{ZH}}, \gamma_2^{\text{ZH}}$



many galaxy SED models available: e.g. ***PROVABGS***

$\log M_*$

$\tau_{\text{BC}}$

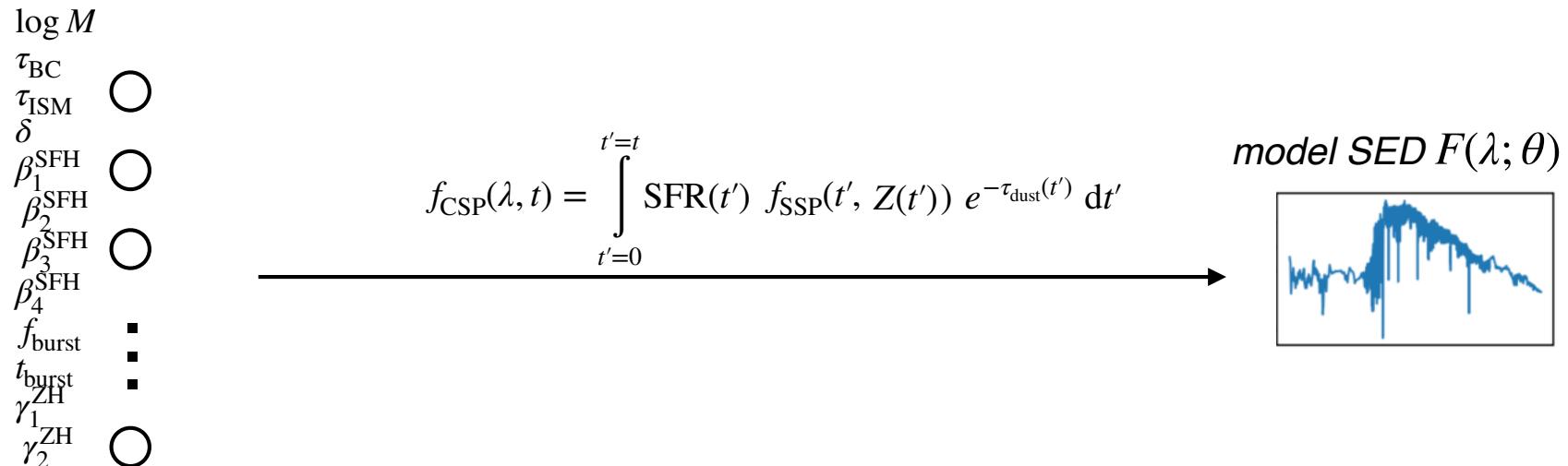
$\tau_{\text{ISM}}$

$\delta$

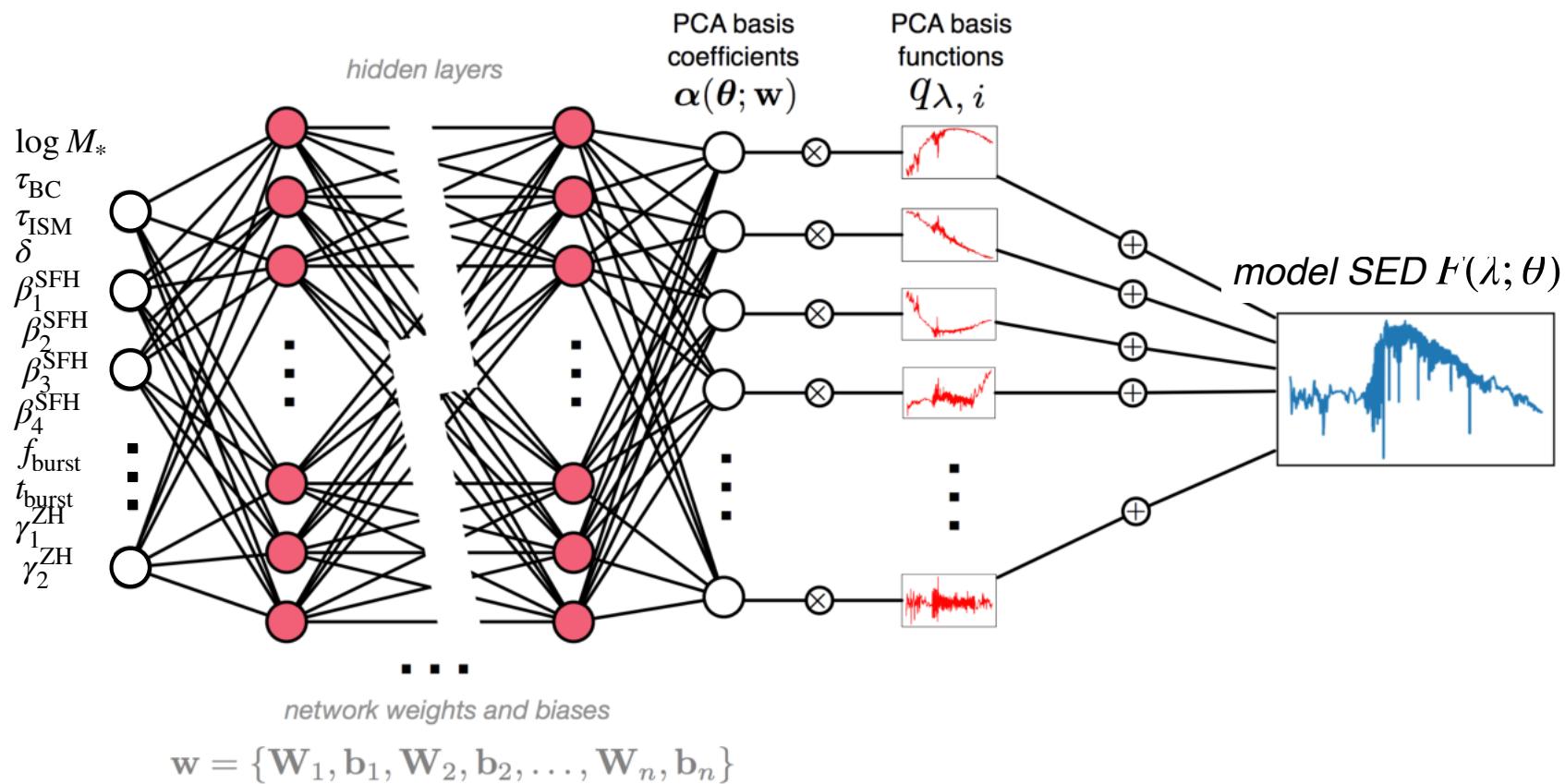
$\beta_1^{\text{SFH}}, \beta_2^{\text{SFH}}, \beta_3^{\text{SFH}}, \beta_4^{\text{SFH}}, f_{\text{burst}}, t_{\text{burst}}$  *star formation history*  
 $\gamma_1^{\text{ZH}}, \gamma_2^{\text{ZH}}$  *metallicity history*



many galaxy SED models available: e.g. *PROVABGS* comes with  
**neural emulator**



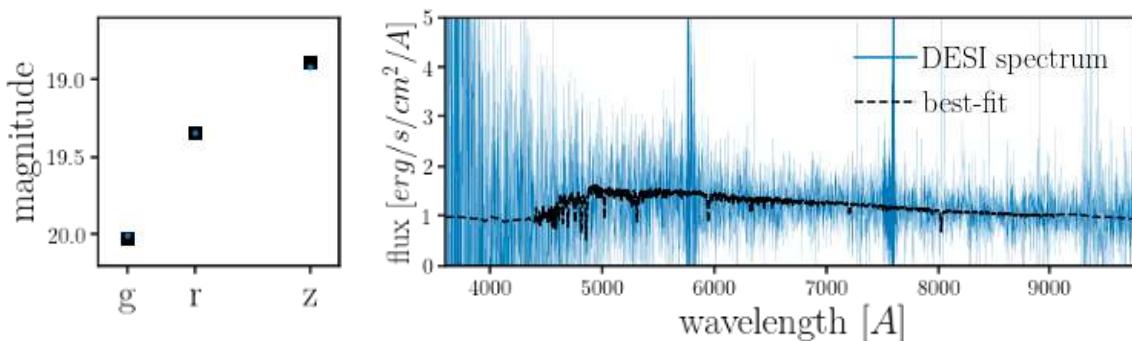
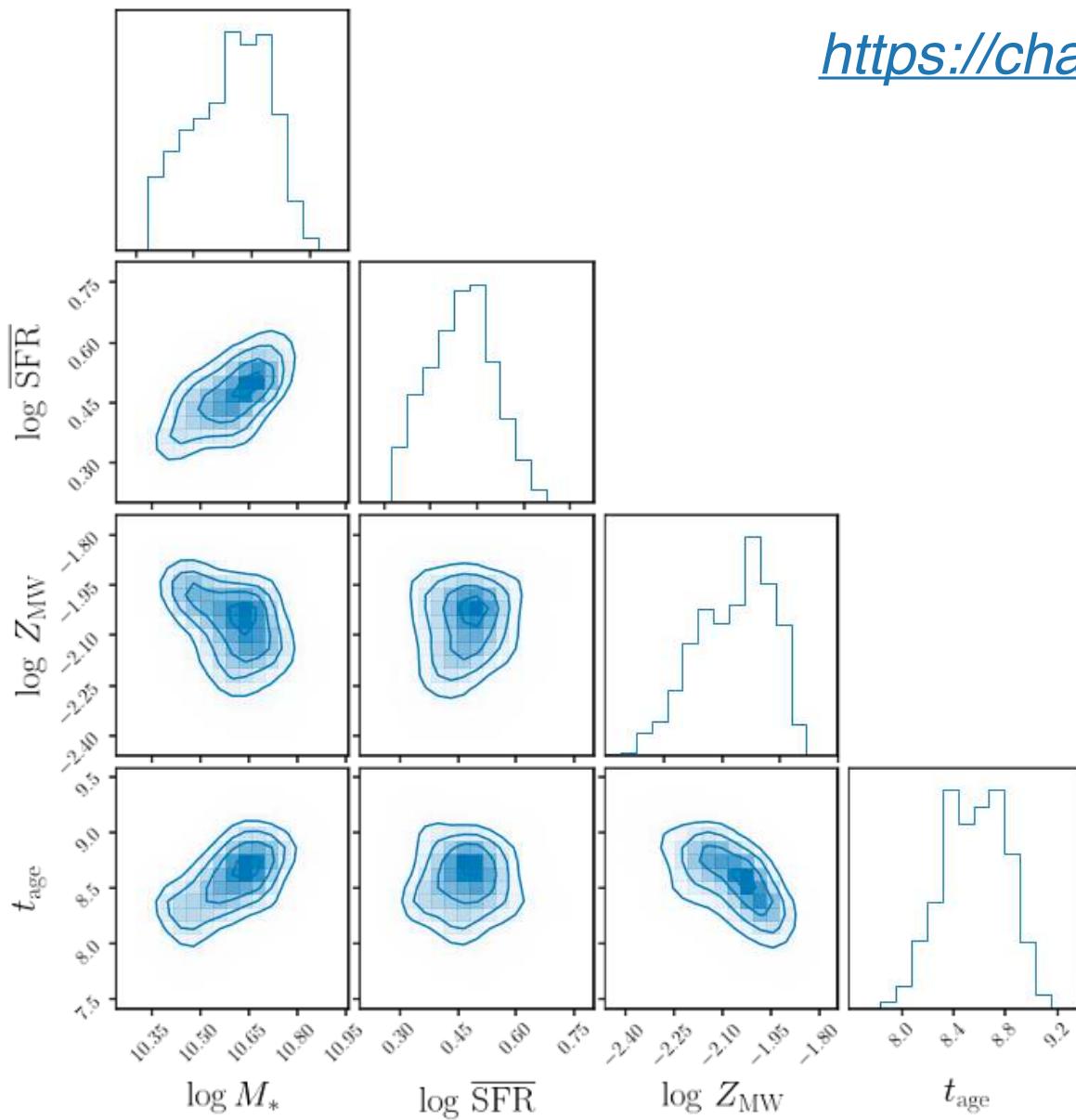
many galaxy SED models available: e.g. *PROVABGS* comes with  
**neural emulator** — 1000x faster



James Kyubin Kwon  
UC Santa Barbara

Alsing+Hahn+(2019), Hahn+(2022b), Kwon,Hahn,Alsing(2022)



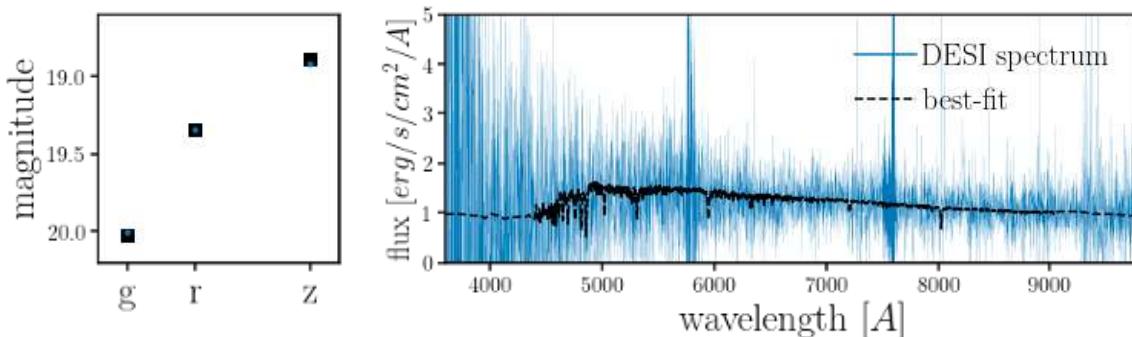
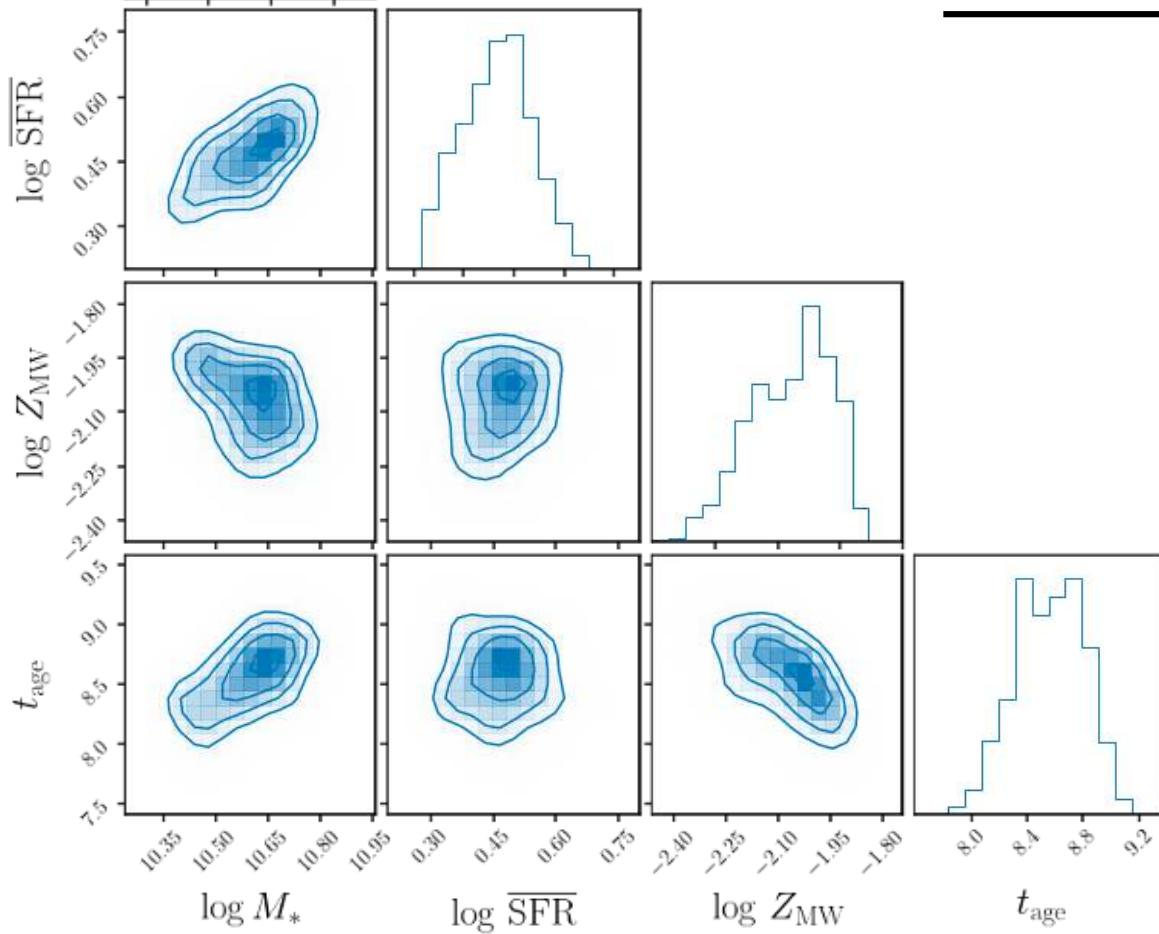


<https://changhoonhahn.github.io/provabgs>

$\sim 10 \text{ CPU min/galaxy}$

$\times (\gg 10 \text{ million galaxies})$

$=\gg 2 \text{ million CPU hours}$



PROVA BGS

Hahn+(2022b), Kwon, Hahn, Alsing(2022)

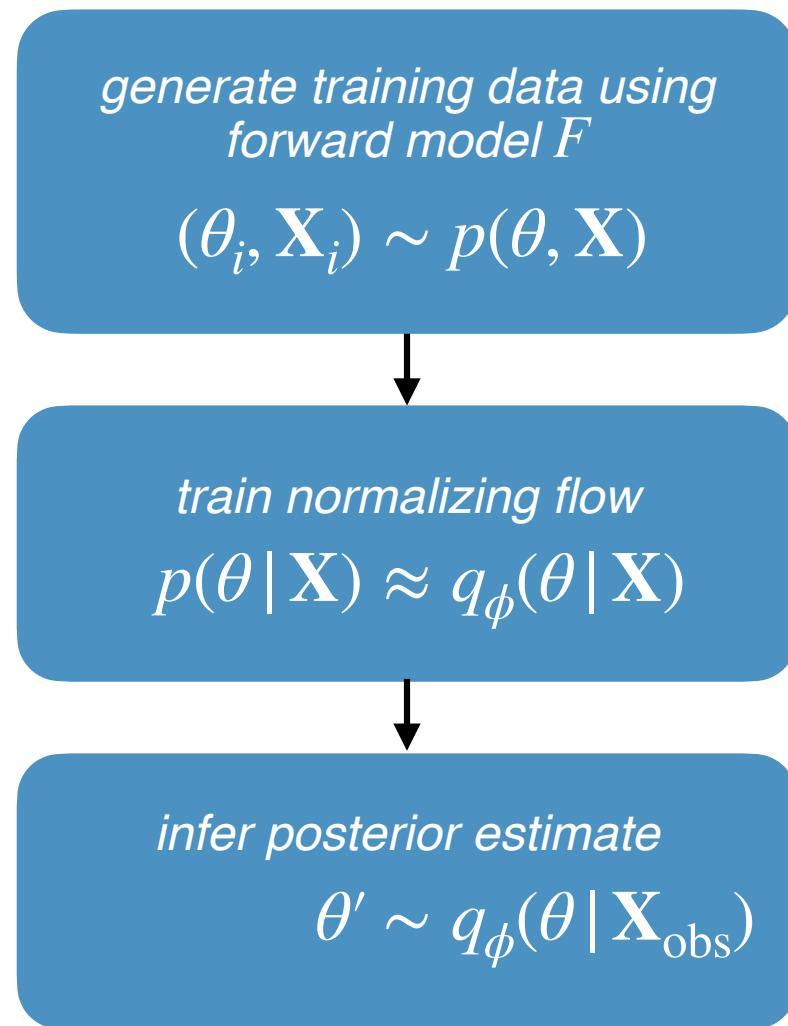
## why **simulation-based inference**?

Gaussian likelihood is often an incorrect assumption

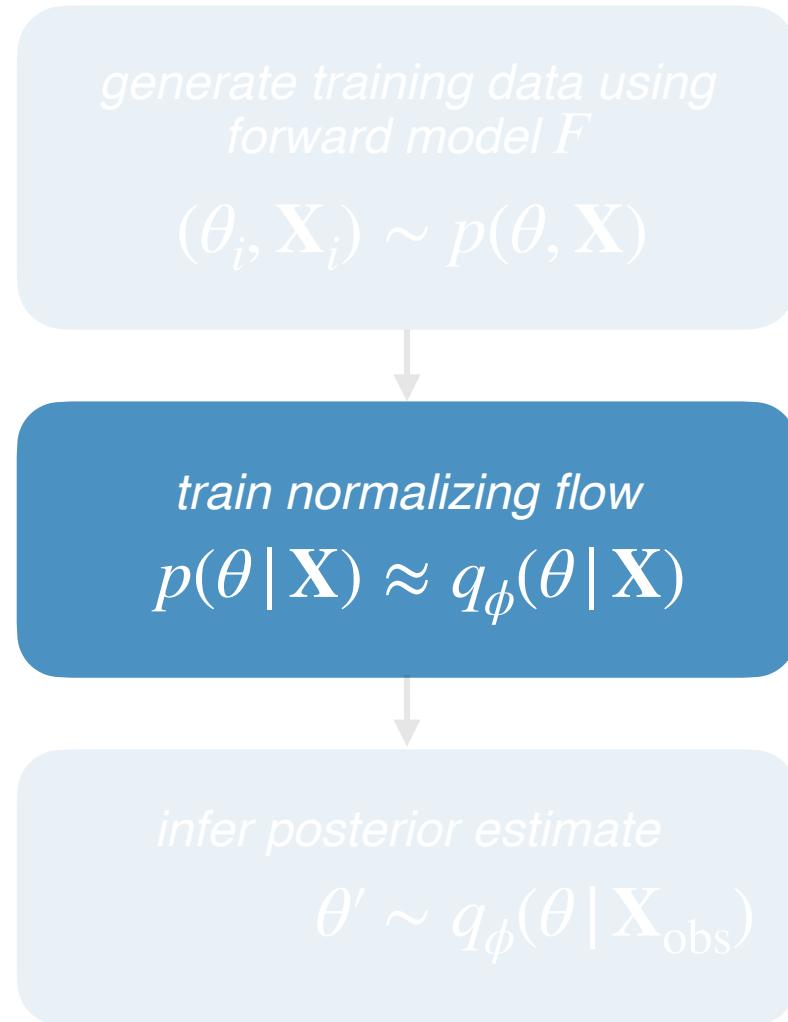
simulations have advantages over analytic models

**amortized** inference to infer *billions* of observations

# simulation-based inference using normalizing flows flowchart

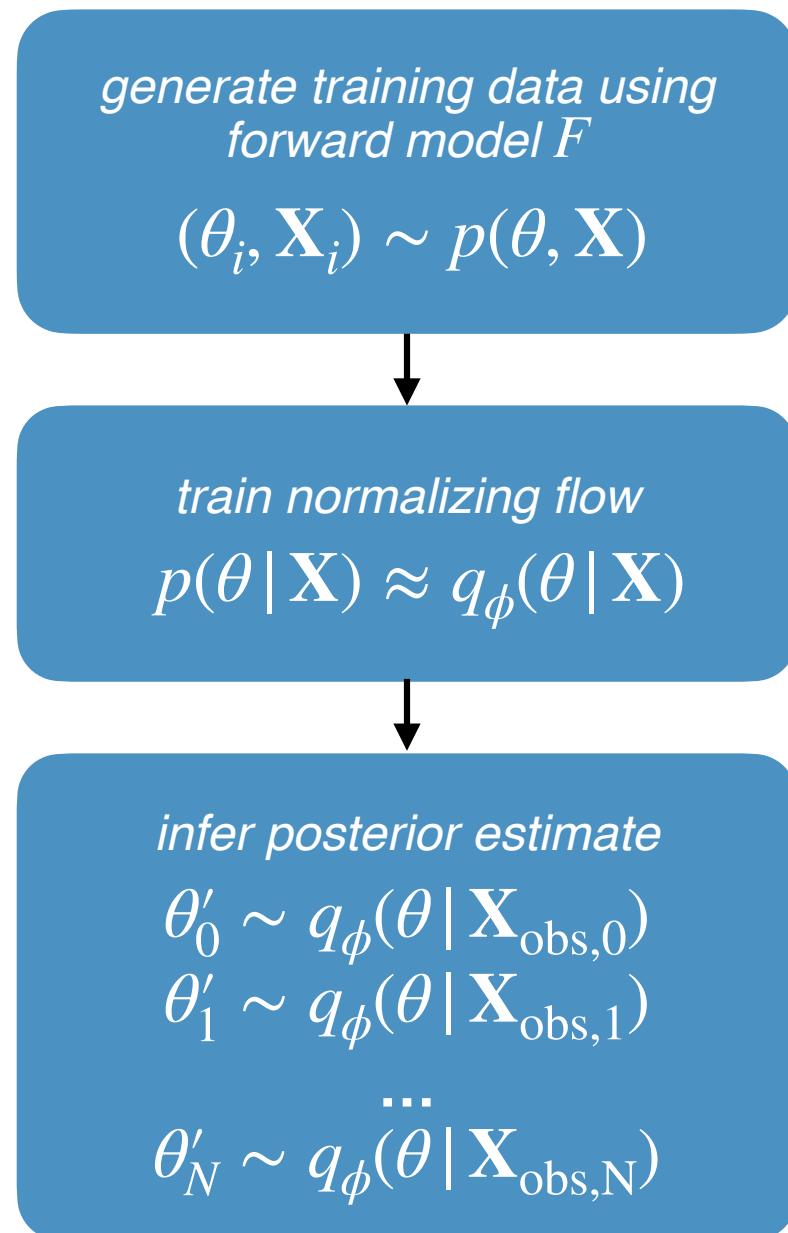


## simulation-based inference using normalizing flows flowchart



we can train  $p(\theta | \mathbf{X}) \approx q_\phi(\theta | \mathbf{X})$  over the **entire  $\mathbf{X}$ -space**

# simulation-based inference using normalizing flows flowchart



*SPS parameters*

$$p(\theta_{\text{SPS}} \mid D, \sigma, z)$$

*SPS parameters*

$$p(\theta_{\text{SPS}} | D, \sigma, z)$$

SDSS photometry

*SPS parameters*

$$p(\theta_{\text{SPS}} | D, \sigma, z)$$

SDSS photometry

measured uncertainties

*SPS parameters*

$$p(\theta_{\text{SPS}} | D, \sigma, z)$$

SDSS photometry

measured uncertainties

measured redshift

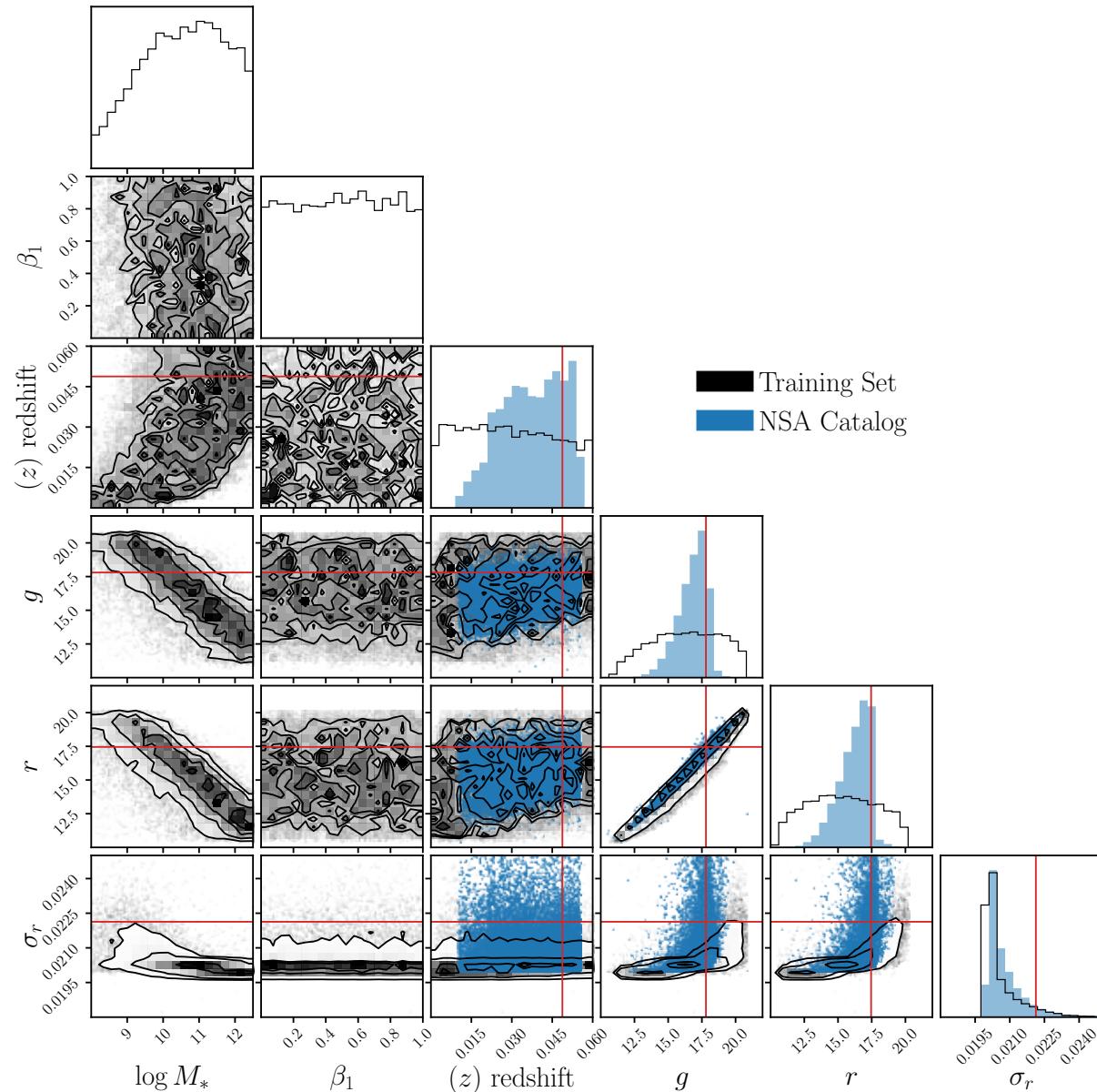
*generate training data using  
forward model  $F$*

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

sample  $\theta' \sim p(\theta_{\text{SPS}}, \sigma, z)$

forward model observations

$$D' = \text{SED model}(\theta') + \text{noise}(\theta')$$



*generate training data using  
forward model  $F$*

$$(\theta_i, \sigma_i, z_i, D_i) \sim p(\theta, \sigma, z, D)$$



*train normalizing flow*

$$p(\theta | D, \sigma, z) \approx q_\phi(\theta | D, \sigma, z)$$



*infer posterior estimates*

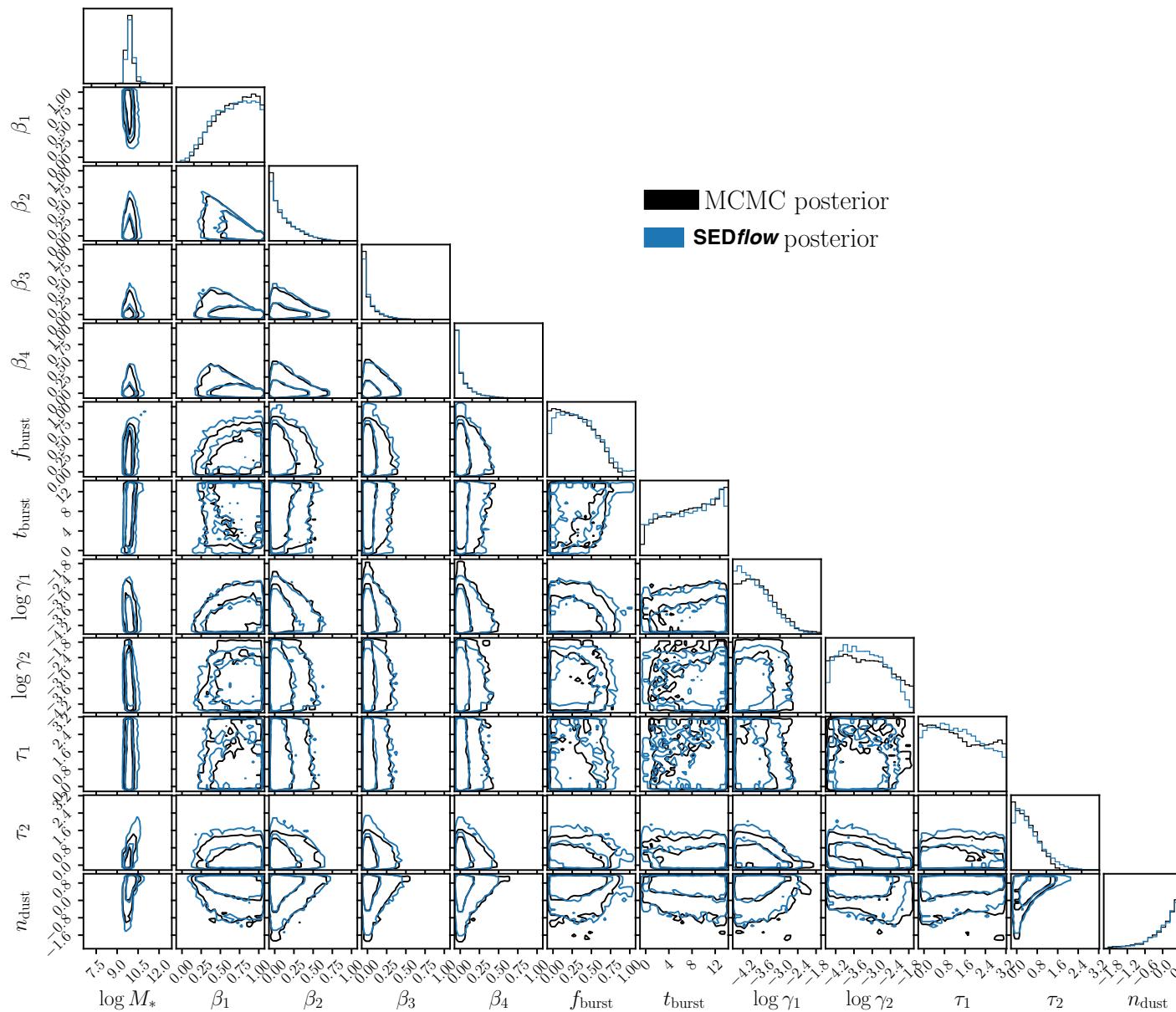
$$\theta'_0 \sim q_\phi(\theta | D_{\text{obs},0}, \sigma_{\text{obs},0}, z_{\text{obs},0})$$

$$\theta'_1 \sim q_\phi(\theta | D_{\text{obs},1}, \sigma_{\text{obs},1}, z_{\text{obs},1})$$

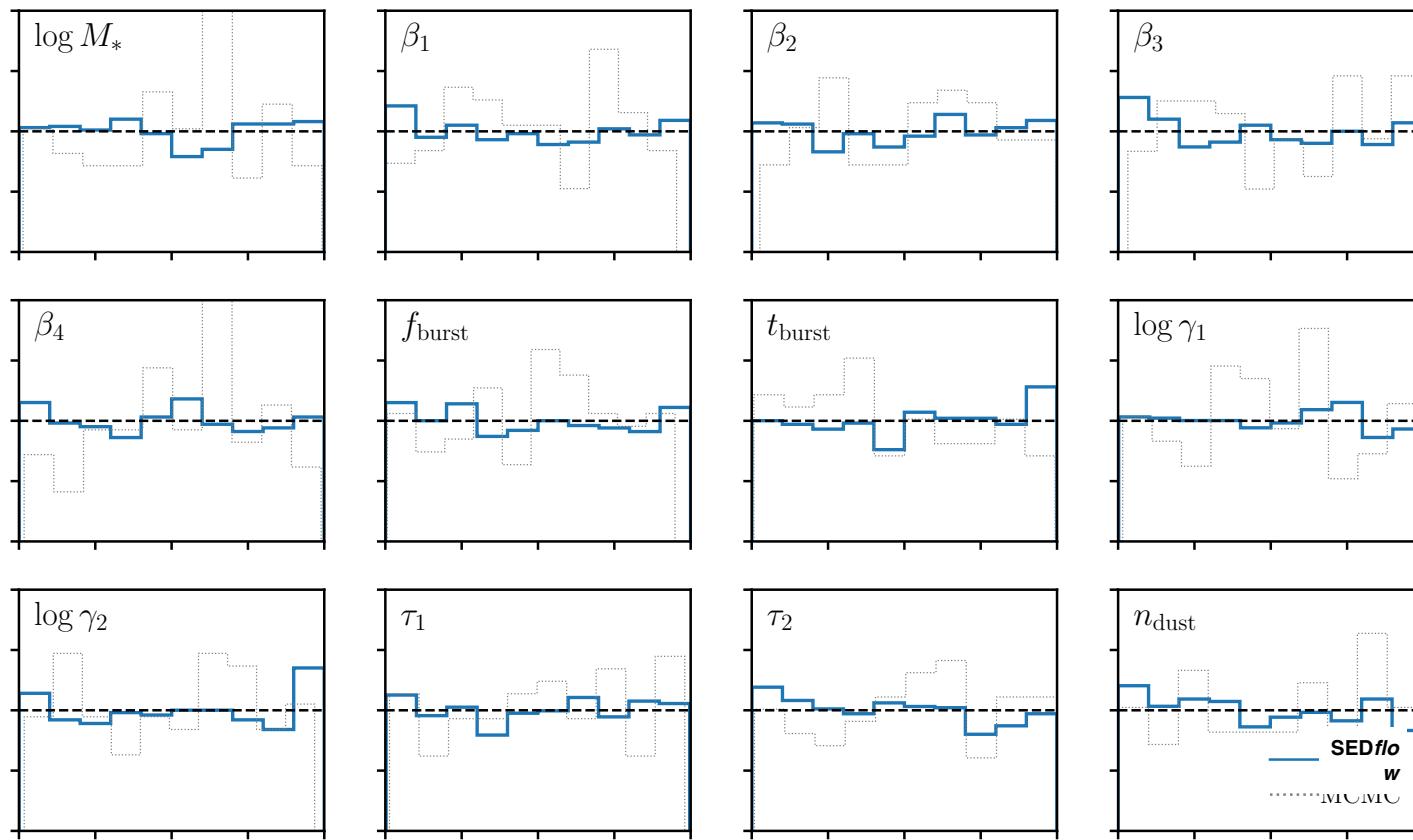
...

$$\theta'_N \sim q_\phi(\theta | D_{\text{obs},N}, \sigma_{\text{obs},N}, z_{\text{obs},N})$$

# **SEDflow** posteriors are *indistinguishable* from MCMC posteriors



# **SEDflow** posteriors accurately estimate the true posteriors — more accurate than MCMC



**SEDflow** posteriors accurately estimate the true posteriors —  
takes ~1 second per galaxy

**SEDflow** posteriors accurately estimate the true posteriors —  
takes ~1 second per galaxy

(~1 second per galaxy)

x (10 million galaxies)

---

<3000 CPU hours

# Thank you!

\* [mhuertas@iac.es](mailto:mhuertas@iac.es)