

# PART V: DIMENSIONALITY REDUCTION, GENERATIVE MODELLING AND DENSITY ESTIMATION

# TYPES OF MACHINE LEARNING

SUPERVISED

LEARNS A MAPPING FROM  
X [FEATURES] TO Y  
[LABELS]

$$P(Y|X)$$

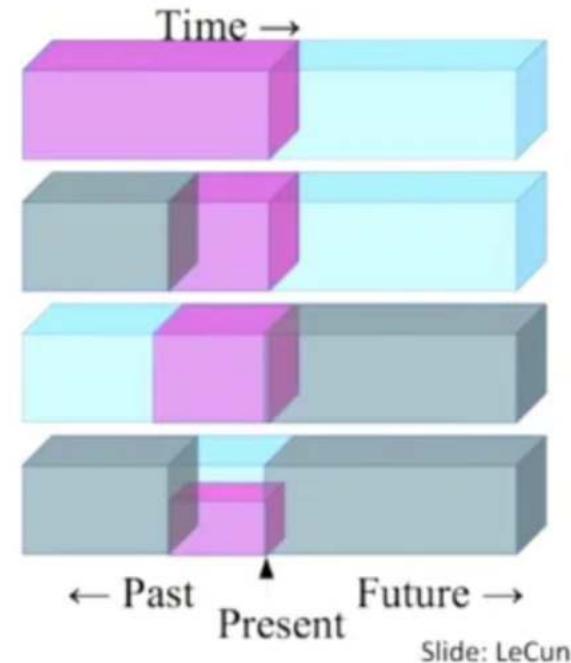
UNSUPERVISED  
AND SELF-SUPERVISED

NO LABELS - DISCOVER  
PATTERNS

$$P(X) \quad P(\hat{Y}|X)$$

## Self-Supervised Learning vs. Unsupervised

- ▶ Predict any part of the input from any other part.
- ▶ Predict the **future** from the **past**.
- ▶ Predict the **future** from the **recent past**.
- ▶ Predict the **past** from the **present**.
- ▶ Predict the **top** from the **bottom**.
- ▶ Predict the **occluded** from the **visible**
- ▶ Pretend there is a part of the input you don't know and predict that.



Slide: LeCun

a self-supervised learning system creates a **pretext task** , which transforms the problem into a “supervised” problem: predict parts of its inputs based on the other parts of its inputs

# Types of Machine Learning

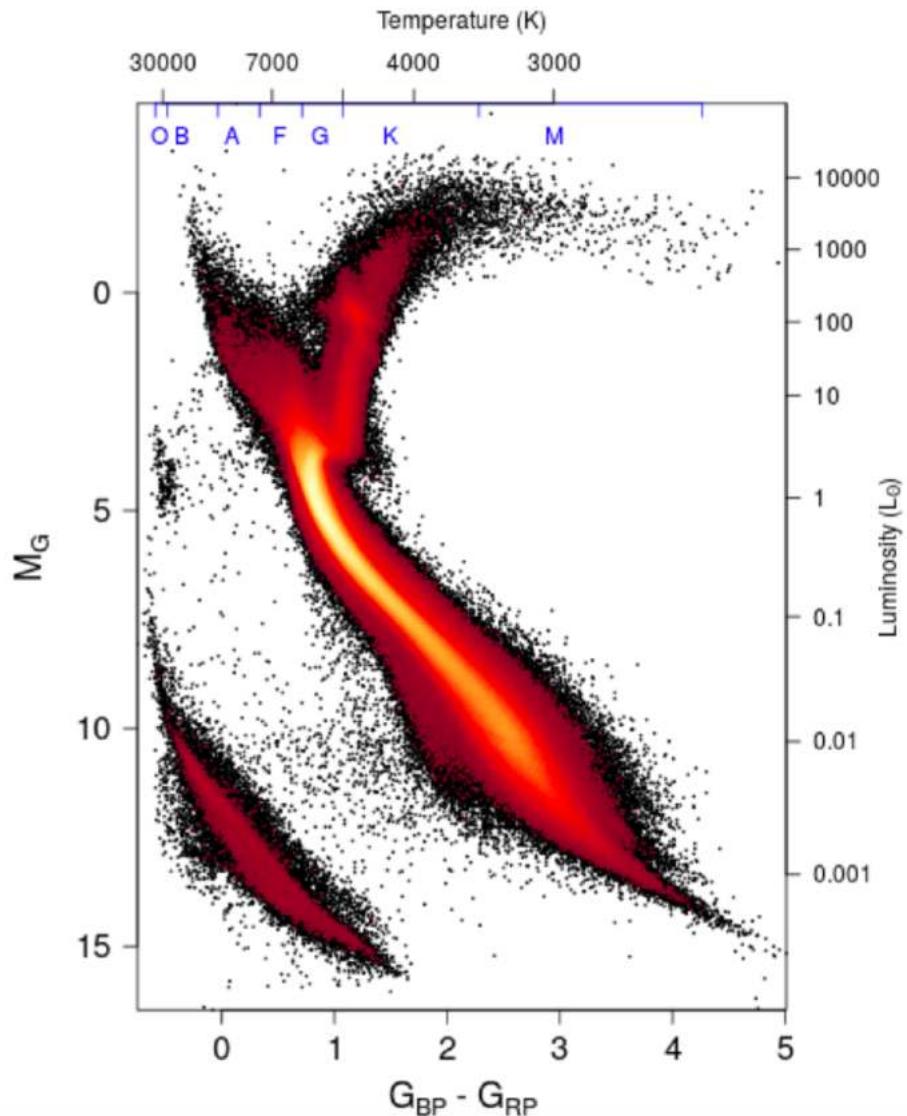
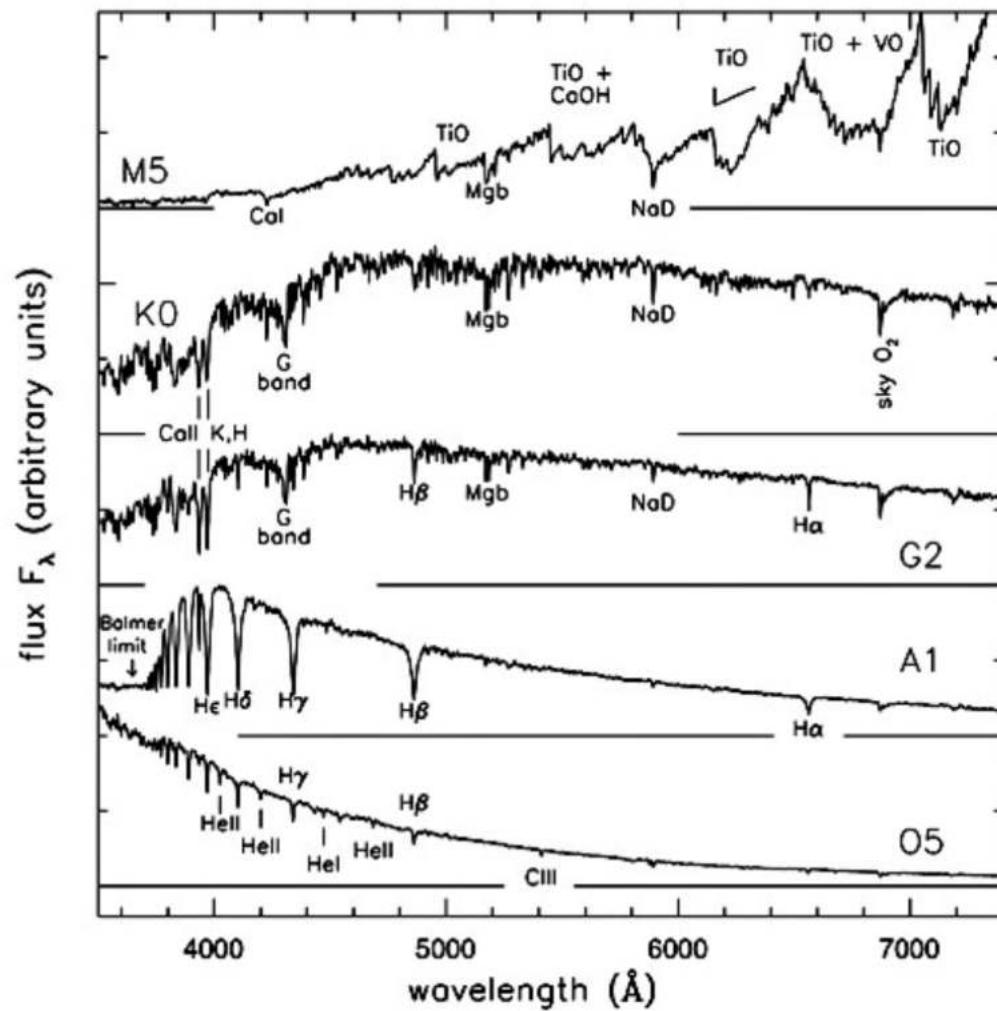
| Approach  | Discriminative  |                                         | Generative      |                                 |
|-----------|-----------------|-----------------------------------------|-----------------|---------------------------------|
|           | Target Function | Method                                  | Target Function | Method                          |
| Labels    | $p(y x)$        | All supervised networks with bottleneck | $p(x y)$        | Conditional generative models   |
| No Labels | $p(\hat{y} x)$  | Self-supervised learning                | $p(x)$          | Generative models, Autoencoders |

There are 3 major applications in astronomy:

- **DIMENSIONALITY REDUCTION:** HOW CAN I REPRESENT MY COMPLEX DATA MORE EFFICIENTLY TO GET NEW INSIGHTS INTO ITS STRUCTURE?
- **SAMPLING:** HOW CAN I INTERPOLATE / EXTRAPOLATE A (SMALL, SPARSE) DATASET TO GENERATE NEW DATA SAMPLED FROM THE SAME (UNKNOWN) DISTRIBUTION?
- **LIKELIHOOD EVALUATION:** WHAT IS THE PROBABILITY THAT A NEW OBSERVATION IS DRAWN FROM THE SAME (UNKNOWN) DISTRIBUTION AS SOME REFERENCE (SMALL, SPARSE) DATASET?

# 1. Dimensionality Reduction

From: Gaia Collaboration et al. 2018



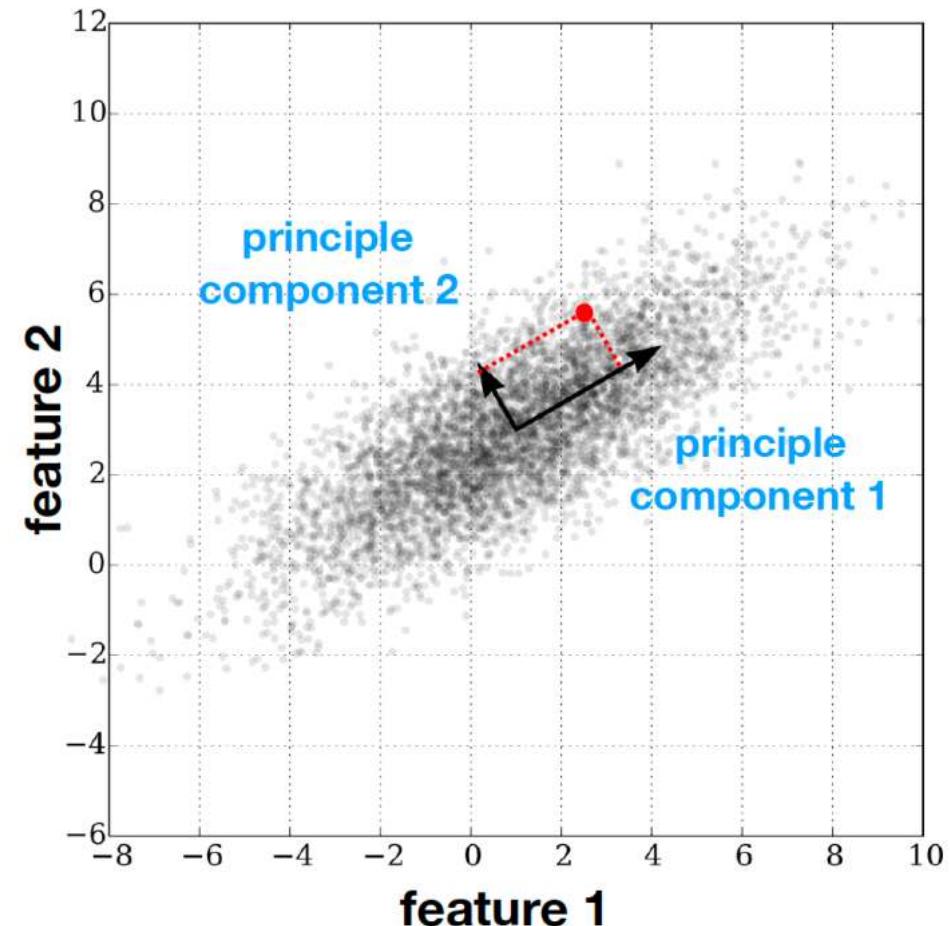
@D. BARON

**THE STELLAR MAIN SEQUENCE IS A  
DIMENSIONALITY REDUCTION**

# PRINCIPAL COMPONENT ANALYSIS

PCA CONVERT A SET OF  
(CORRELATED) VARIABLES INTO A  
SET  
OF VALUES LINEARLY  
UNCORRELATED

1. THE FIRST PRINCIPLE COMPONENT (“PROTOTYPE”), HAS THE LARGEST POSSIBLE VARIANCE
2. THE FOLLOWING COMPONENTS HAVE THE LARGEST VARIANCES WITH THE ADDITIONAL CONSTRAIN THAT THEY ARE ORTHOGONAL TO THE PRECEDING COMPONENTS



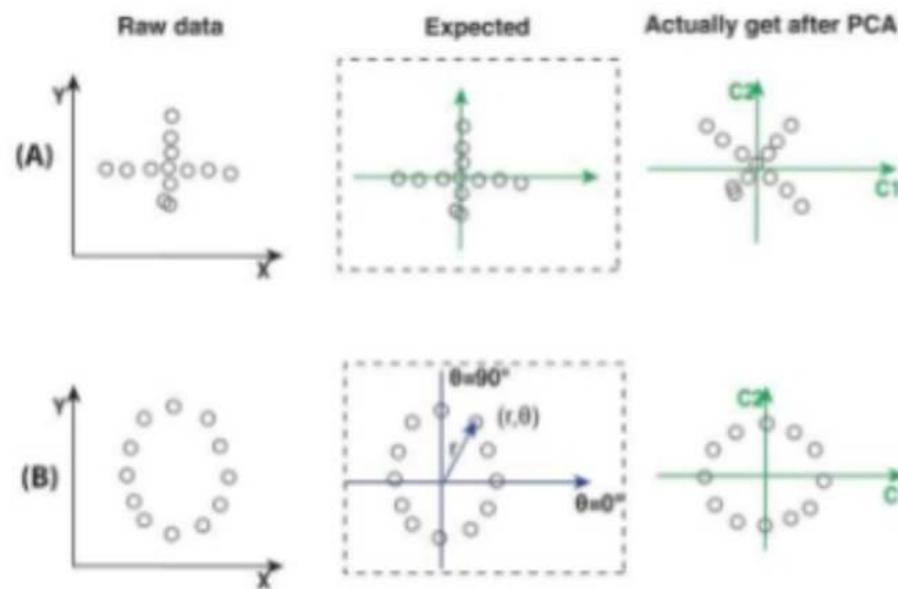
# PRINCIPAL COMPONENT ANALYSIS

1. COMPUTE THE COVARIANCE MATRIX OF YOUR DATA
2. COMPUTE THE EIGENVALUES AND EIGENVECTORS OF THE COVARIANCE MATRIX
3. THE EIGENVECTORS PROVIDE THE DIRECTION OF MAXIMUM VARIANCE AND THE EIGENVALUES THE 'IMPORTANCE' OF THAT PARTICULAR FEATURE

# LIMITATIONS OF PCA

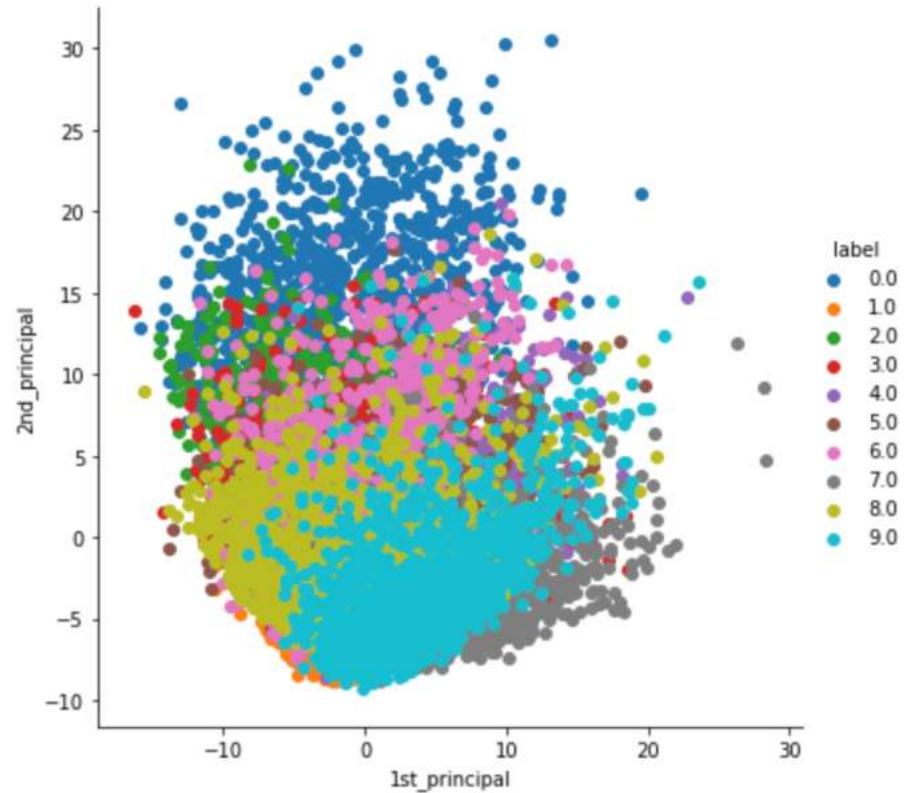
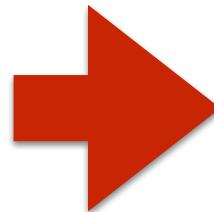
PCA APPLY LINEAR TRANSFORMATIONS

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



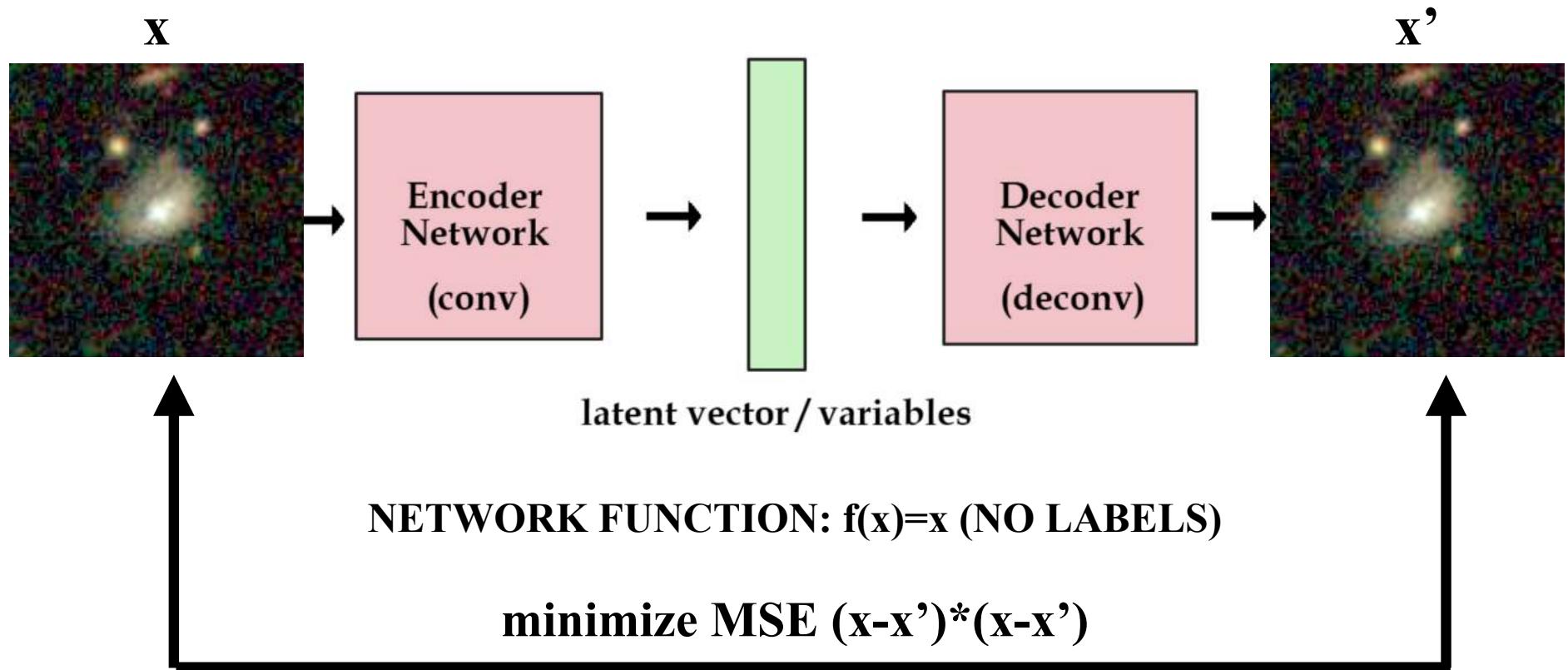
# LIMITATIONS OF PCA

AND DATA IS NOT ALWAYS GAUSSIAN....



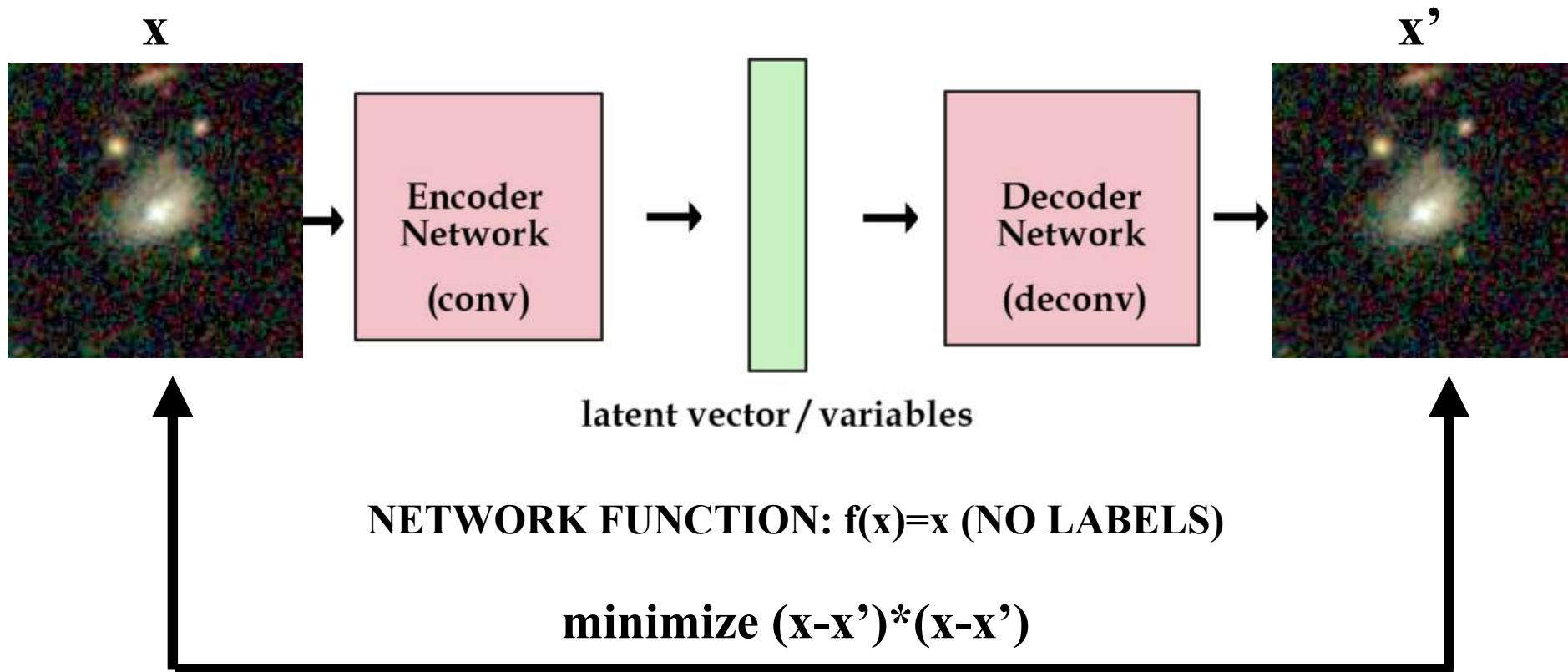
**CAN WE GENERALIZE THAT?**

# CONVOLUTIONAL AUTO-ENCODER



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

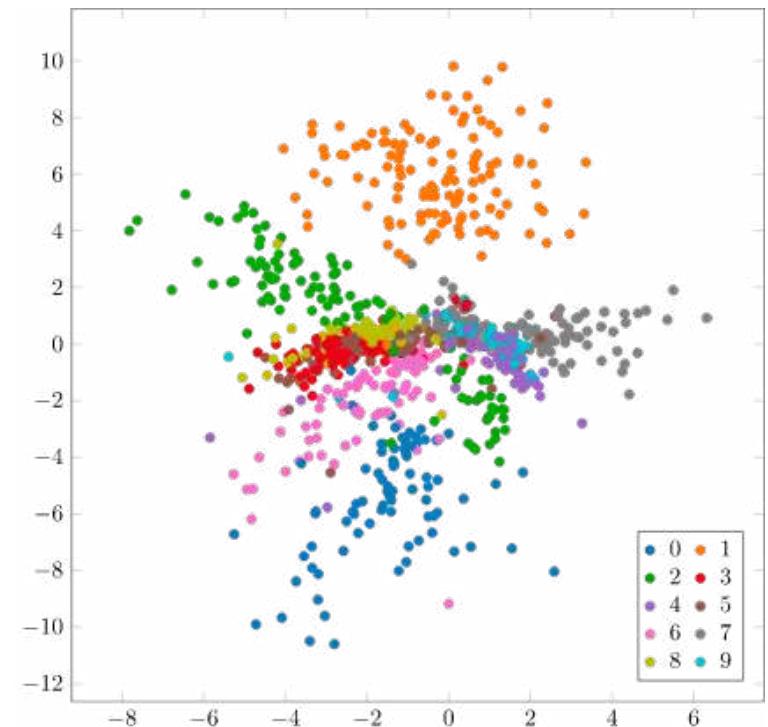
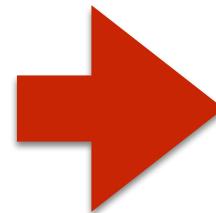
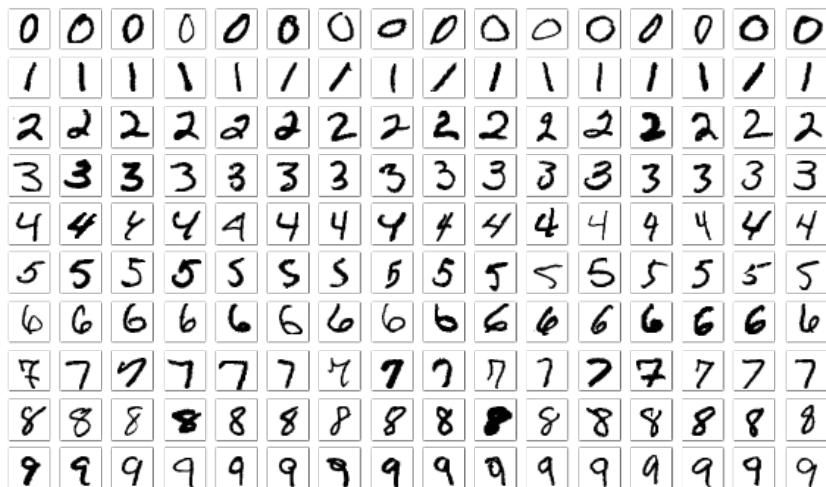
# CONVOLUTIONAL AUTO-ENCODER



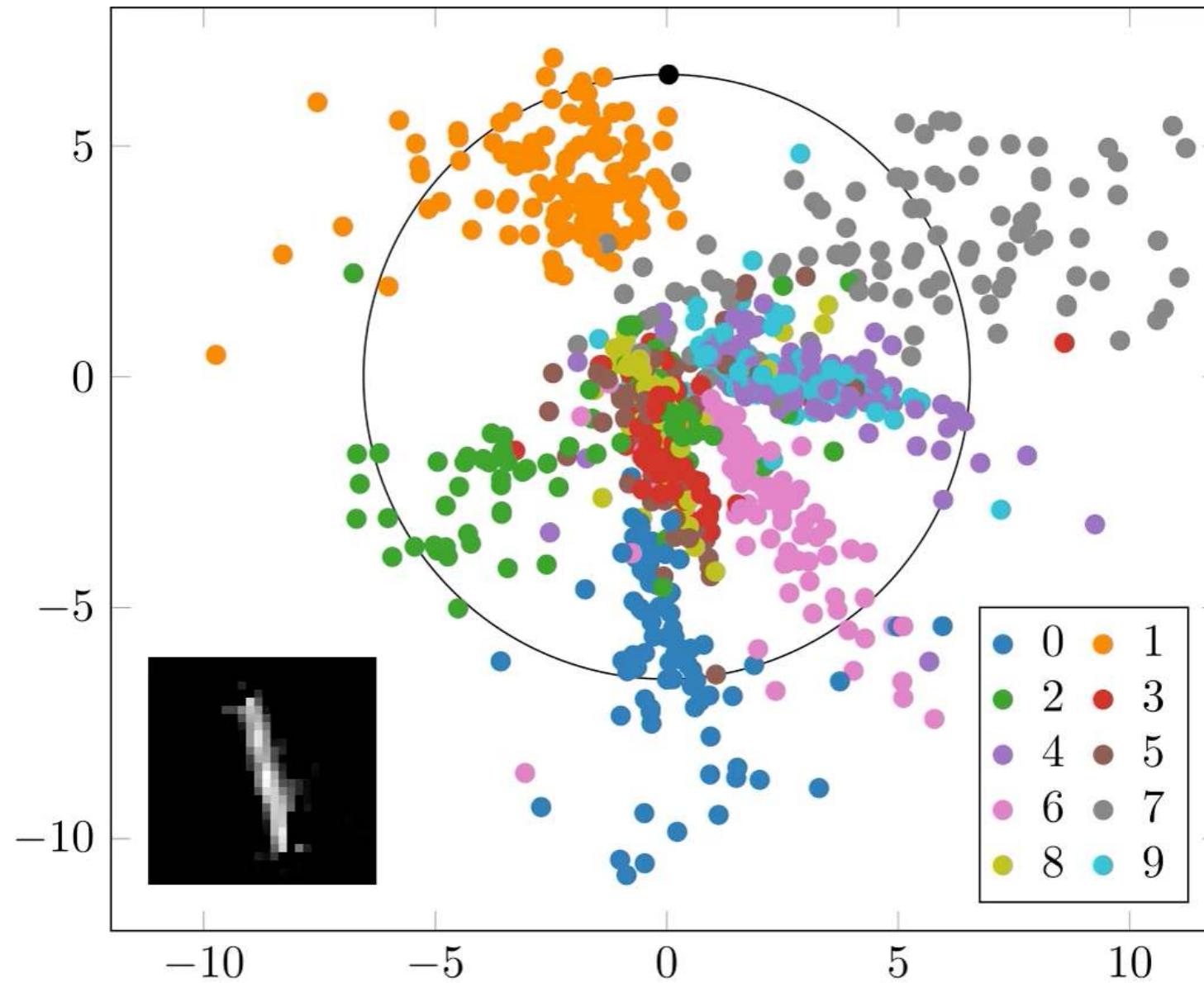
BY REDUCING THE DIMENSIONALITY IN THE LATENT SPACE WE FORCE THE NETWORK TO LEARN A REPRESENTATION OF THE INPUT DATA IN A LOWER DIMENSIONALITY SPACE

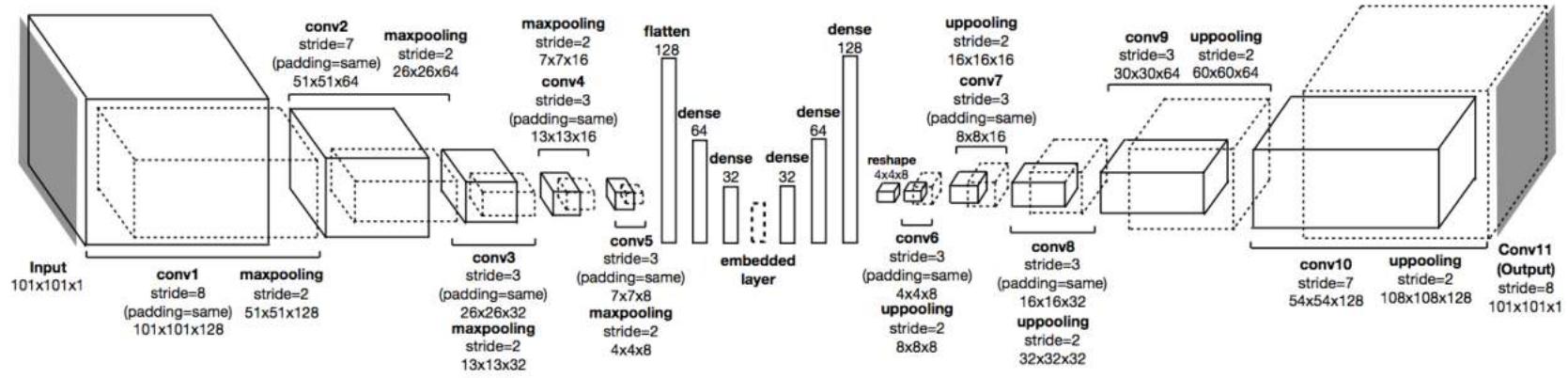
- \* NO NEED TO BE CONVOLUTIONAL - ANY NEURAL NETWORK WITH A BOTTLENECK WILL DO THE JOB
- \* **QUESTION:** WHAT WOULD HAPPEN IF WE SET AN AUTOENCODER WITH NO ACTIVATION FUNCTIONS?

# AUTOENCODER REPRESENTATION OF MNIST

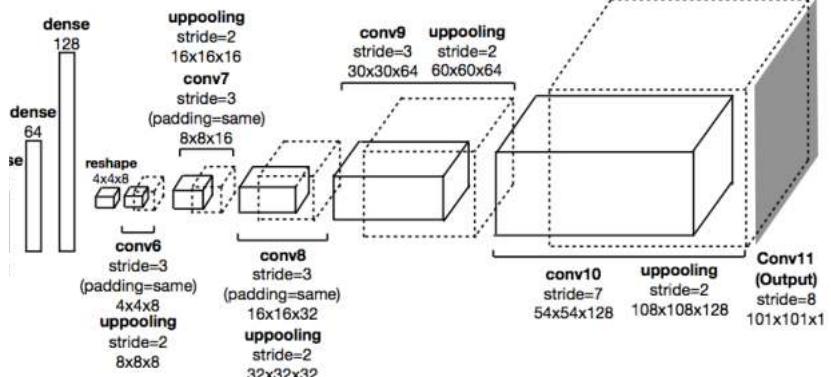
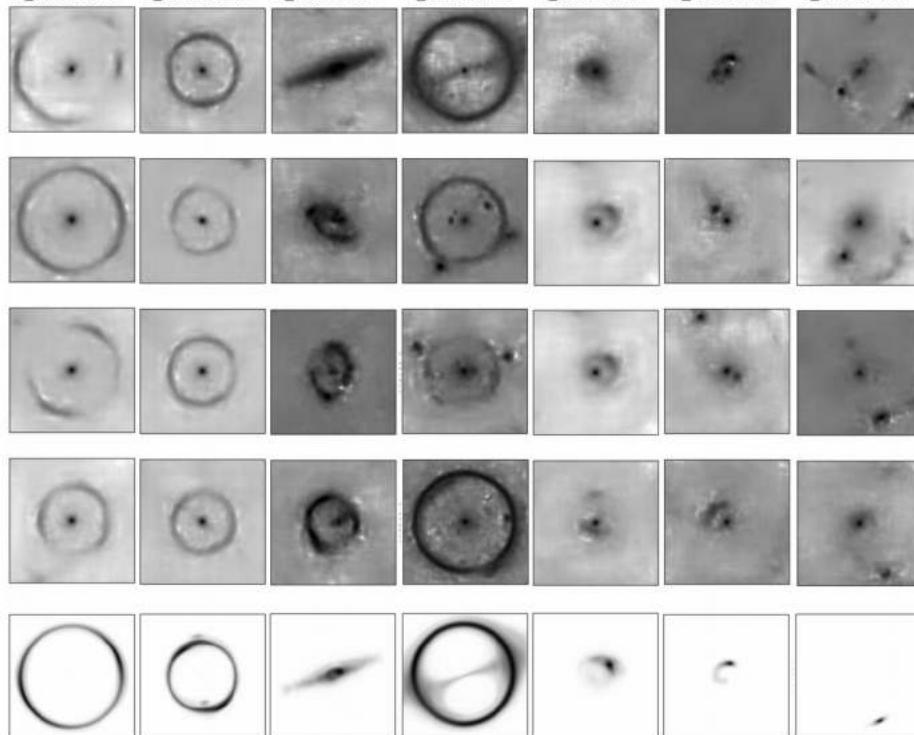


# AUTOENCODER REPRESENTATION OF MNIST

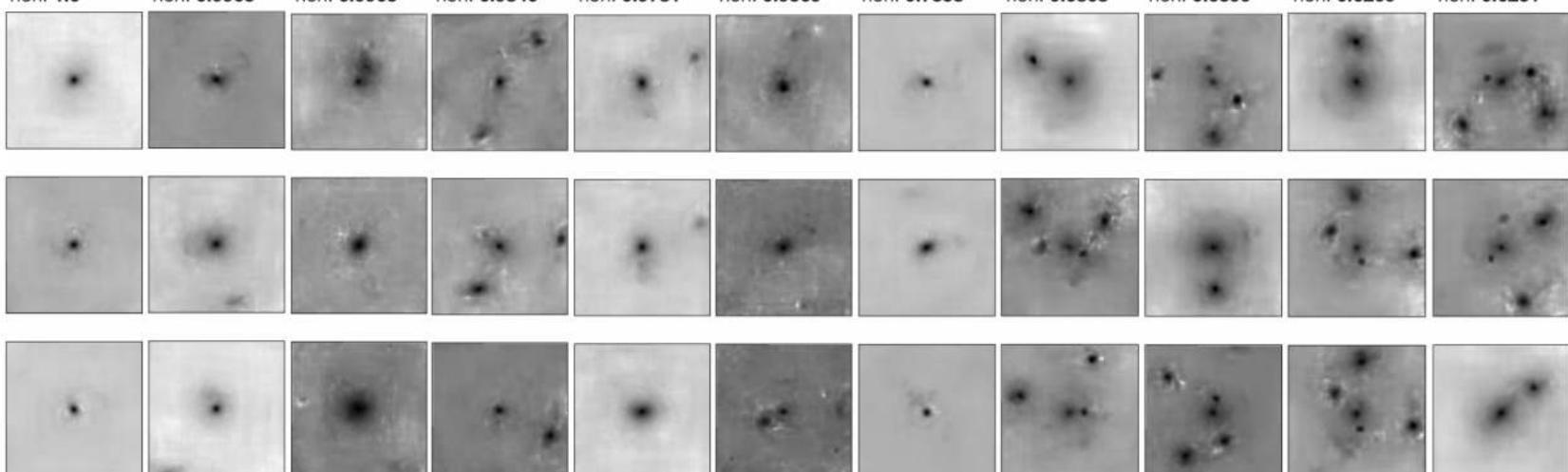




|                                                                       |                                                                       |                                                                      |                                                                      |                                                                    |                                                                       |                                                                      |
|-----------------------------------------------------------------------|-----------------------------------------------------------------------|----------------------------------------------------------------------|----------------------------------------------------------------------|--------------------------------------------------------------------|-----------------------------------------------------------------------|----------------------------------------------------------------------|
| <b>Cluster 17:</b><br>lensing: 0.9873<br>non: 0.0127<br>F_len: 0.0914 | <b>Cluster 21:</b><br>lensing: 0.9448<br>non: 0.0552<br>F_len: 0.0695 | <b>Cluster 1:</b><br>lensing: 0.9159<br>non: 0.0841<br>F_len: 0.0731 | <b>Cluster 6:</b><br>lensing: 0.8997<br>non: 0.1003<br>F_len: 0.0729 | <b>Cluster 2:</b><br>lensing: 0.803<br>non: 0.197<br>F_len: 0.1945 | <b>Cluster 20:</b><br>lensing: 0.6206<br>non: 0.3794<br>F_len: 0.0575 | <b>Cluster 5:</b><br>lensing: 0.6170<br>non: 0.3830<br>F_len: 0.0734 |
|-----------------------------------------------------------------------|-----------------------------------------------------------------------|----------------------------------------------------------------------|----------------------------------------------------------------------|--------------------------------------------------------------------|-----------------------------------------------------------------------|----------------------------------------------------------------------|



|                                                |                                                     |                                                     |                                                      |                                                      |                                                     |                                                     |                                                      |                                                      |                                                      |                                                     |
|------------------------------------------------|-----------------------------------------------------|-----------------------------------------------------|------------------------------------------------------|------------------------------------------------------|-----------------------------------------------------|-----------------------------------------------------|------------------------------------------------------|------------------------------------------------------|------------------------------------------------------|-----------------------------------------------------|
| <b>Cluster 14:</b><br>lensing: 0.0<br>non: 1.0 | <b>Cluster 3:</b><br>lensing: 0.0037<br>non: 0.9963 | <b>Cluster 8:</b><br>lensing: 0.0037<br>non: 0.9963 | <b>Cluster 22:</b><br>lensing: 0.0154<br>non: 0.9846 | <b>Cluster 16:</b><br>lensing: 0.0219<br>non: 0.9781 | <b>Cluster 4:</b><br>lensing: 0.0431<br>non: 0.9569 | <b>Cluster 0:</b><br>lensing: 0.2642<br>non: 0.7358 | <b>Cluster 18:</b><br>lensing: 0.3132<br>non: 0.6868 | <b>Cluster 19:</b><br>lensing: 0.3601<br>non: 0.6399 | <b>Cluster 13:</b><br>lensing: 0.3731<br>non: 0.6269 | <b>Cluster 9:</b><br>lensing: 0.3769<br>non: 0.6231 |
|------------------------------------------------|-----------------------------------------------------|-----------------------------------------------------|------------------------------------------------------|------------------------------------------------------|-----------------------------------------------------|-----------------------------------------------------|------------------------------------------------------|------------------------------------------------------|------------------------------------------------------|-----------------------------------------------------|



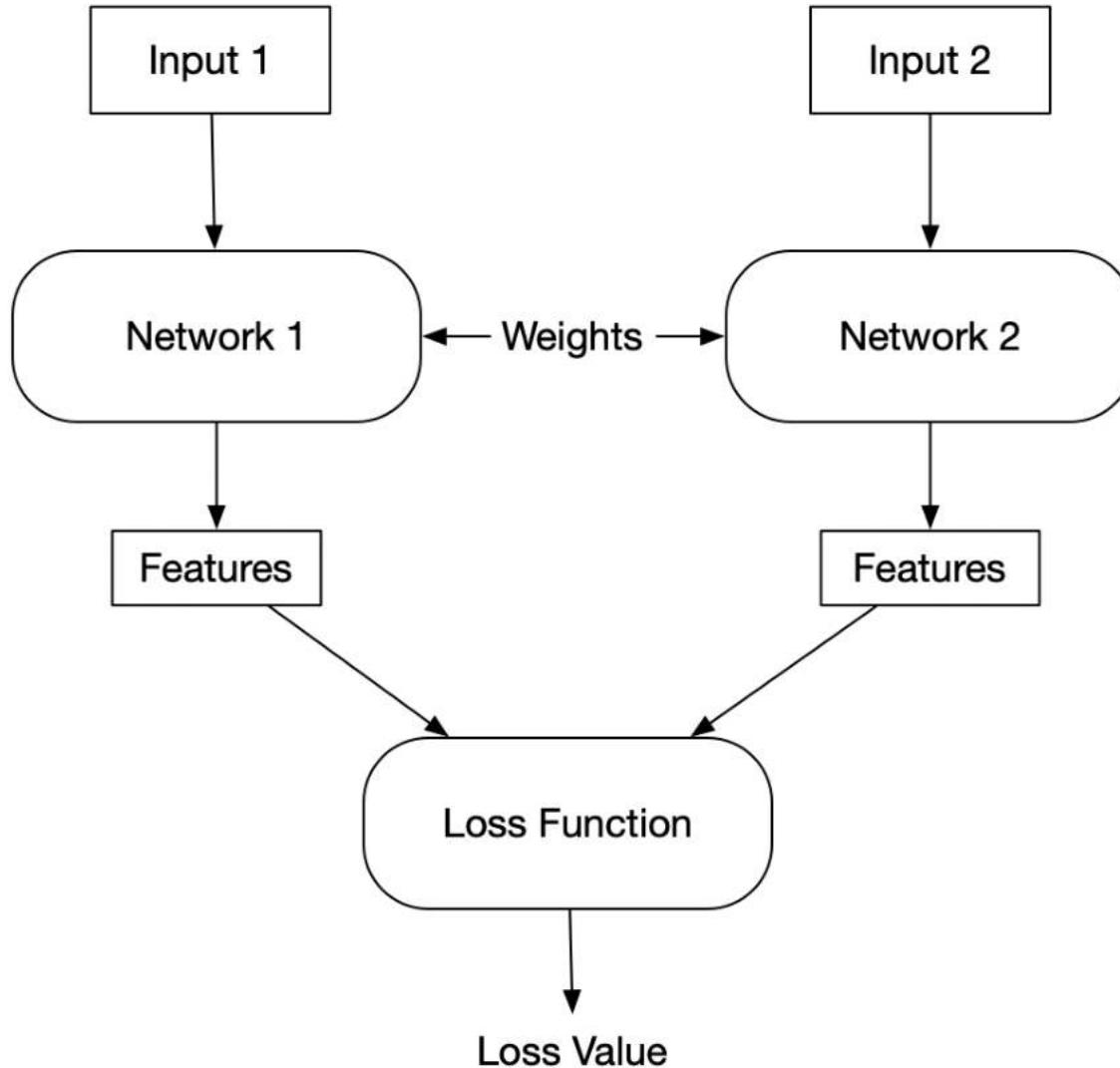
# **SELF-SUPERVISED CONTRASTIVE LEARNING**

Humans tend to identify objects without remembering all the details, by creating some abstract representations which are used to identify new objects of similar time.

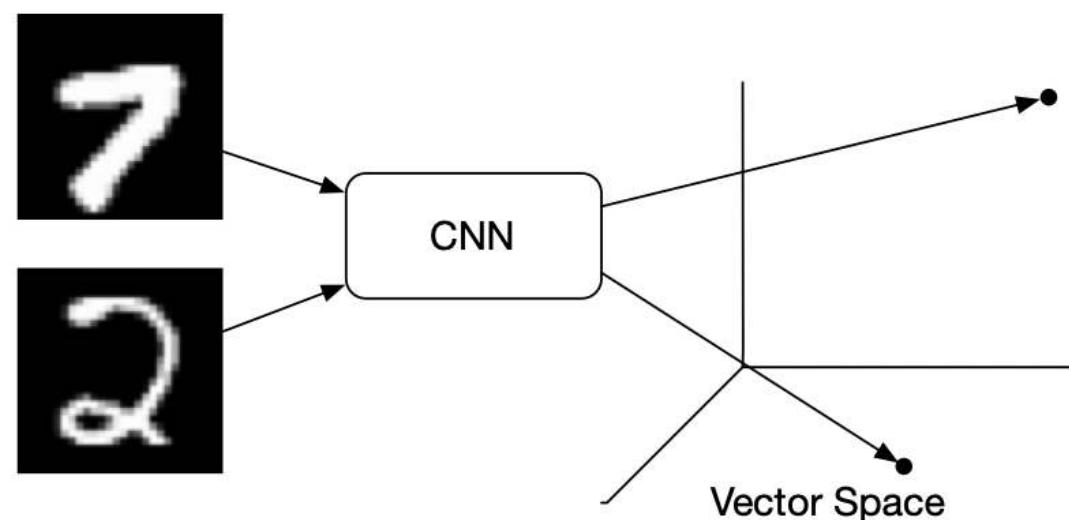
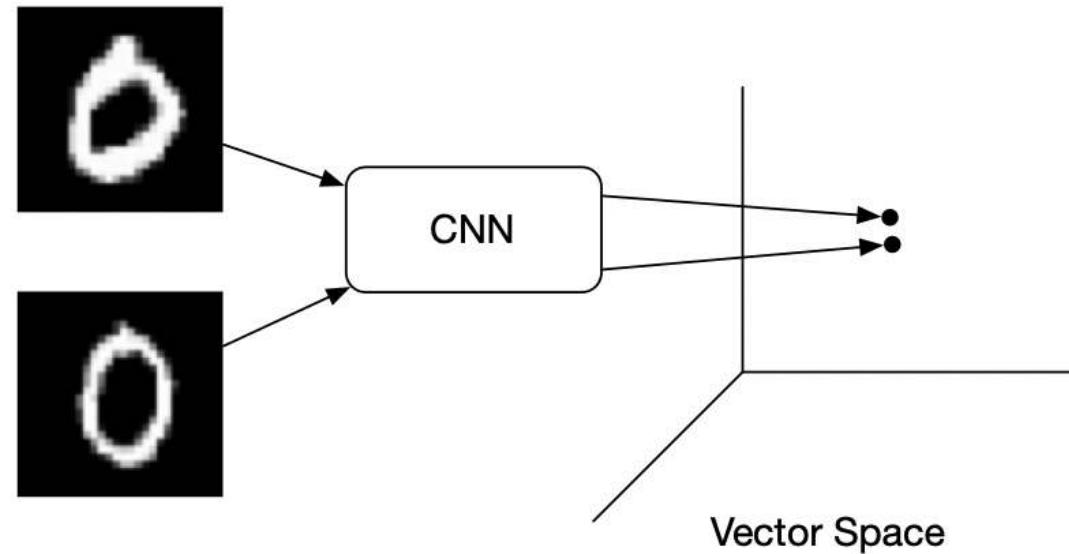
The main goal of self-supervised contrastive learning is to create and generalize these representations.

Therefore, it aims at building some general representations that can be used for other “downstream tasks”

Siamese networks are two networks sharing weights and common loss function



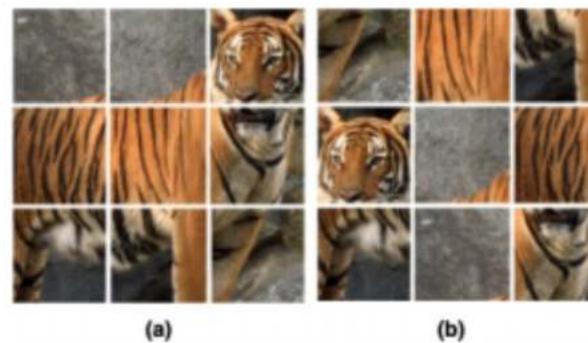
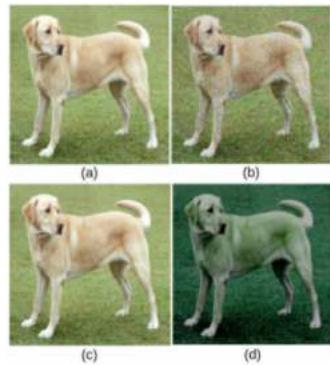
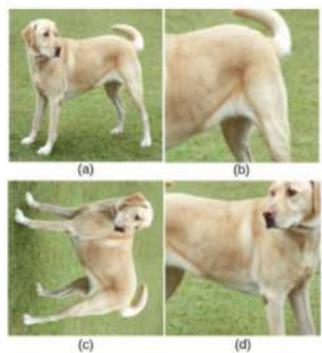
In a supervised setting, similar objects are used to be together by simply using a crossentropy loss



In most cases, labels are not available.

How do you recognise similar objects?

What would be the loss function?

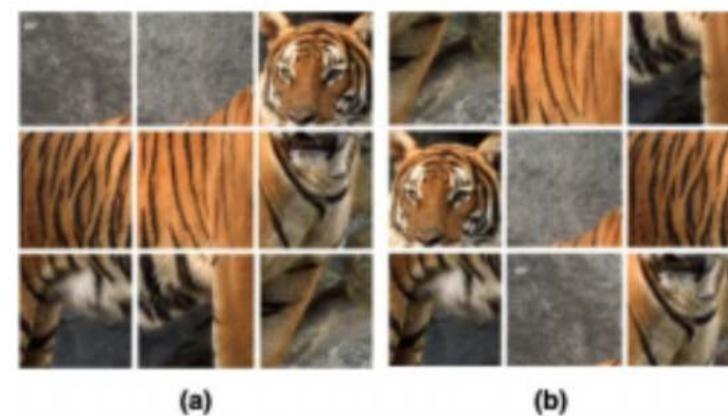
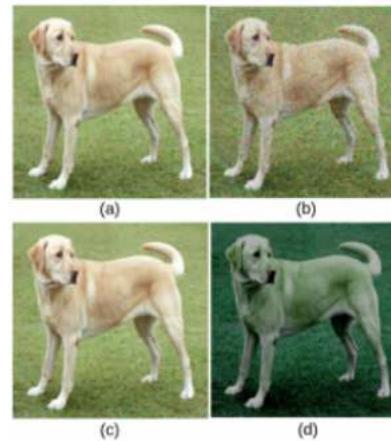
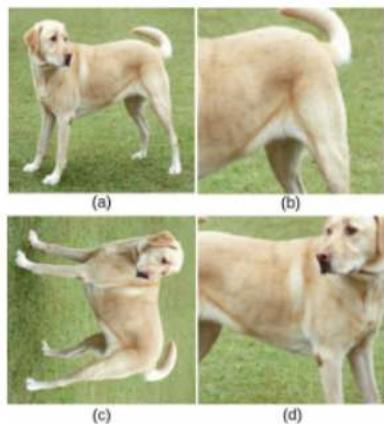


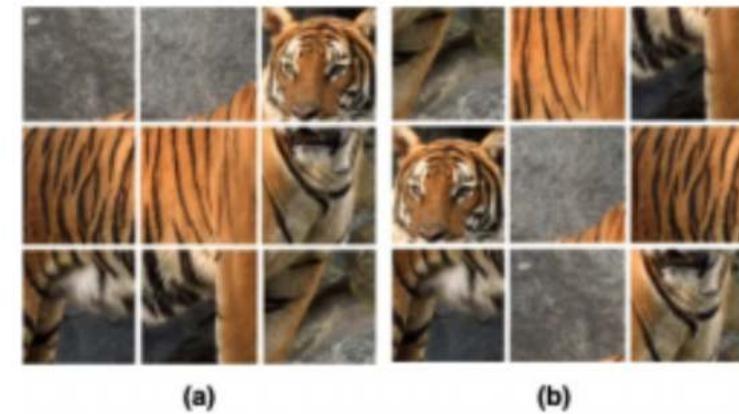
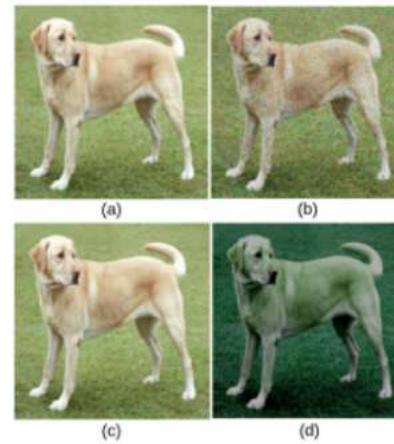
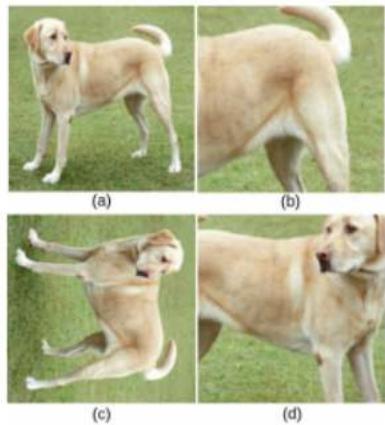
In most cases, labels are not available.

How do you recognise similar objects?

What would be the loss function?

We create “pretext tasks” from a unique image:





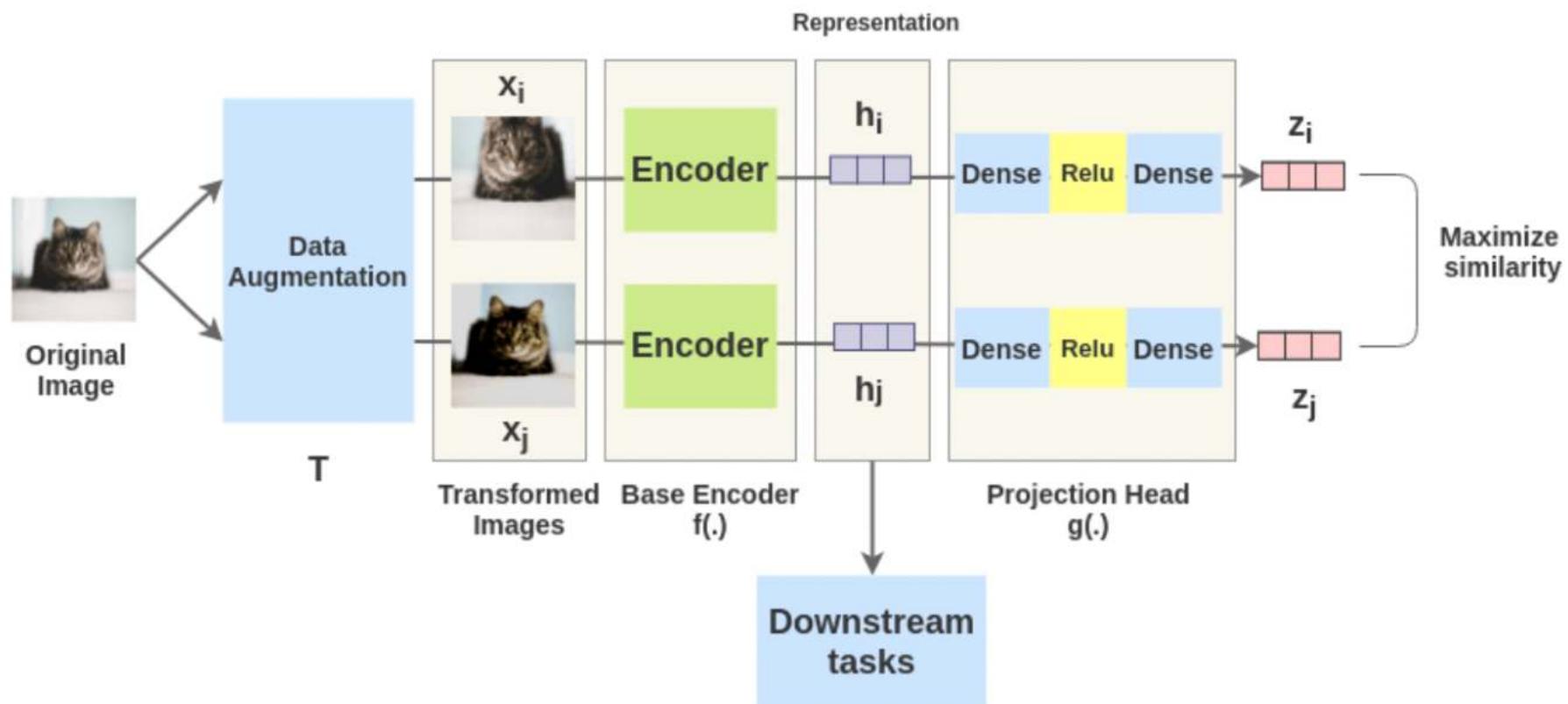
## Color Augmentation

## Image Rotation

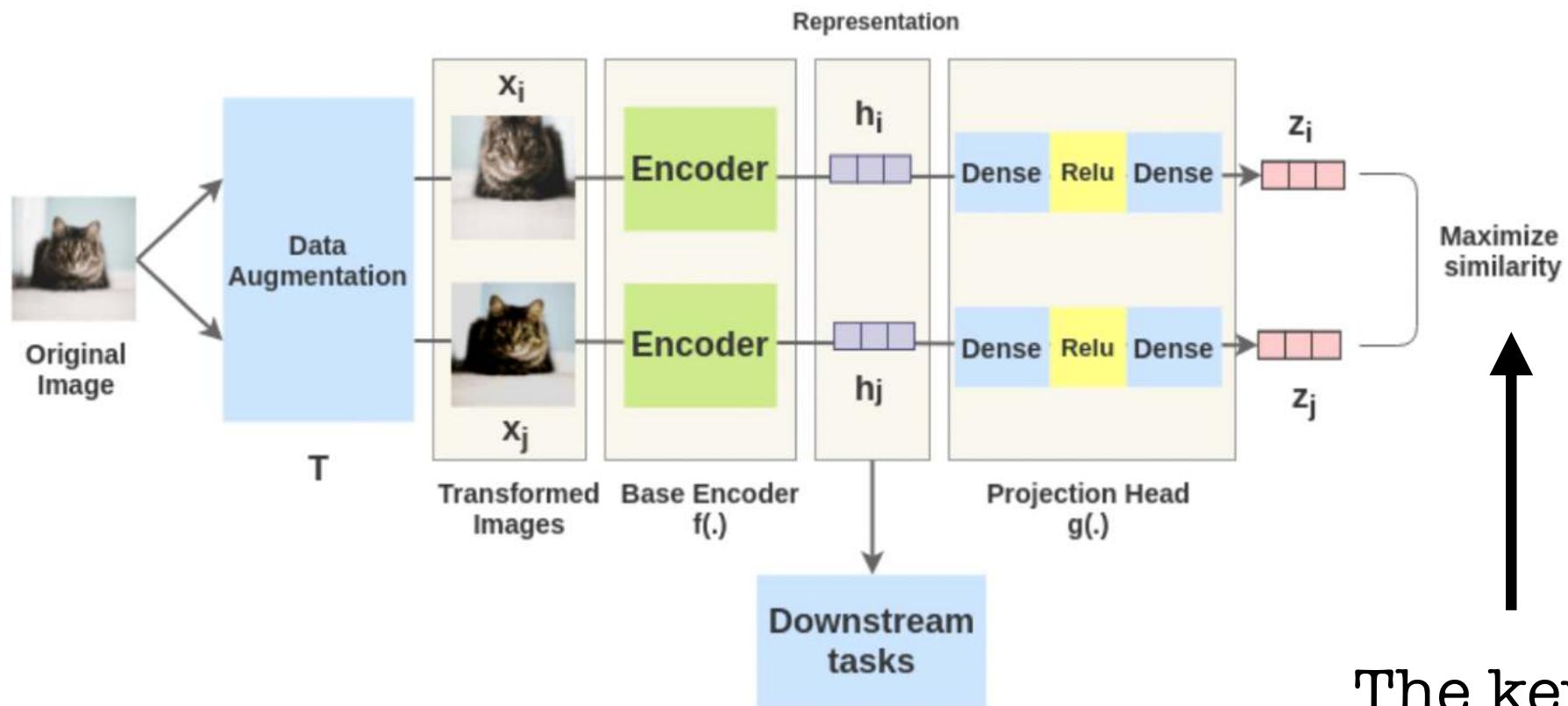
## Image Cropping

## 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

Similarly between  
two representations of positive pairs

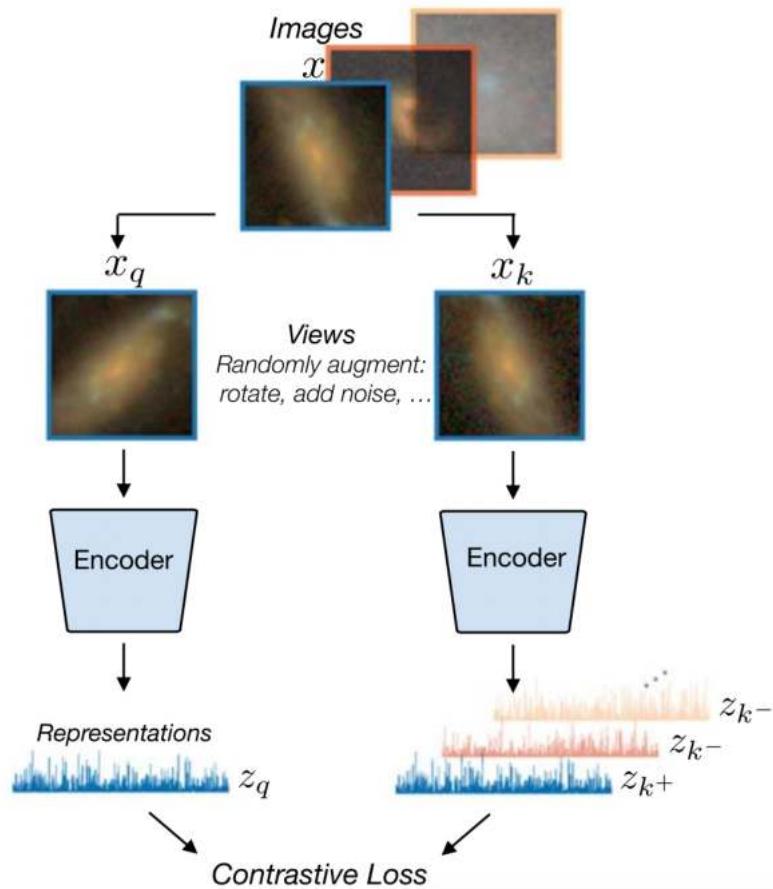
The contrastive loss:

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

Sum of all similarities between  
negative pairs

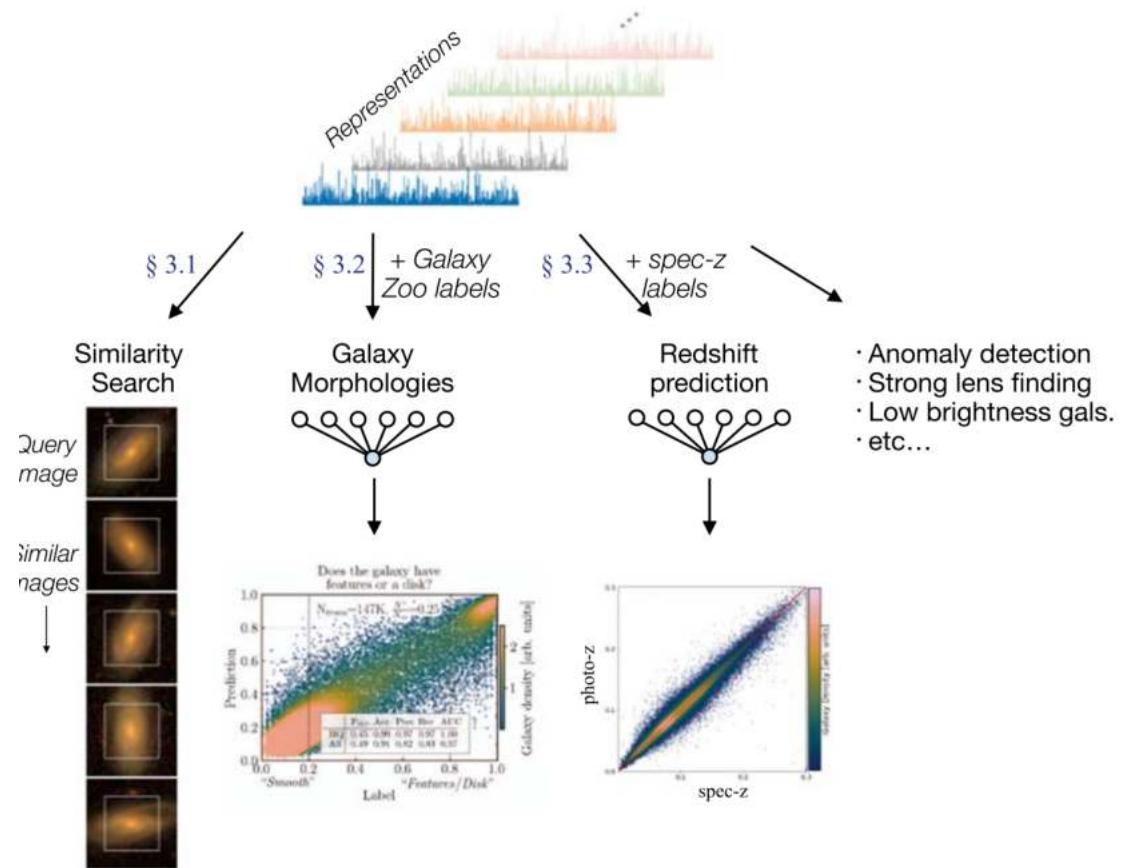
## 1. Self-supervised contrastive representation learning

Learn representations in an unsupervised manner



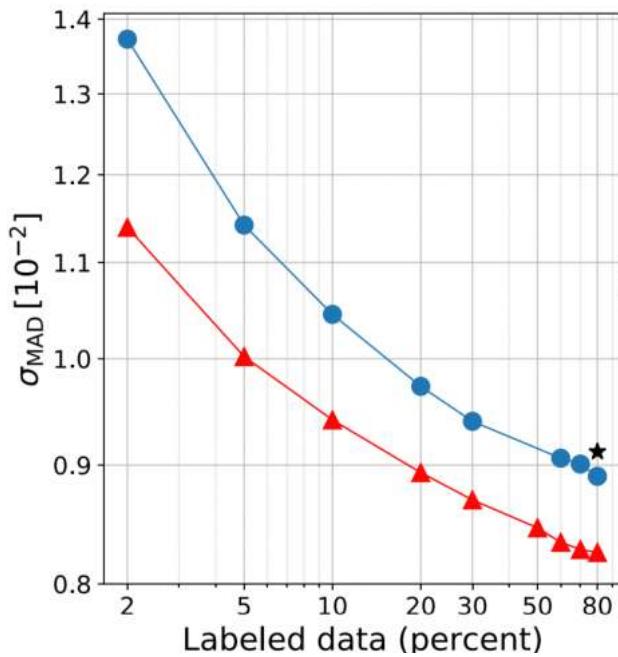
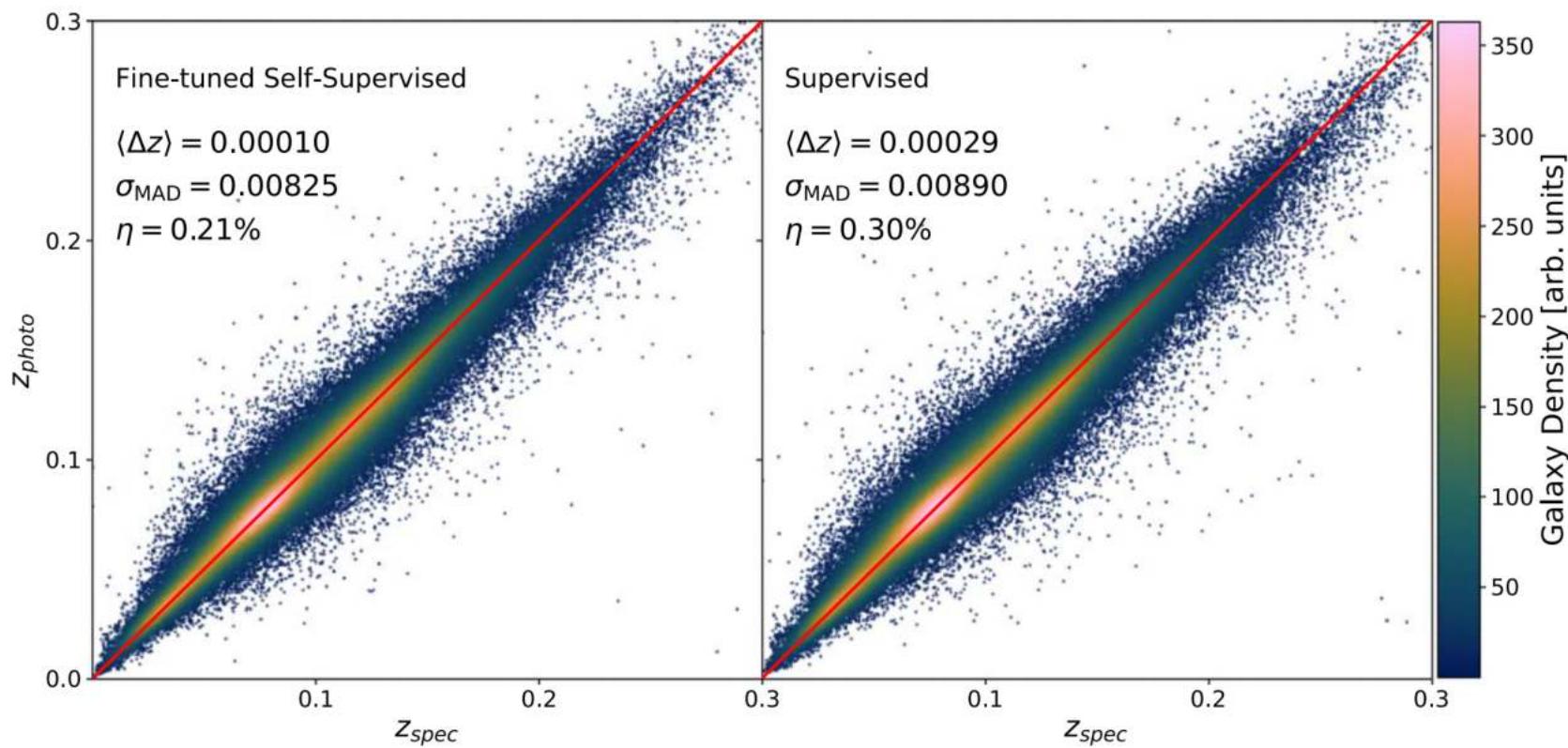
## 2. Downstream tasks

Use representations for a variety of applications

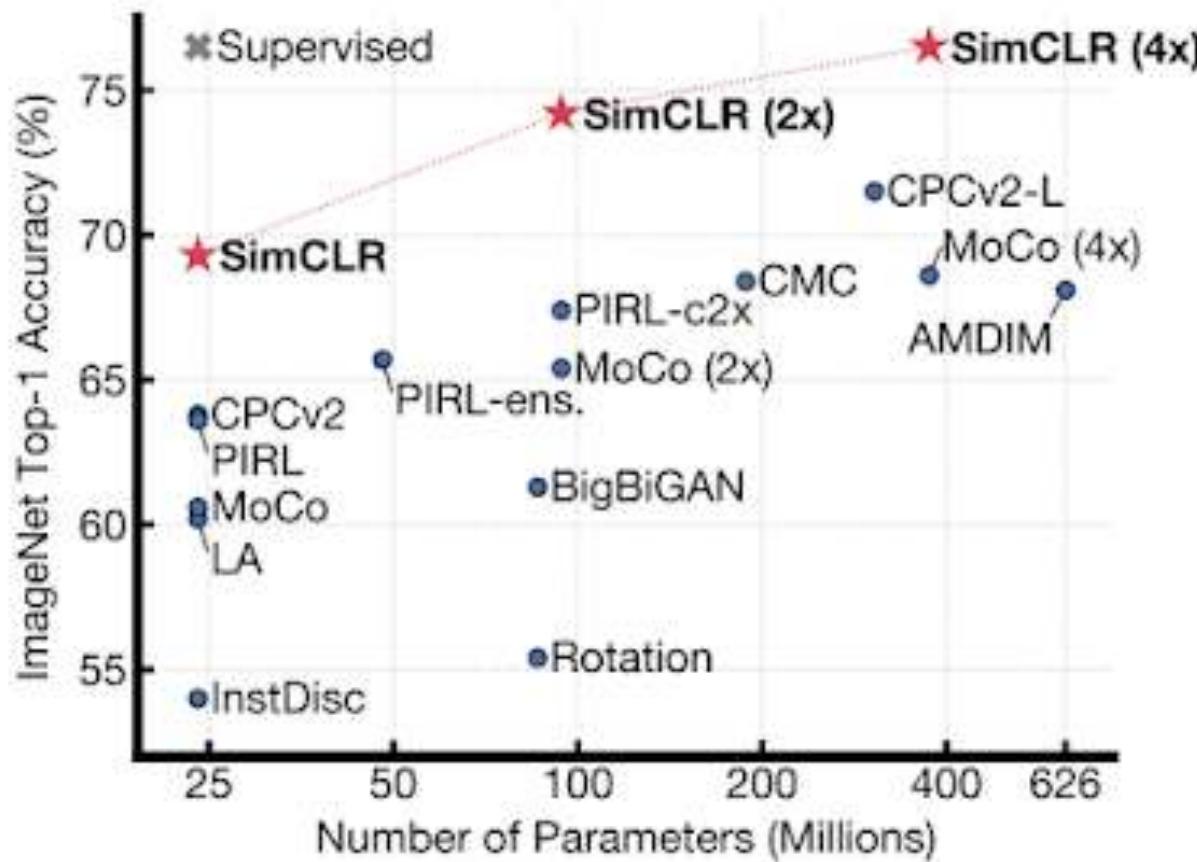


Hayat+21

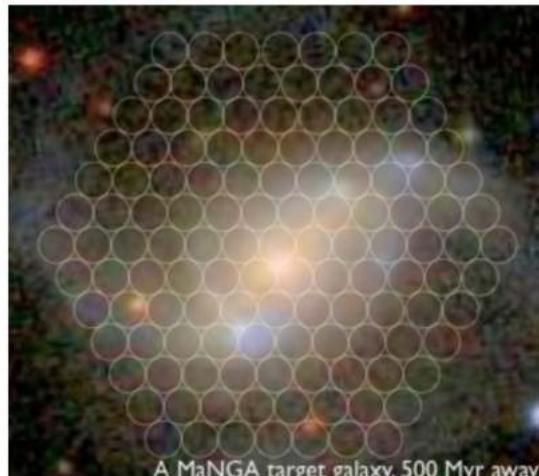
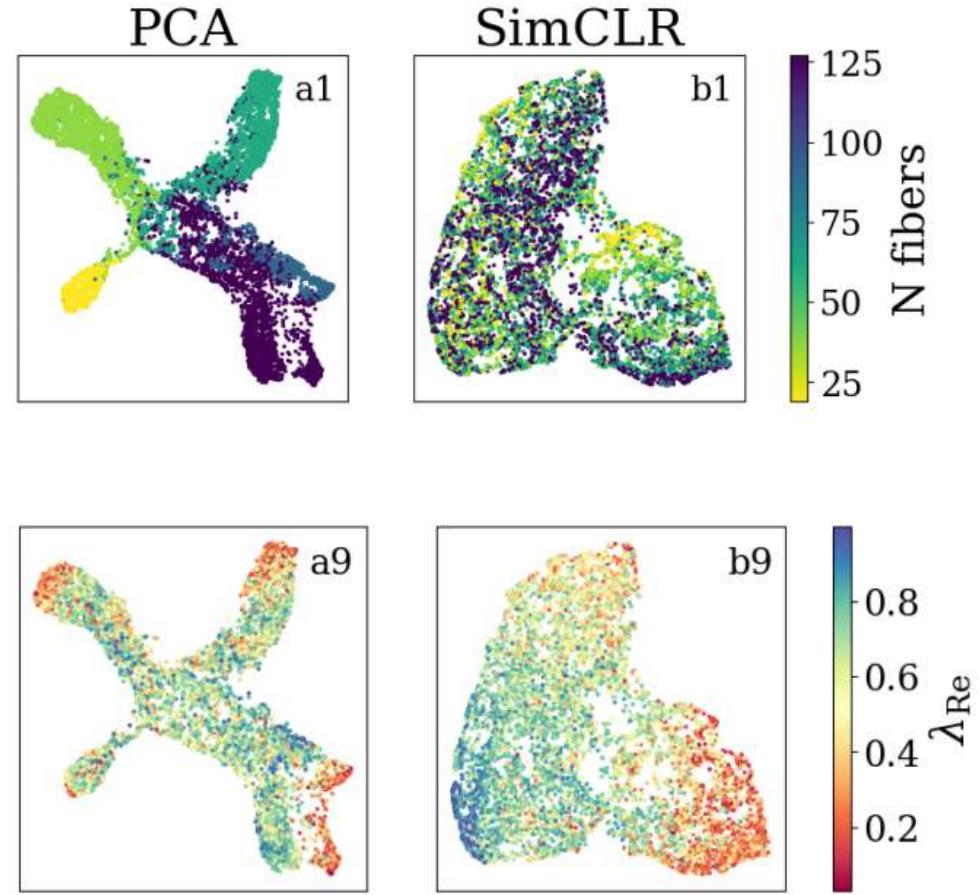
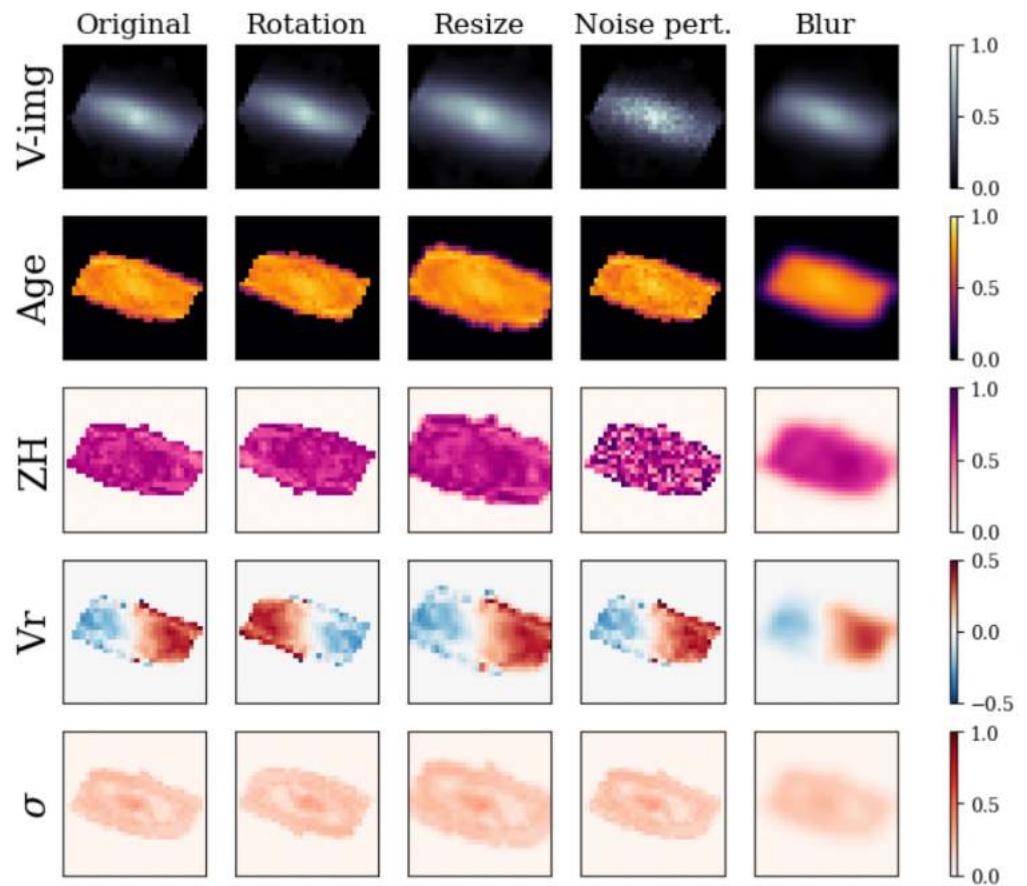
# Hayat+21



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



## 2. Sampling



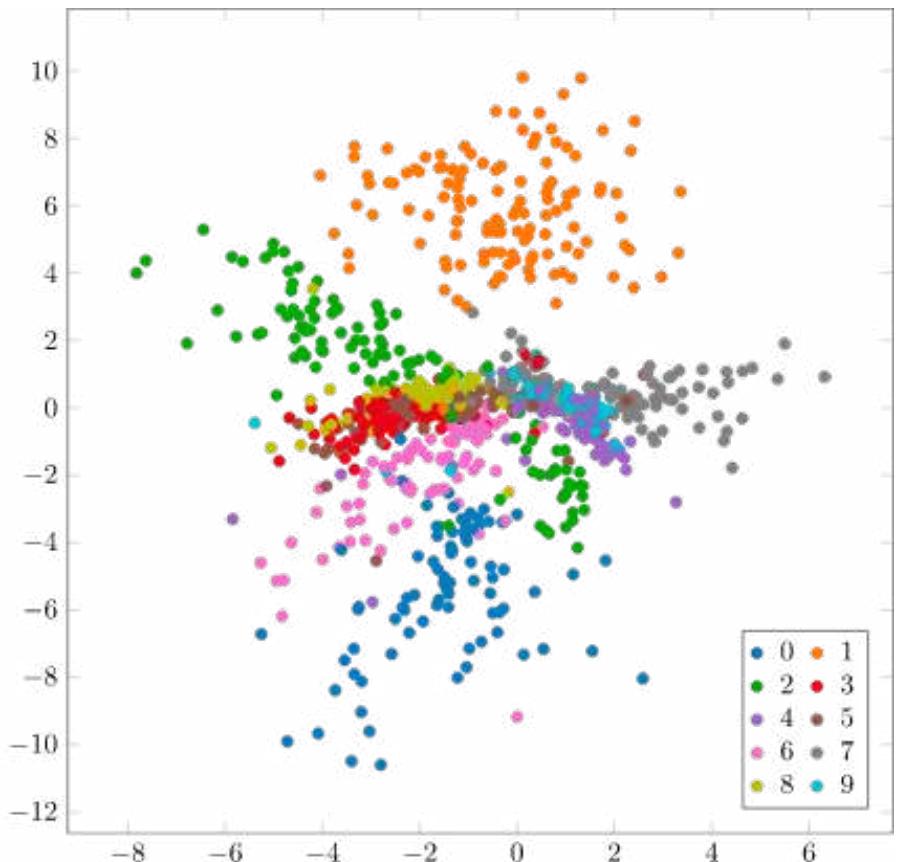
## Contrastive learning representation of Manga galaxies

Sarmiento+21

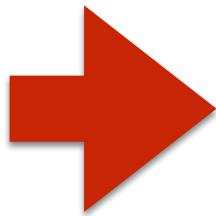
HOW MIGHT YOU SOLVE THESE  
RELATED PROBLEMS ONCE YOU HAVE  
THE LATENT SPACE COORDINATES FOR  
YOUR TRAINING SAMPLE?

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**”)

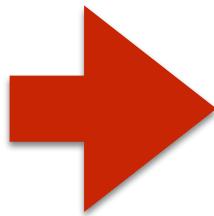
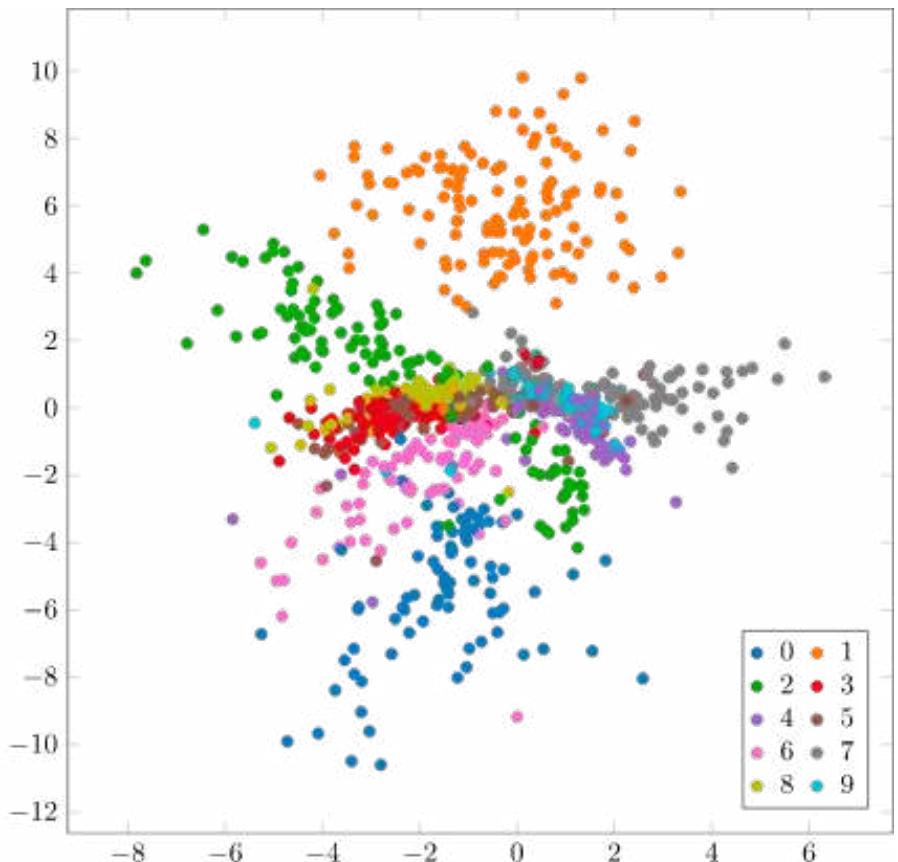


z space (latent space)



HOW CAN I  
ESTIMATE  $P(X)$ ?

$(P(z|x))$

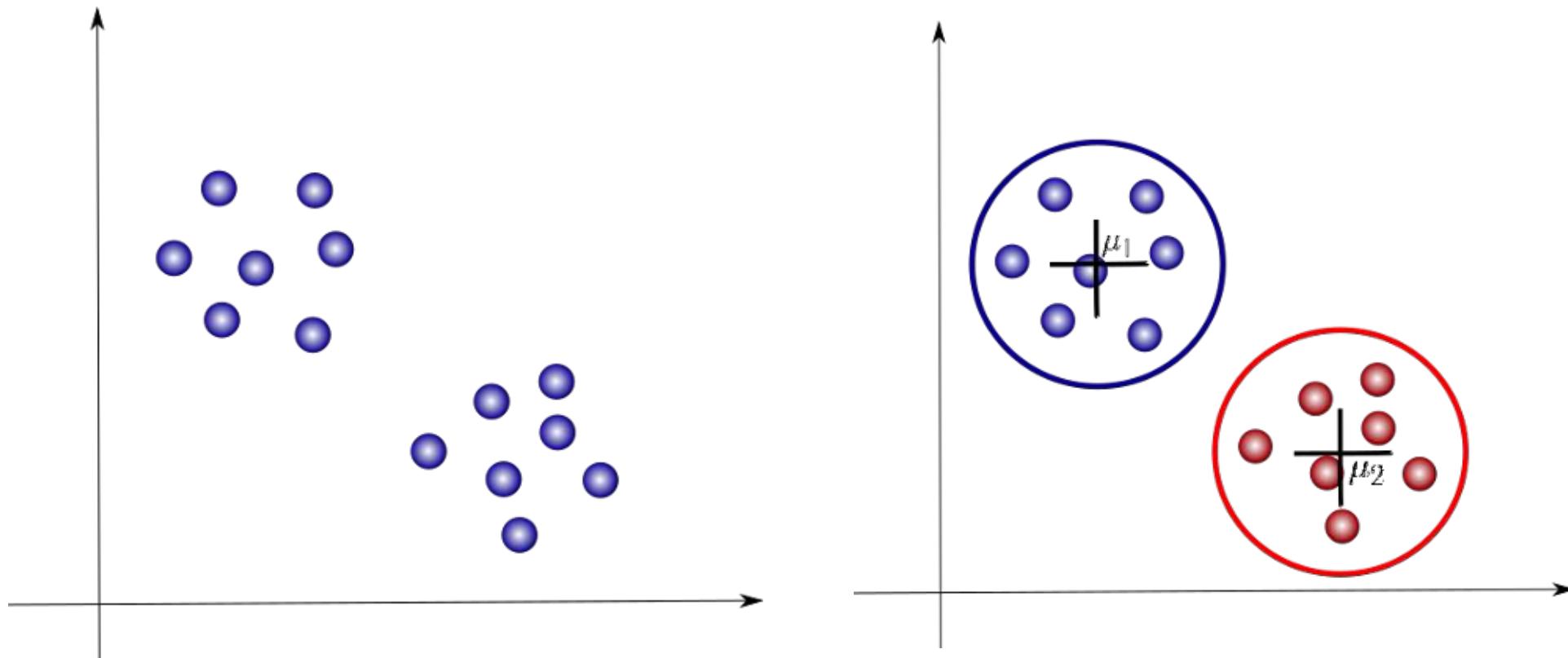


HOW CAN I  
ESTIMATE  $P(X)$ ?

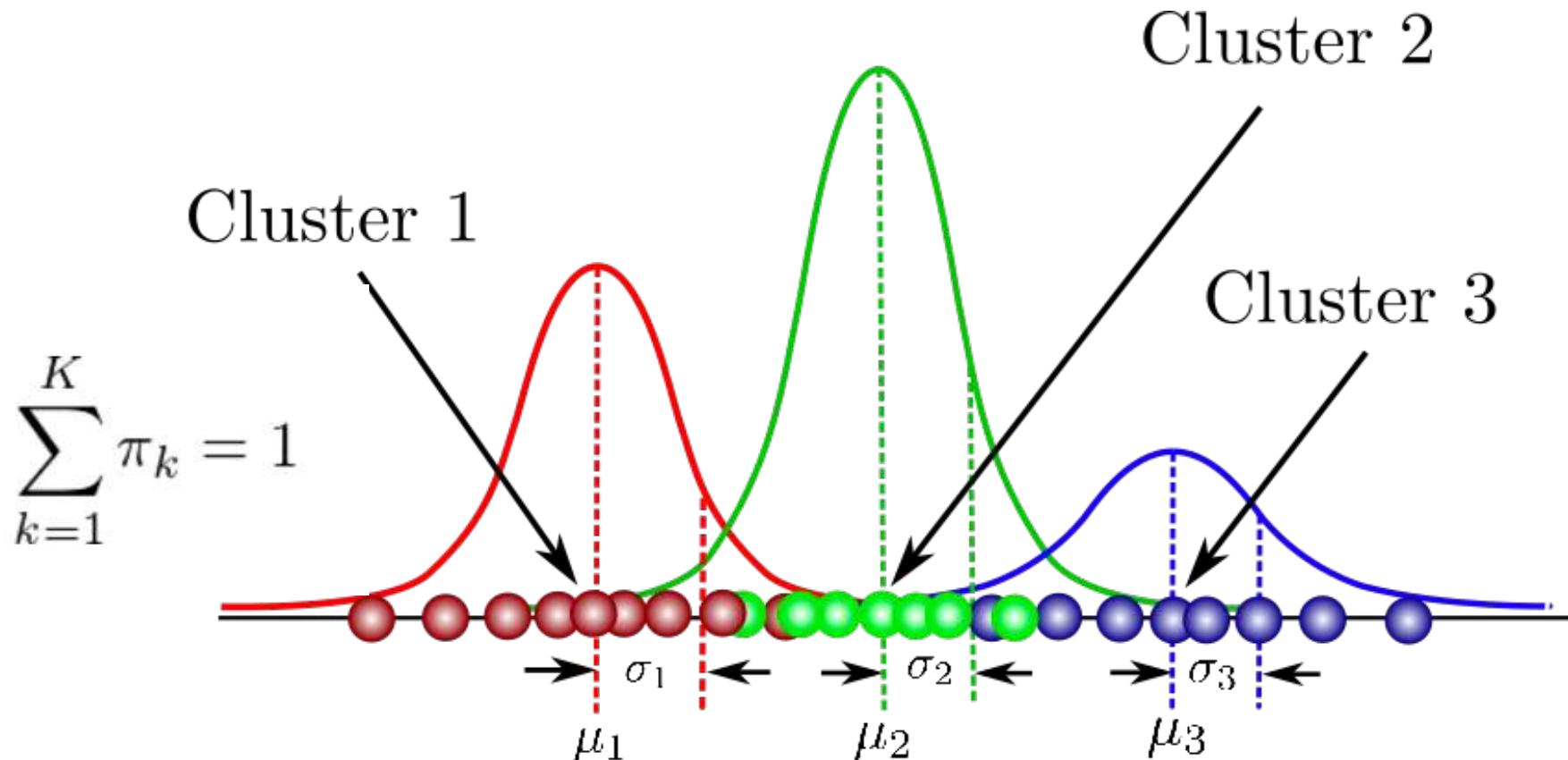
$(P(z|x))$

WHEN YOU DON'T KNOW, ASSUME IT IS GAUSSIAN....

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

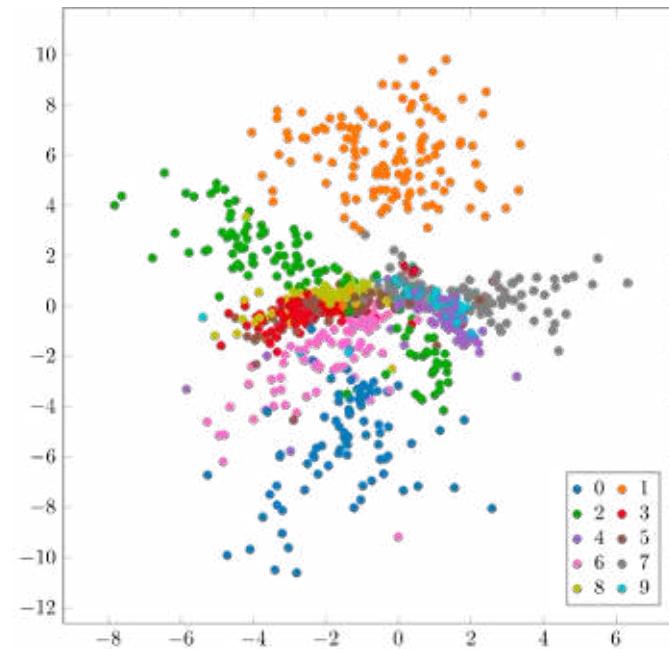
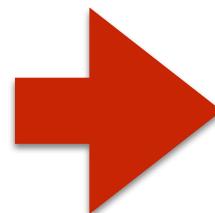
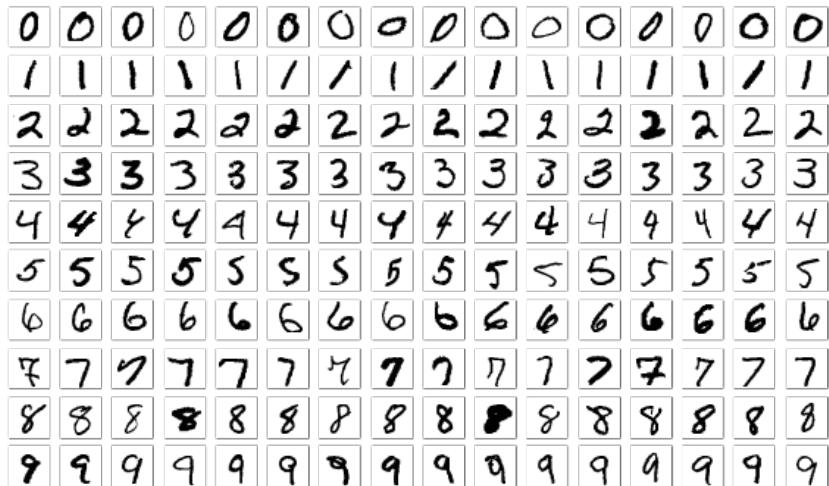


# GAUSSIAN MIXTURE MODELS (GMMs) ARE DENSITY ESTIMATOR METHODS THAT FIT MULTIPLE GAUSSIANS TO THE REPRESENTATION

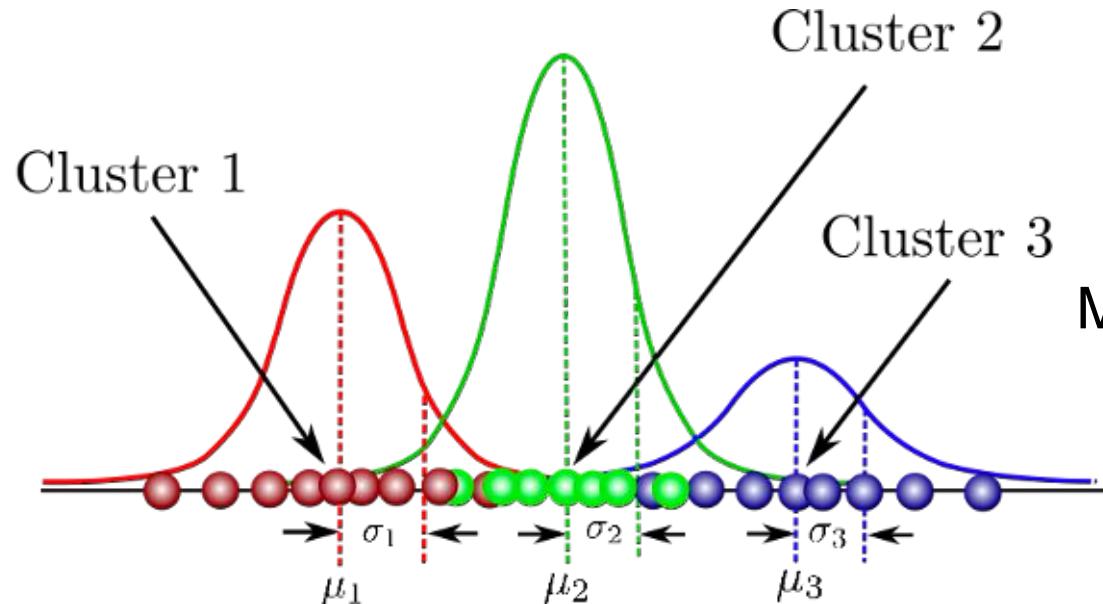


means, sigmas and scale factors of each gaussian are free parameters

# DATA



REPRESENTATION  
(AUTOENCODER -  $z$ )



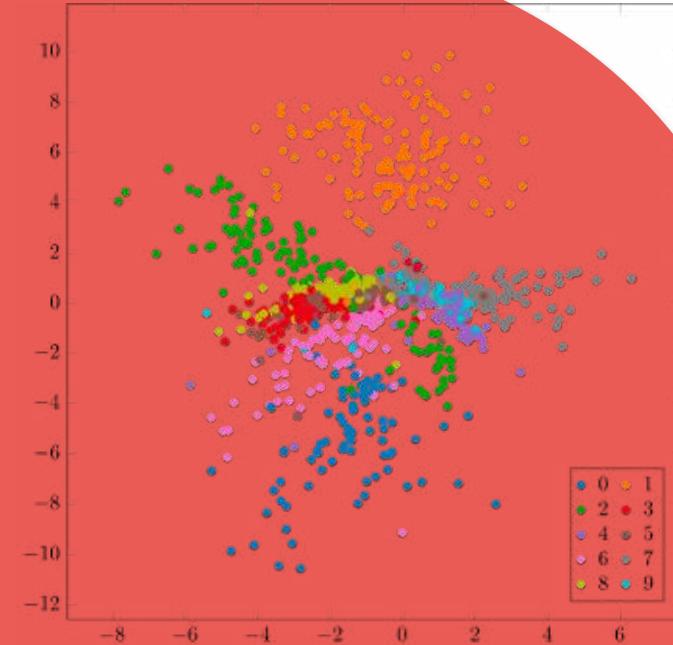
MODELING OF  $P(z|x)$  WITH GMMs

# GMMs

- Both sampling and evaluating are straightforward with GMMs
- However, performance scales poorly with increasing dimensionality
- Depends on initial choices...

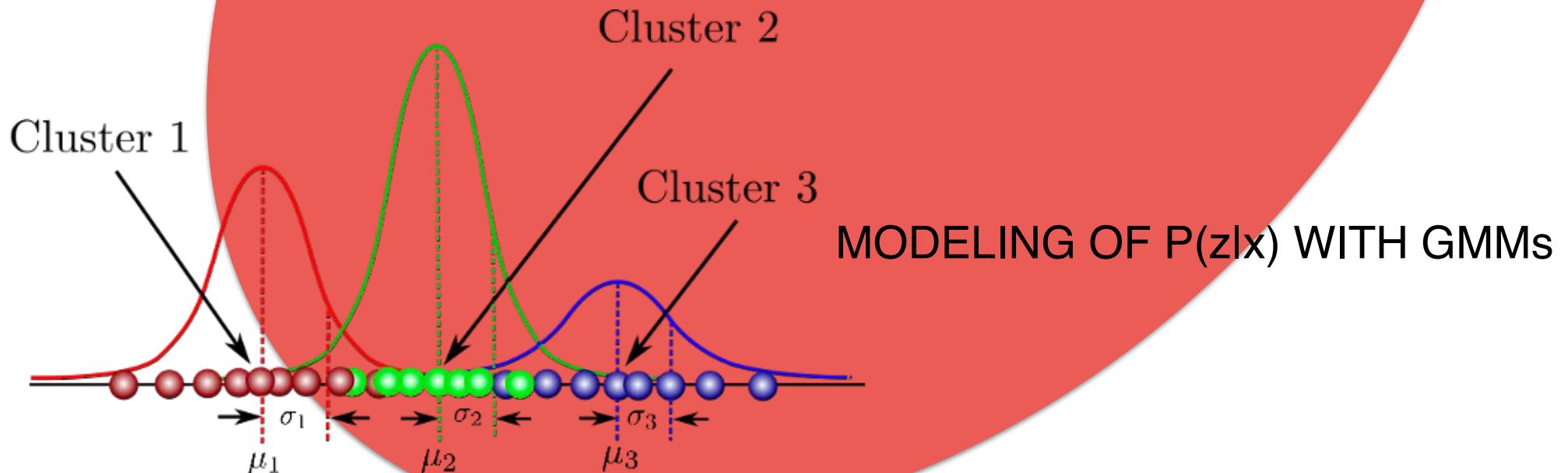
# DATA

|   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 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)



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 FLOWS, AUTOREGRESSIVE FLOWS)**

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

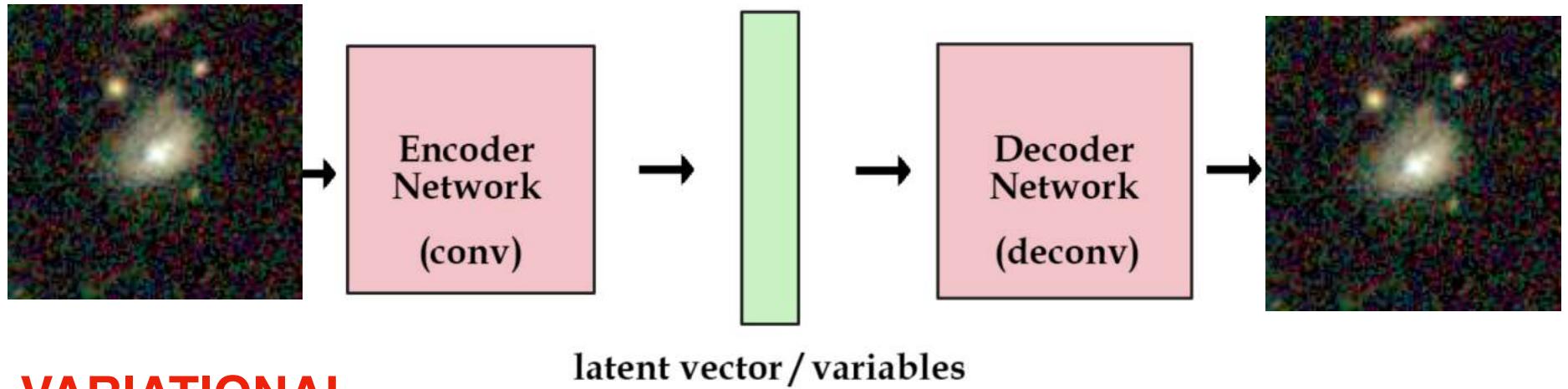
Much of the recent progress in unsupervised deep learning has been to invent network architectures that are capable of solving either or both of these related problems directly, without resorting to any auxiliary methods

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

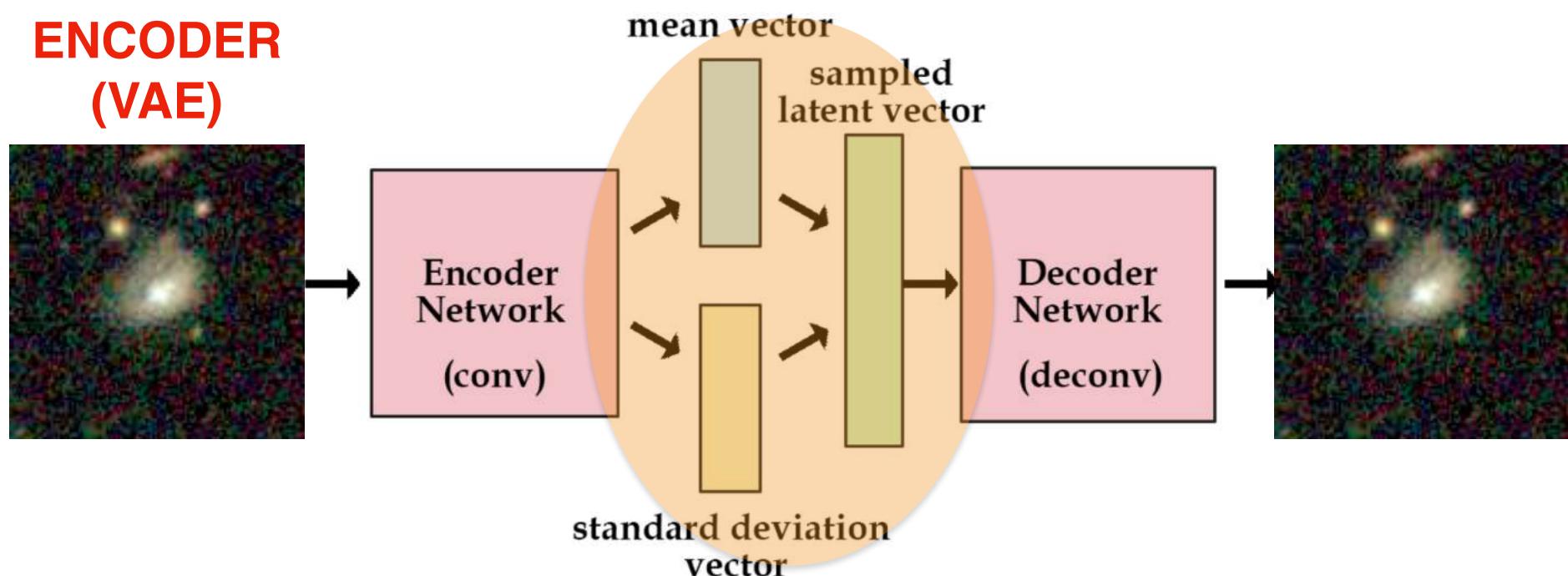
**GAN**  
**(GENERATIVE ADVERSARIAL NETWRK)**

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

# AUTO-ENCODER

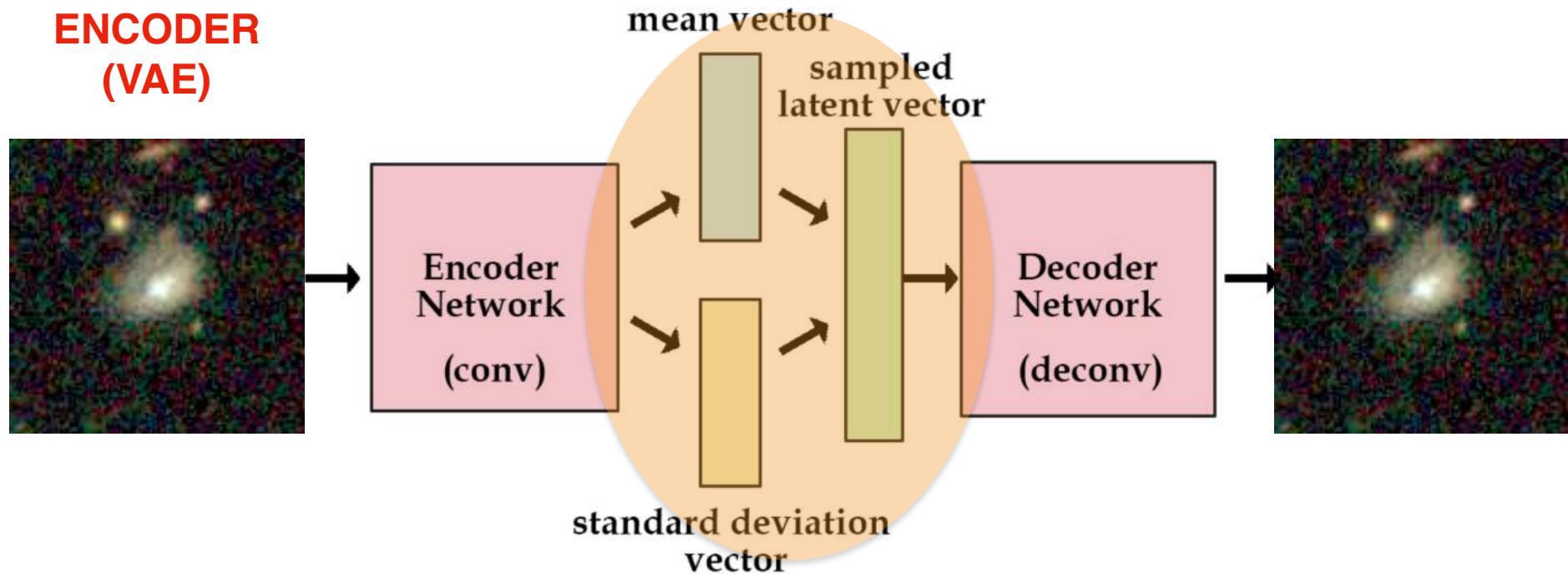


# VARIATIONAL AUTO-ENCODER (VAE)



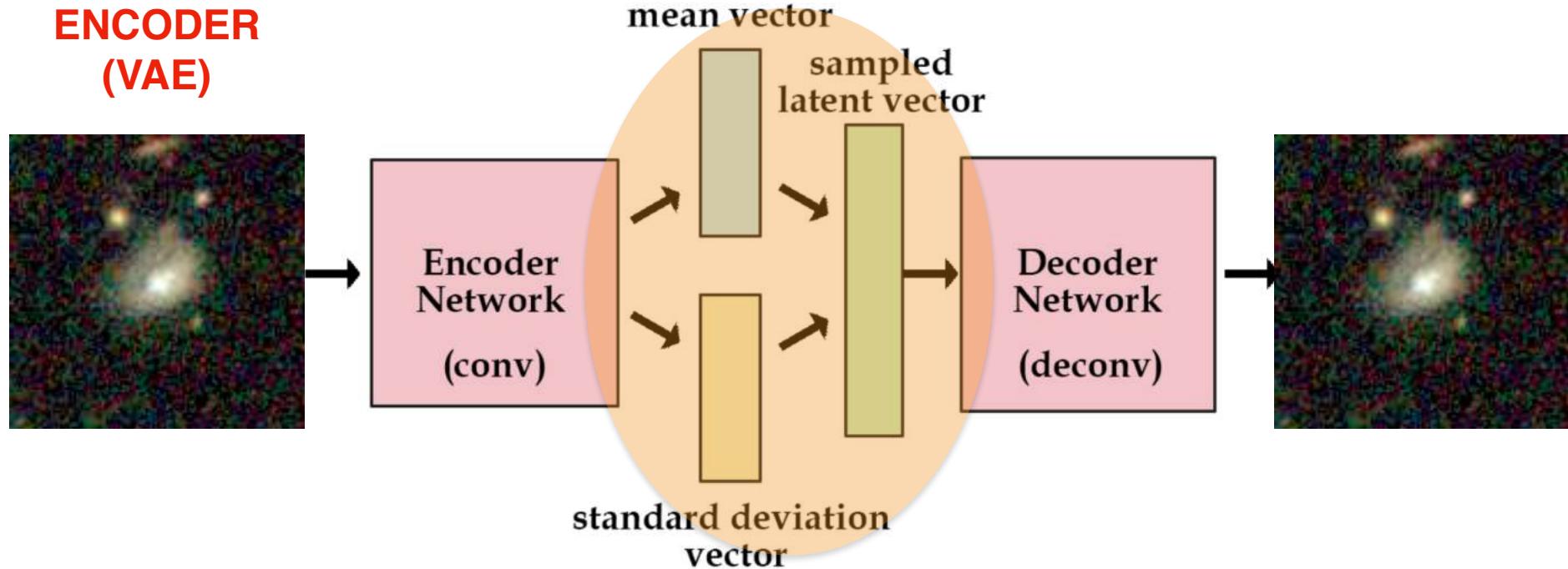
LET'S MODEL THE LATENT SPACE WITH A MIXTURE OF GAUSSIANS

## VARIATIONAL AUTO- ENCODER (VAE)



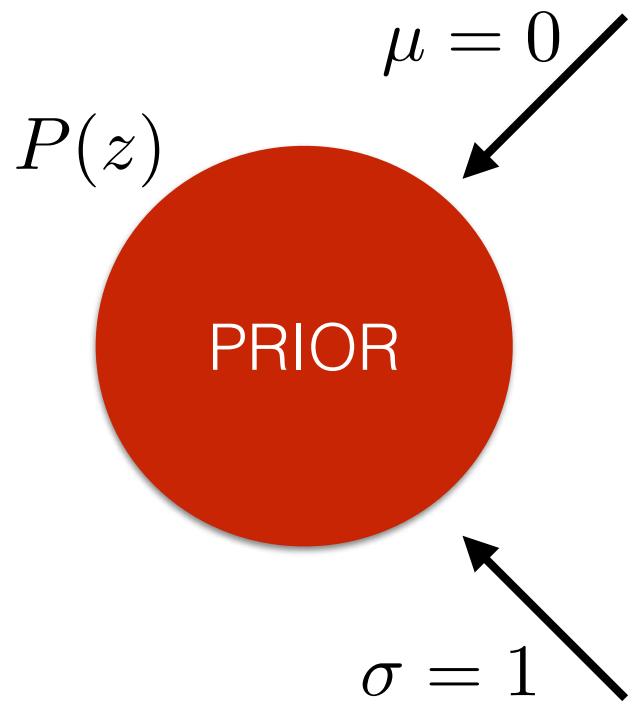
HOWEVER, NOTHING GUARANTEES US THAT  $P(z|x)$  CAN BE ACCURATELY MODELLED BY A MIXTURE OF GAUSSIANS....

## VARIATIONAL AUTO- ENCODER (VAE)

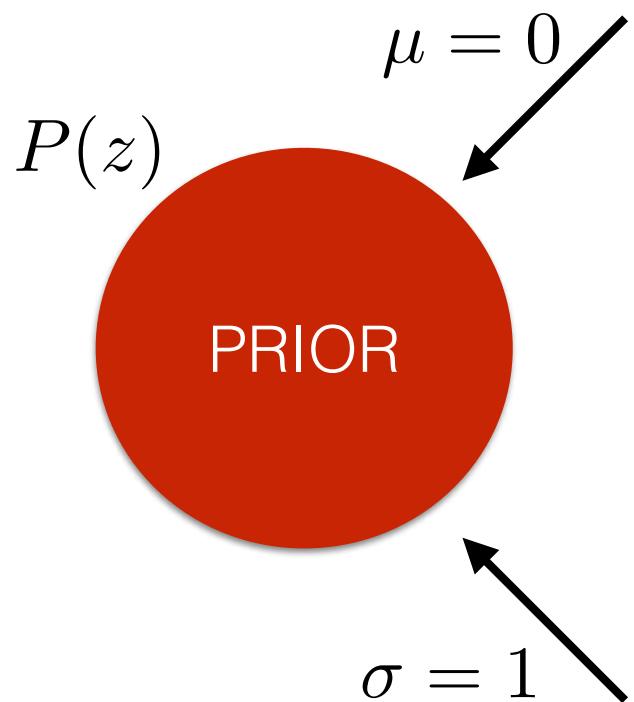


HOWEVER, NOTHING GUARANTEES US THAT THE LATENT  
SPACE CAN BE MODELLED BY A MIXTURE OF  
GAUSSIANS....

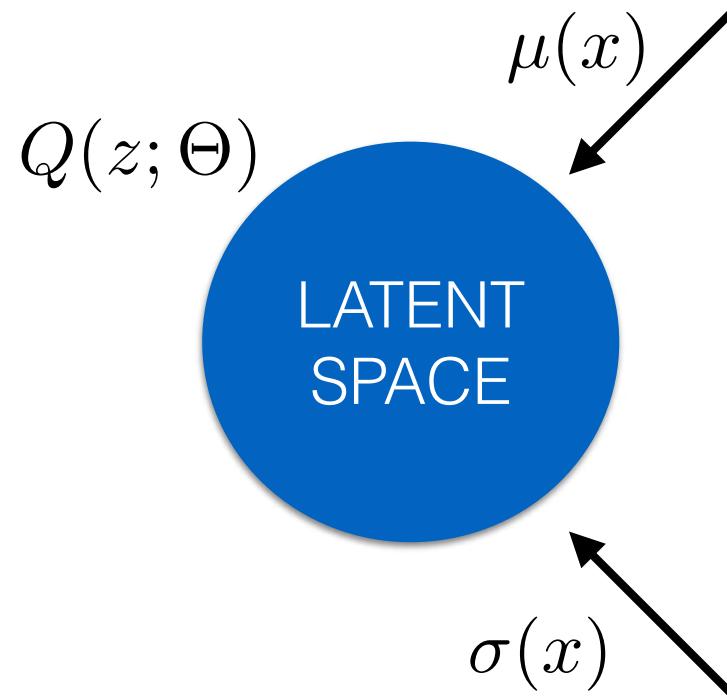
... LET'S FORCE IT TO BE GAUSSIAN LIKE!



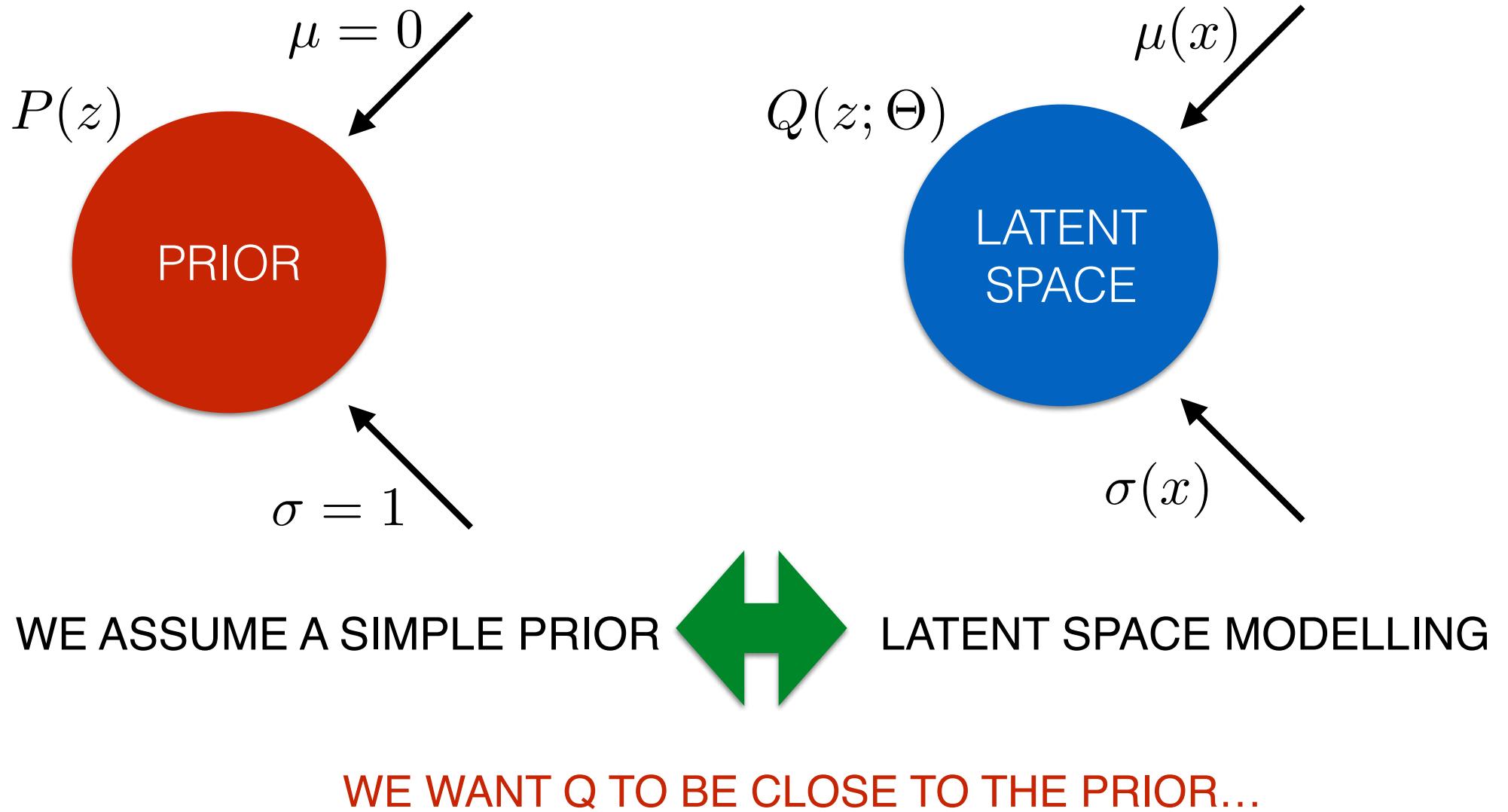
WE ASSUME A SIMPLE PRIOR

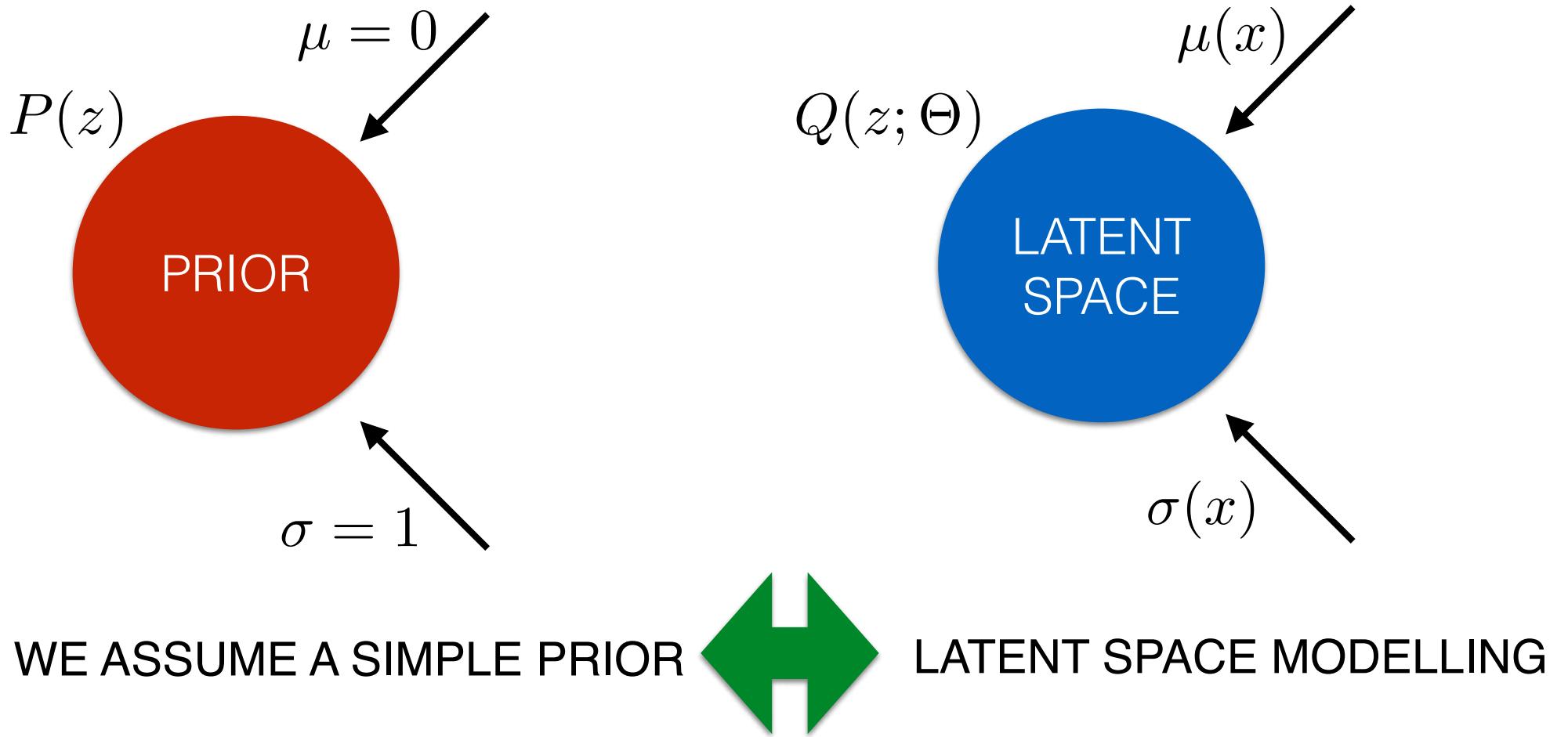


WE ASSUME A SIMPLE PRIOR



LATENT SPACE MODELING





WE 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...

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

L2 LOSS

REGULARIZATION TERM  
~Gaussian

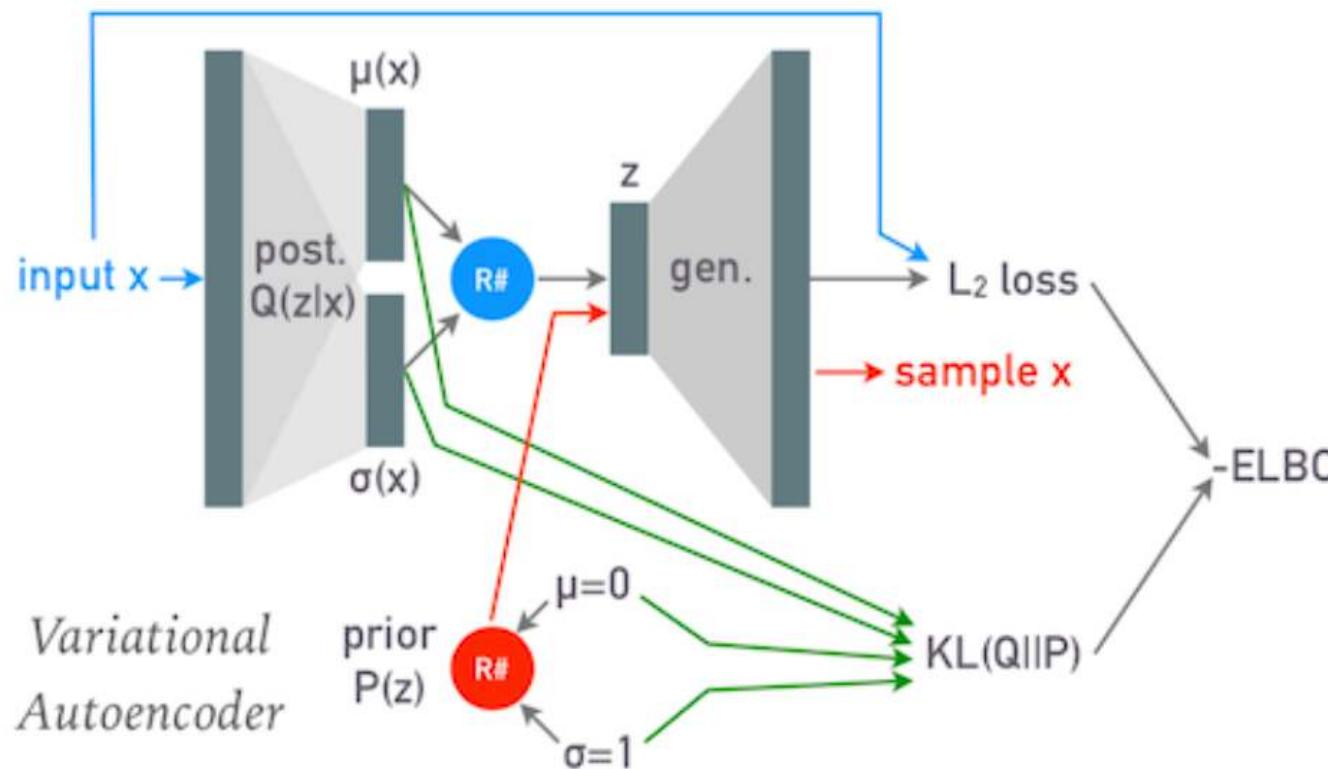
# WHAT WOULD BE THEN THE LOSS FUNCTION OF A VAE?

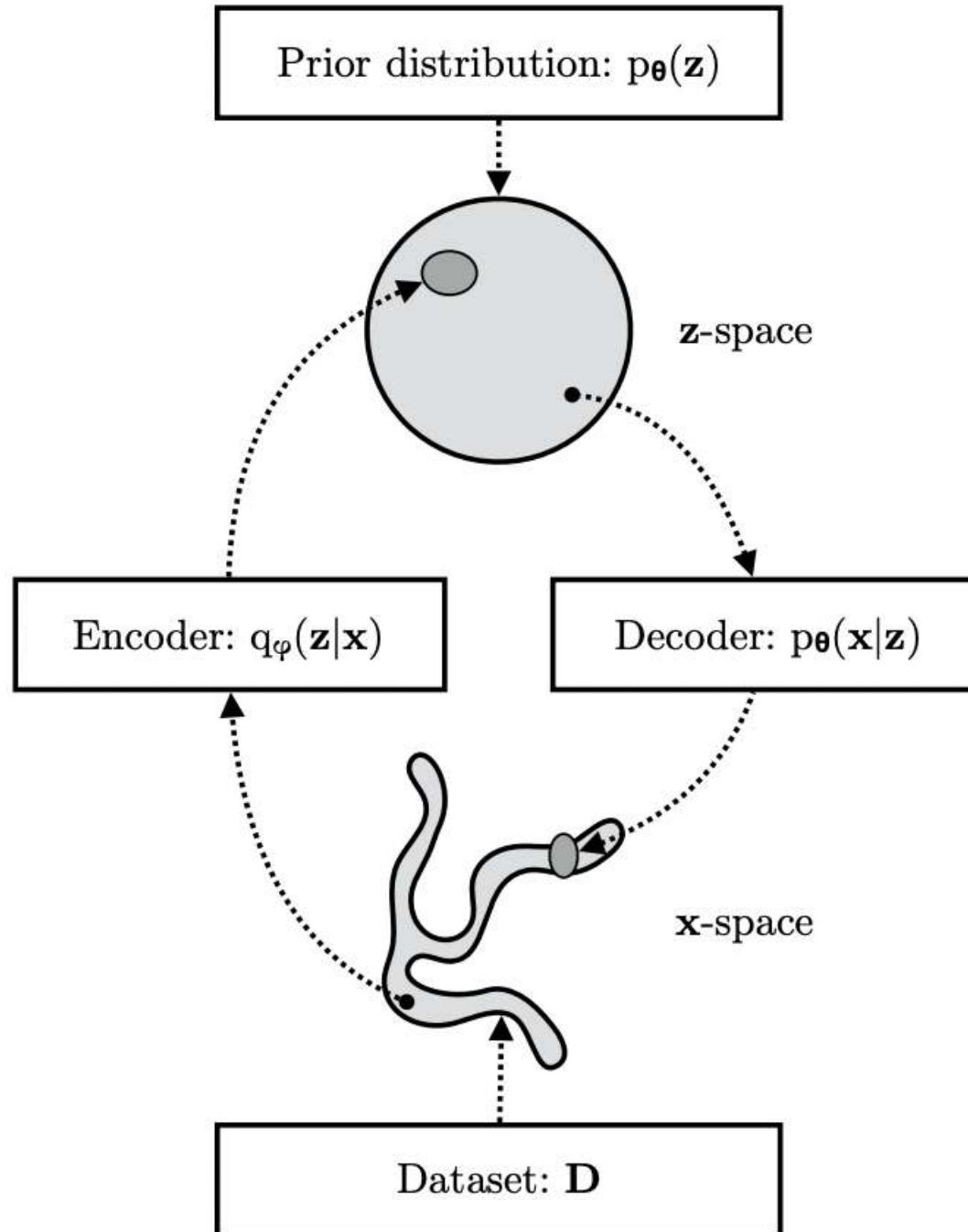
The key insight of VAE is that we are actually performing **variational inference** here, which then tells us what the loss function should be...

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

L2 LOSS

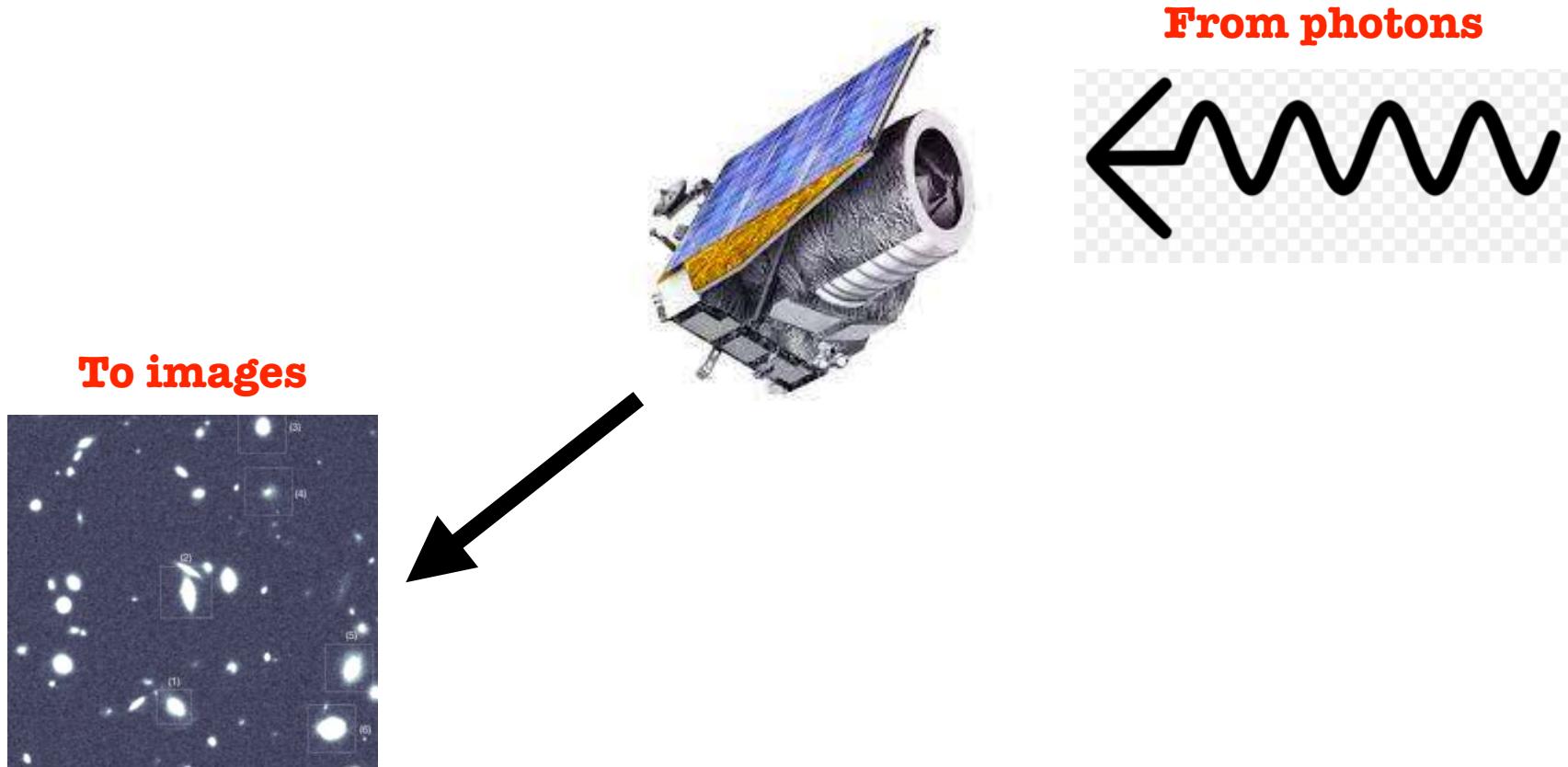
REGULARIZATION TERM  
~Gaussian





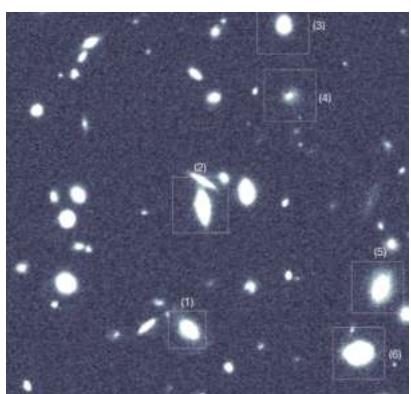
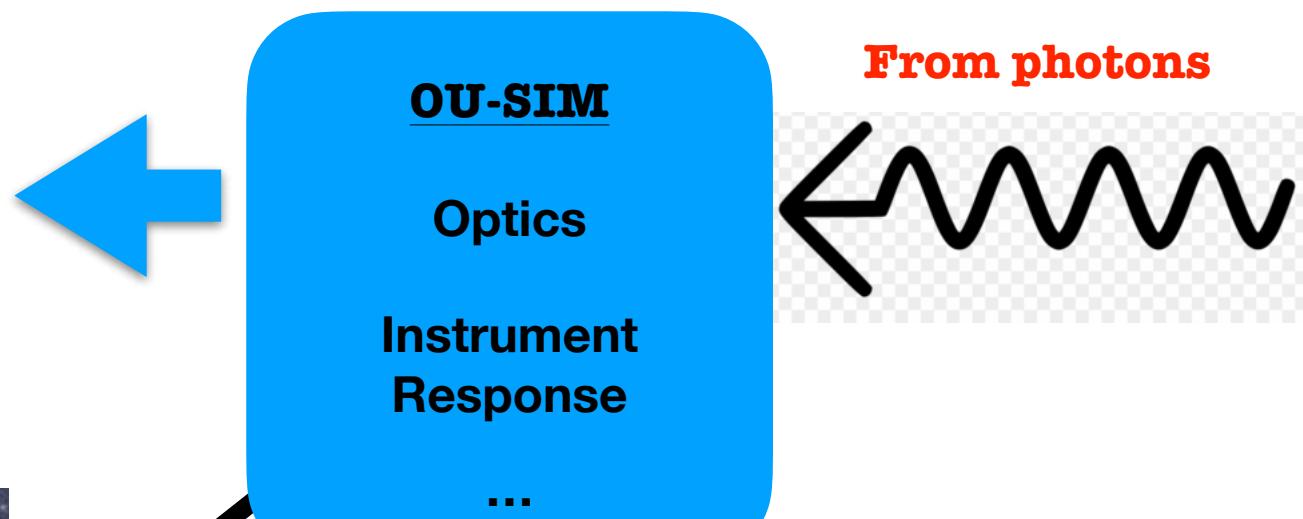
[Kingma and Welling, 2019]

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



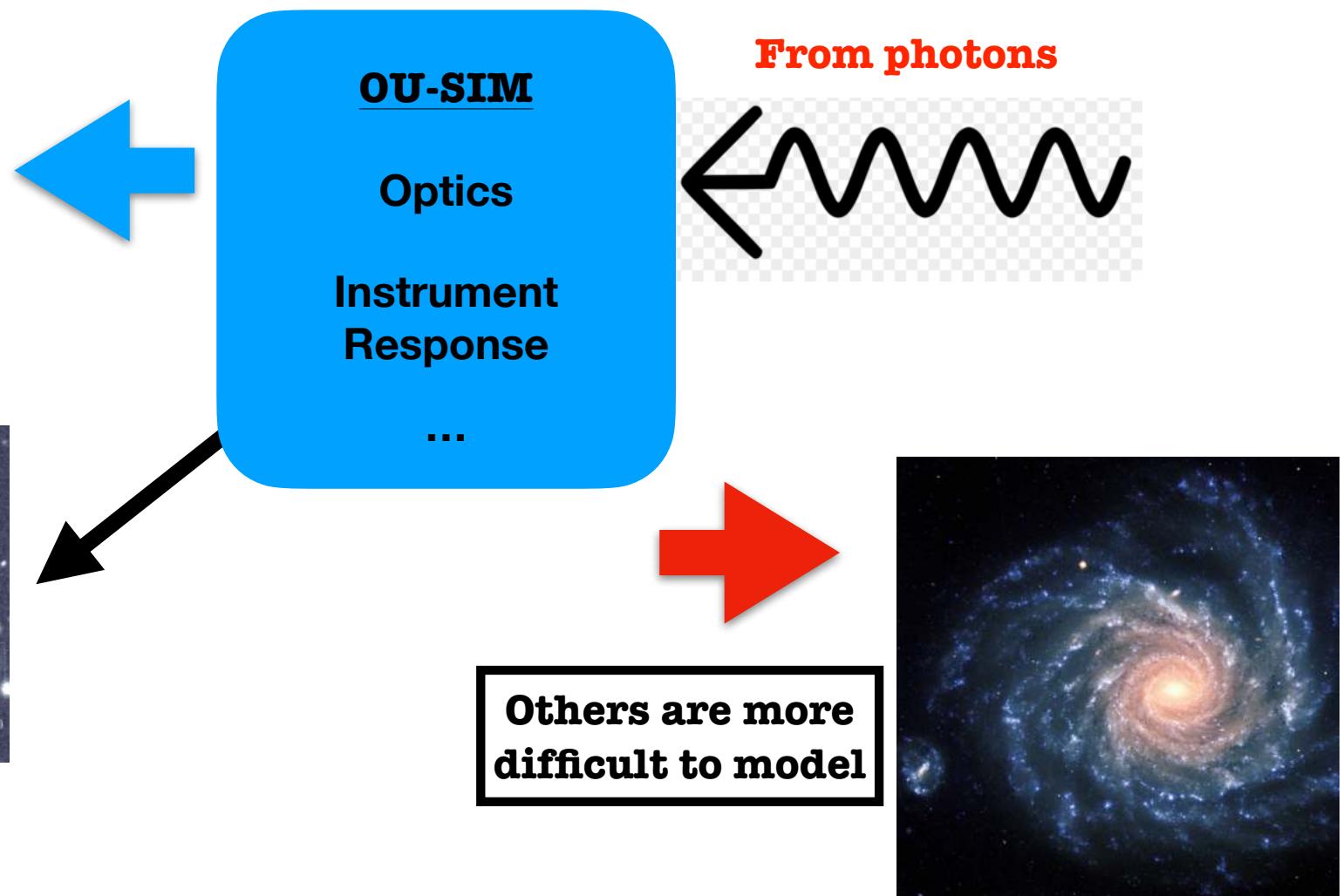
# Preparing a space mission requires very detailed simulations

There is a physical model for many of these effects



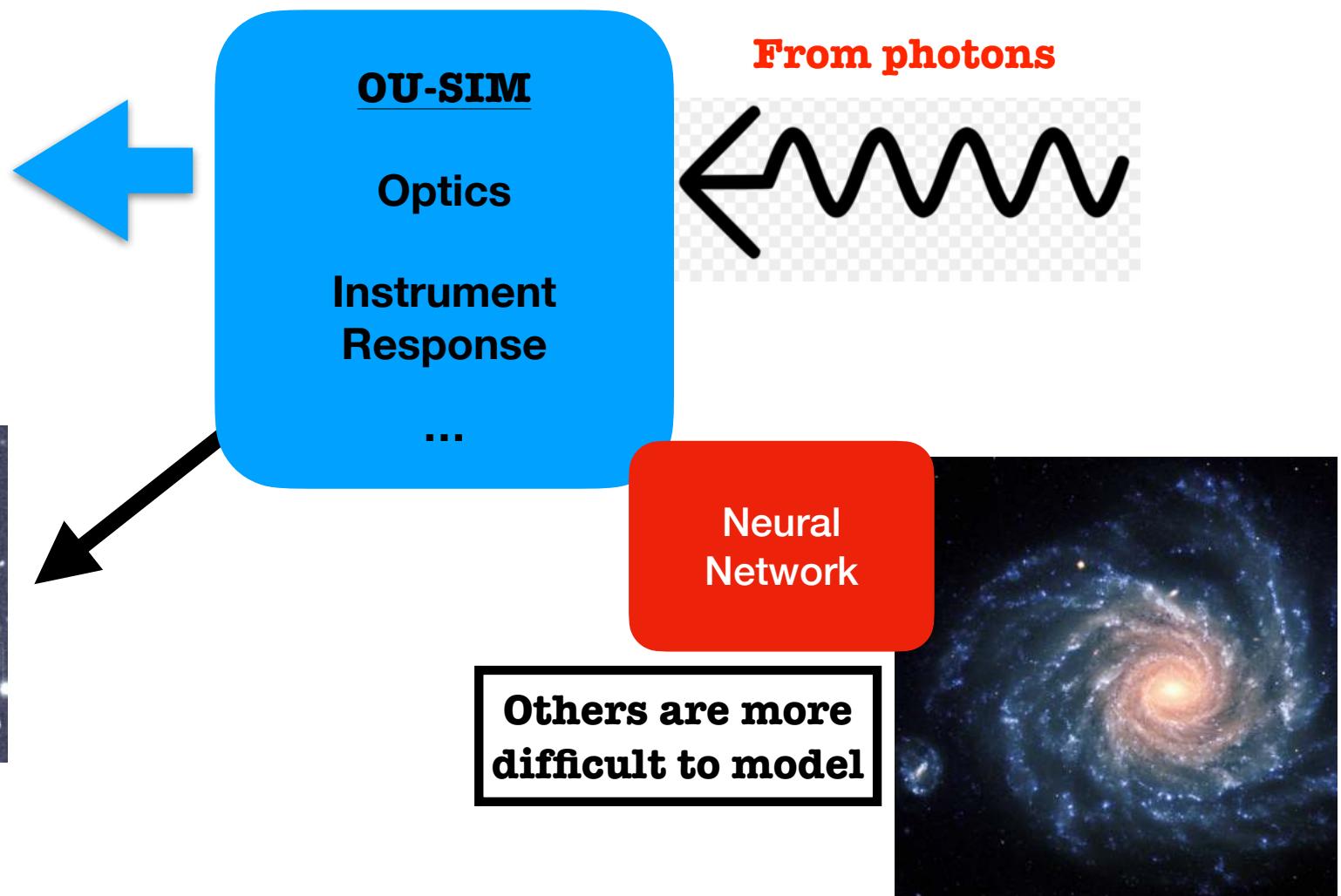
# **Preparing a space mission requires very detailed simulations**

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

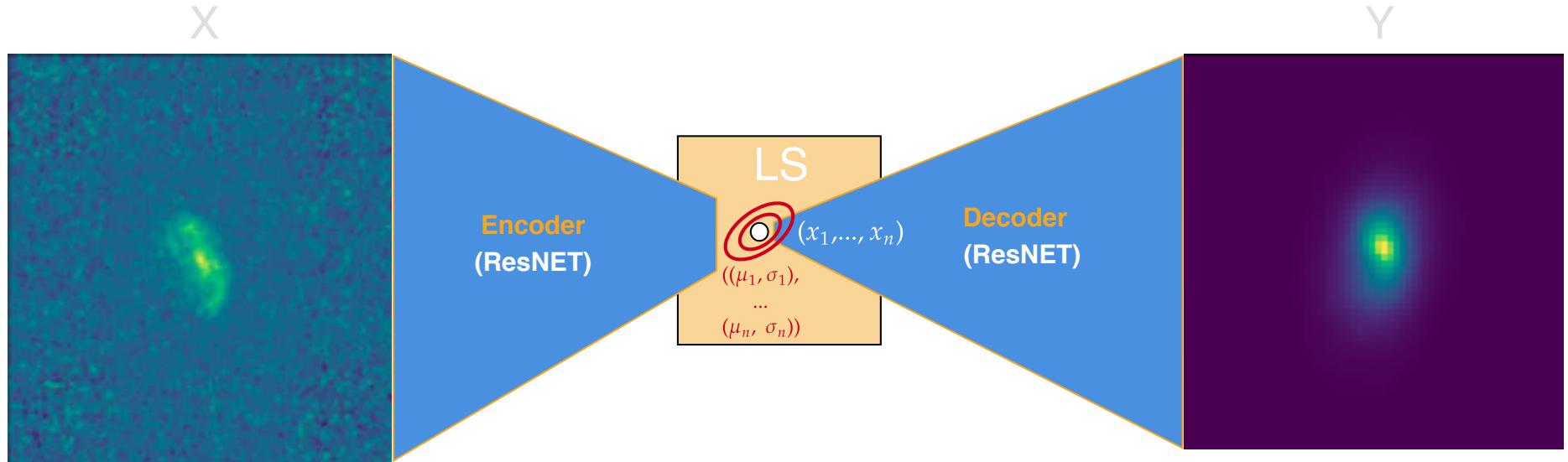


# **Preparing a space mission requires very detailed simulations**

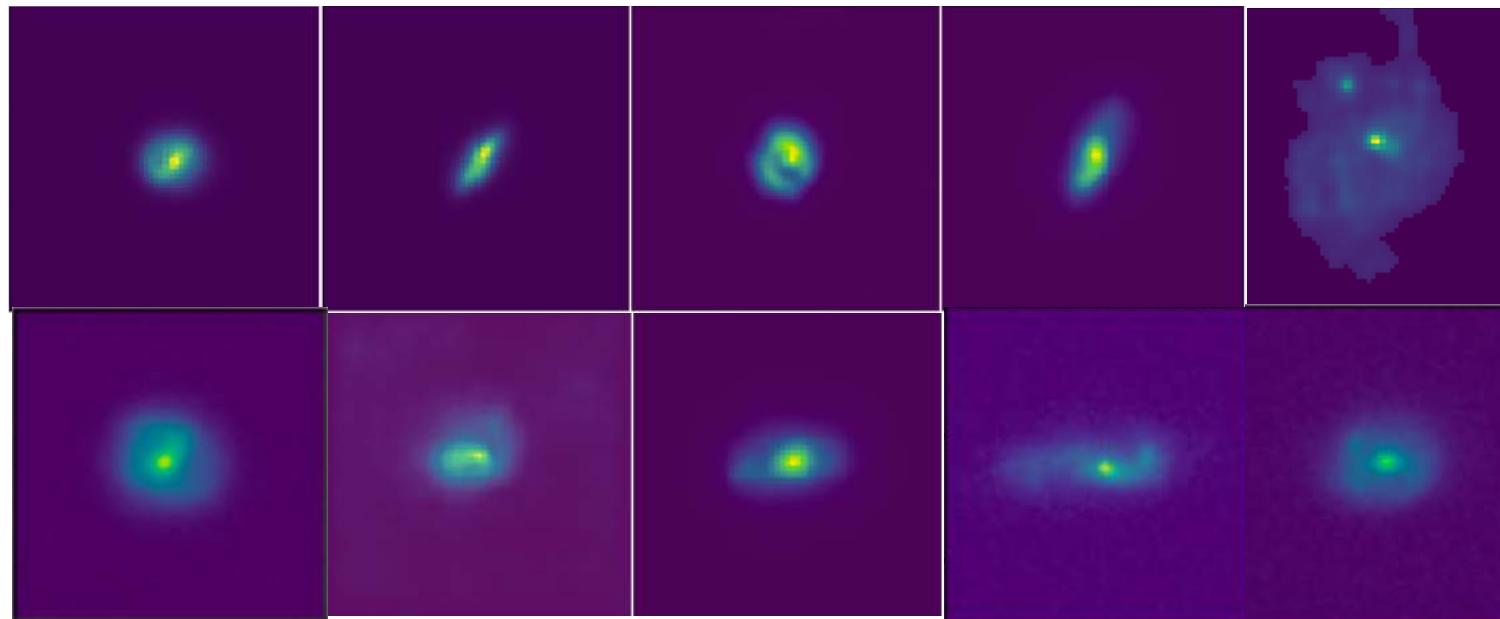
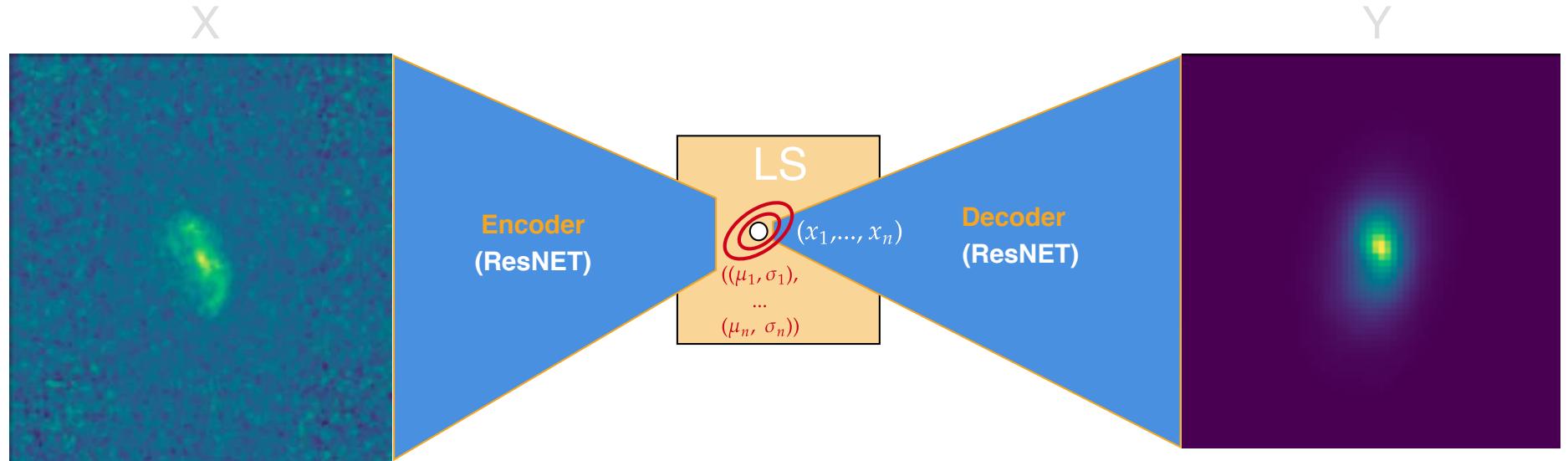
**There is a physical model for many of these effects**

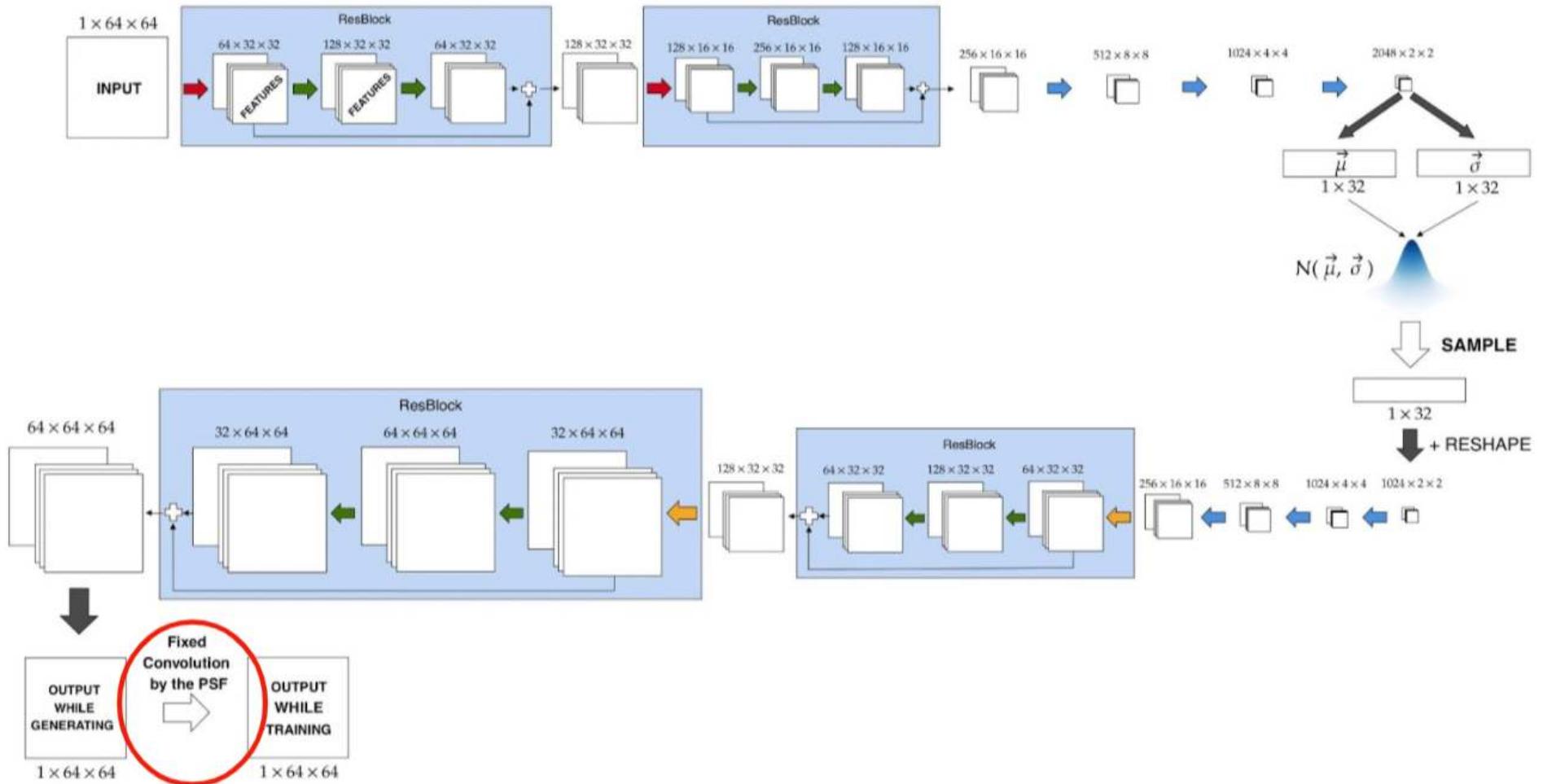
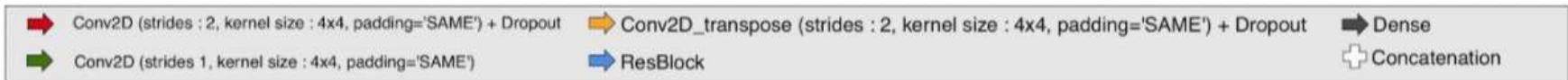


# Include a Generative Model for galaxy generation in the Euclid Simulation Pipeline to model process for which we do not have a physical model



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





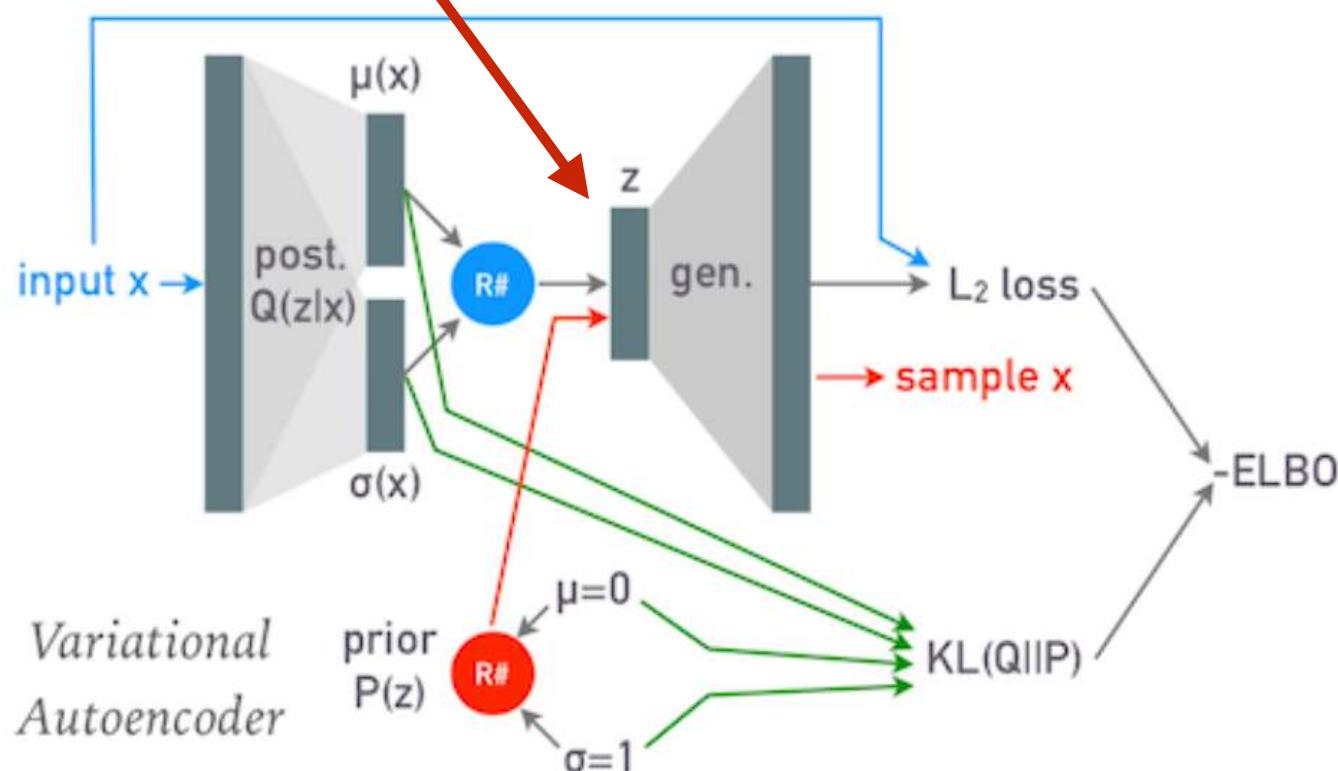
**Physical information is included as well within the architecture**

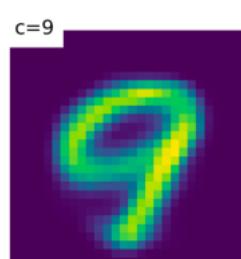
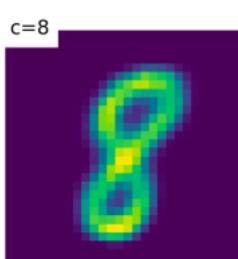
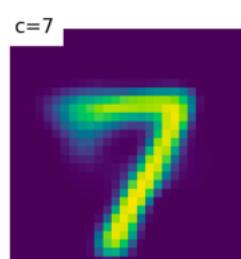
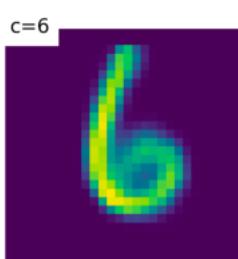
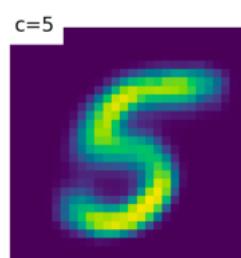
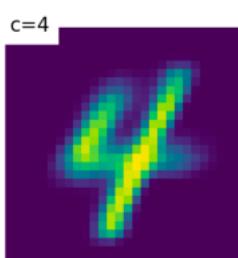
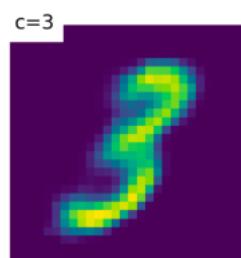
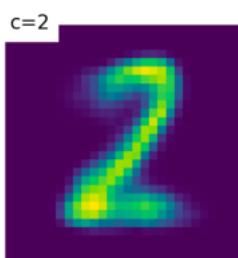
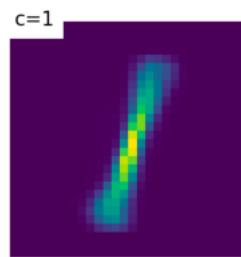
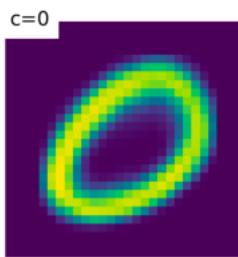
VAE's can also be conditioned, to generate a given class

concatenate labels

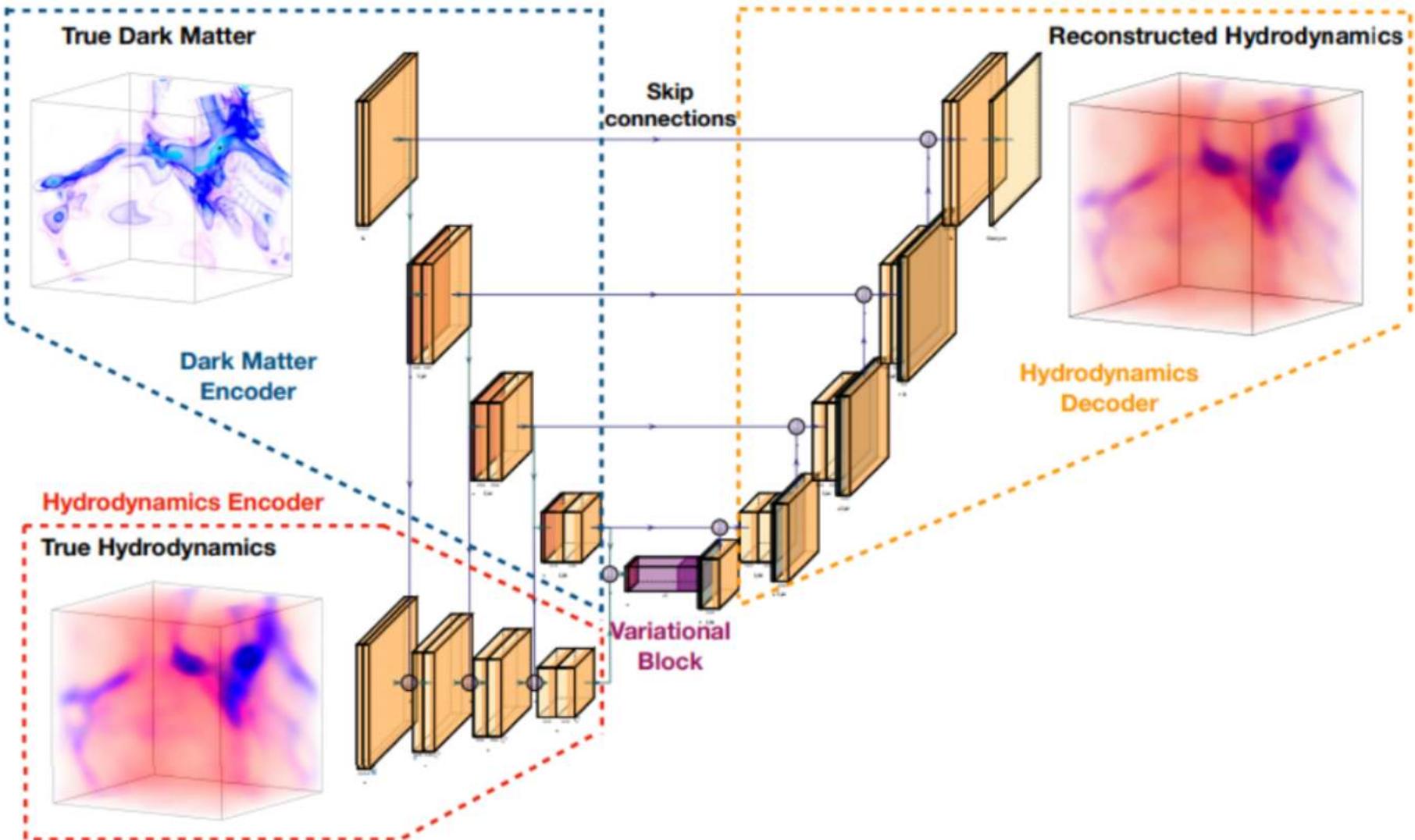


$$P(X|Y)$$





# Painting Baryons



Rodriguez+19, Modi+18, Berger+18, He+18, Zhang+19, Troster+19, Zamudio-Fernandez+19, Perraudin+19, Charnock+19, List+19, Giusarma+19, Bernardini+19, Chardin+19, Mustafa+19, Ramanah+20, Tamasiunas+20, Feder+20, Moster+20, Thiele+20, Wadekar+20, Dai+20, Li+20, Lucie-Smith+20, Kasmanoff+20, Ni+21, Rouhainen+21, Harrington+21, Horowitz+21, Horowitz+21, Bernardini+21, Schaurecker+21, Etezad-Razavi+21, Curtis+21

Horowitz+21

# Generative Adversarial Networks

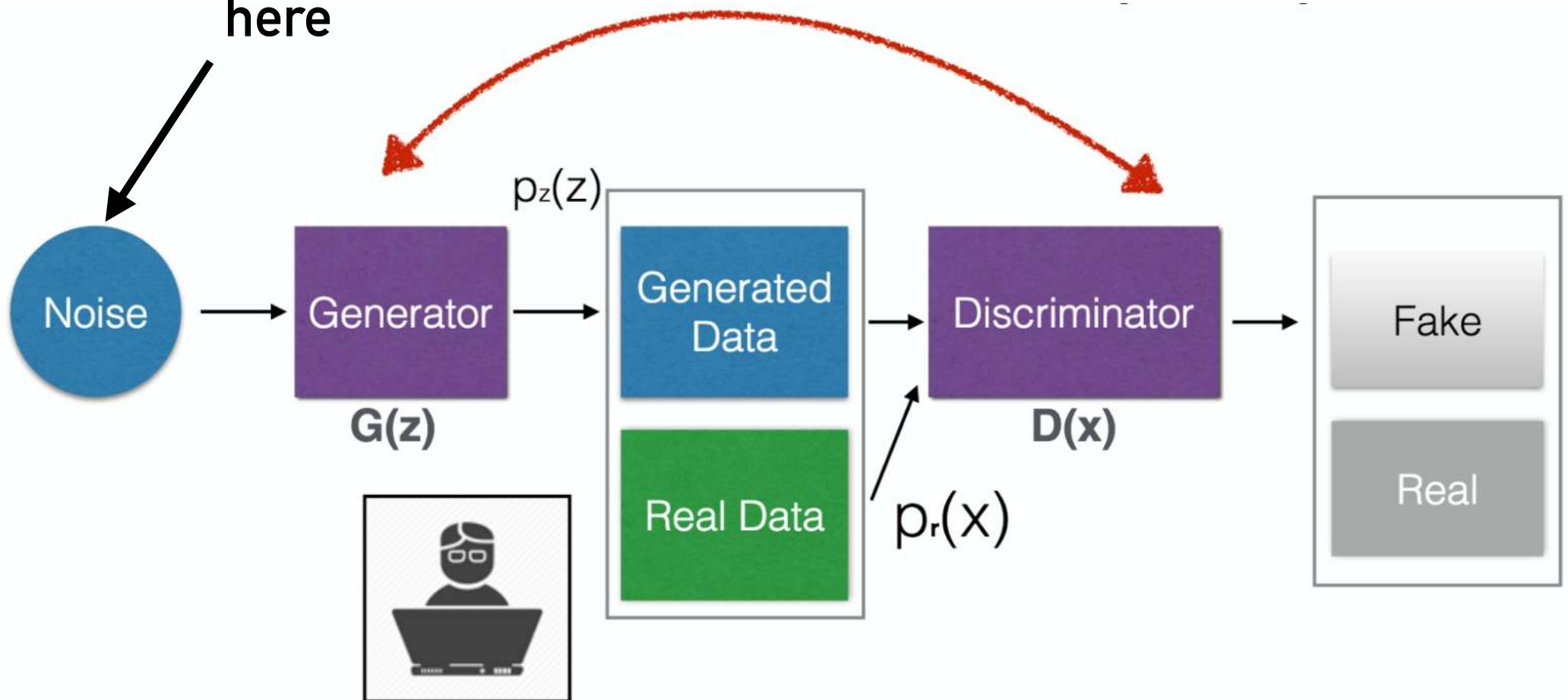
As for VAEs,  
the goal of generative Adversarial Networks (GANs) is to estimate  
 $p(z|x)$

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

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+14)

The latent variable is  
here

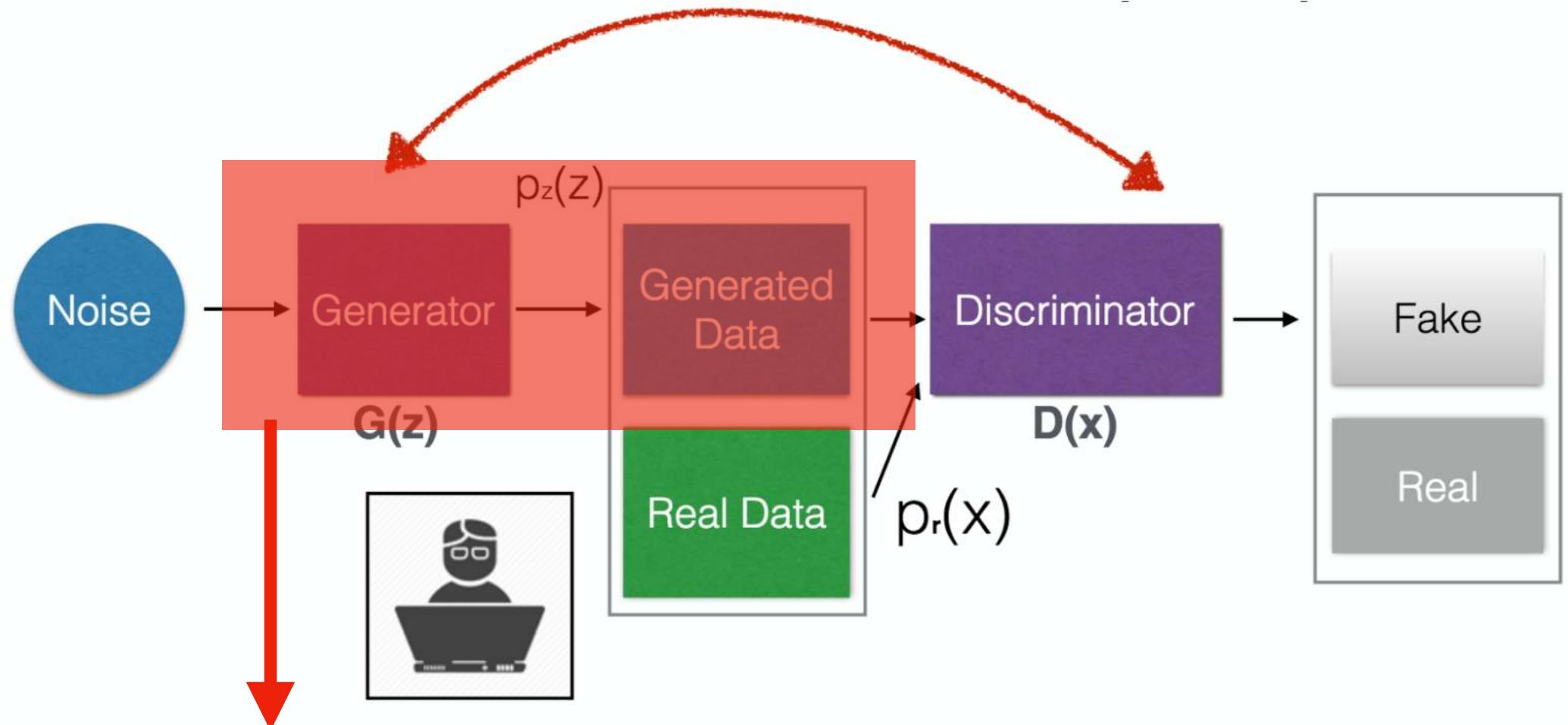


TWO COMPETING NETWORKS

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+)

## TWO COMPETING NETWORKS

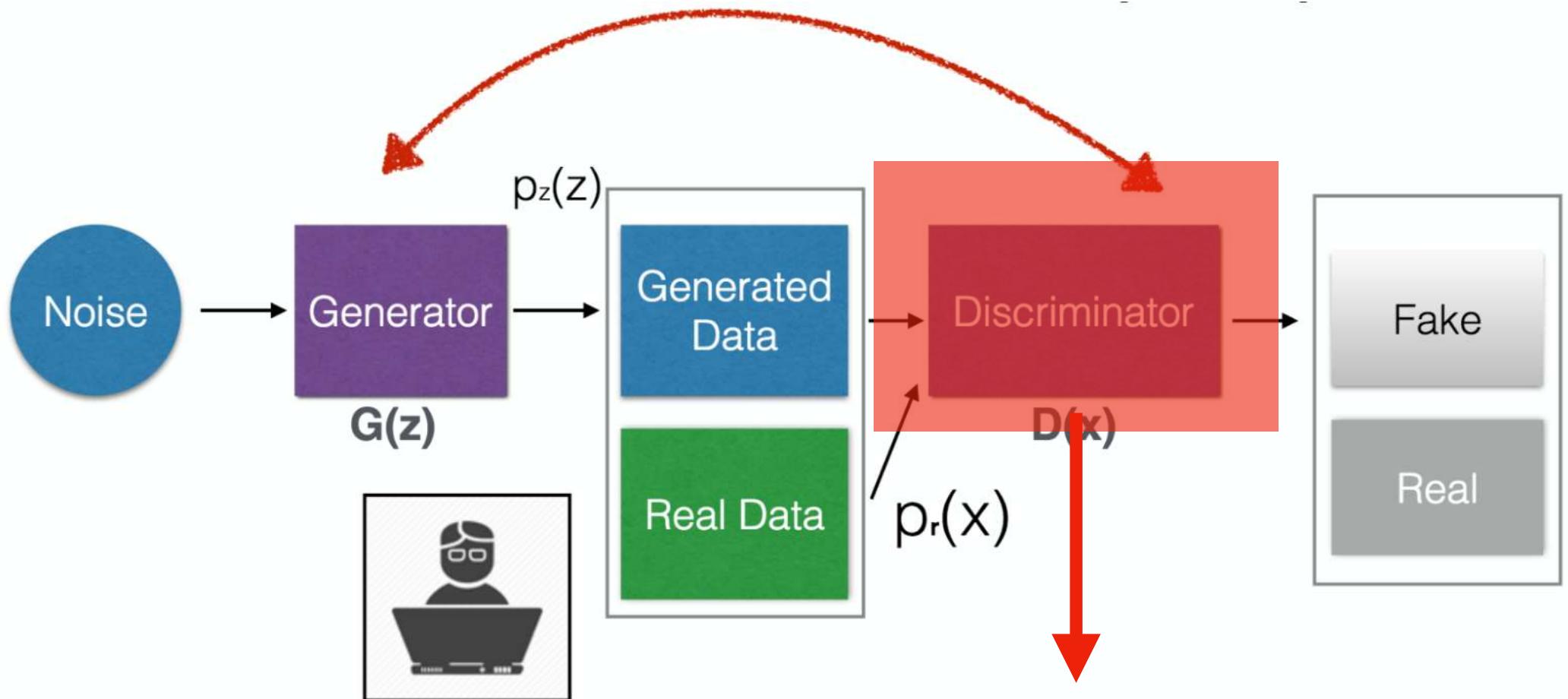


Every N iterations the generator  
is trained to force the discriminator  
to classify as real

# GENERATIVE ADVERSARIAL NETWORKS

(Goodfellow+)

## TWO COMPETING NETWORKS

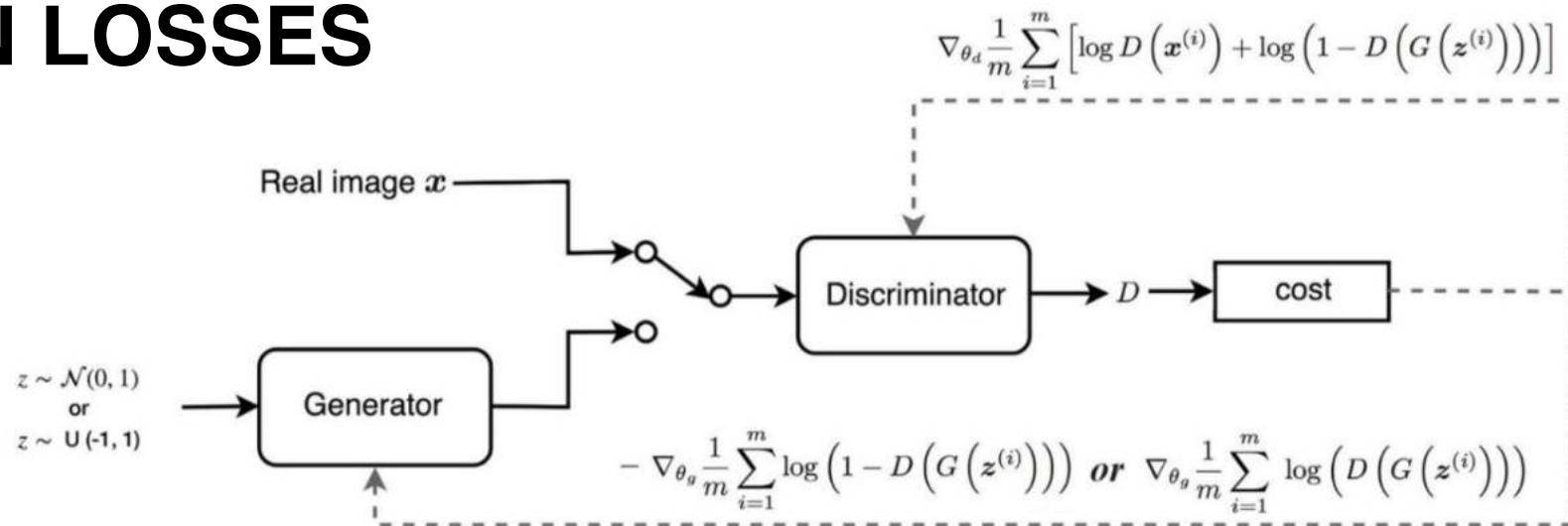


Every N iterations the discriminator  
is trained to force to distinguish between  
real and fake

## IN PRACTICE

### DISCRIMINATOR LOSS (CROSS-ENTROPY)

## GAN LOSSES



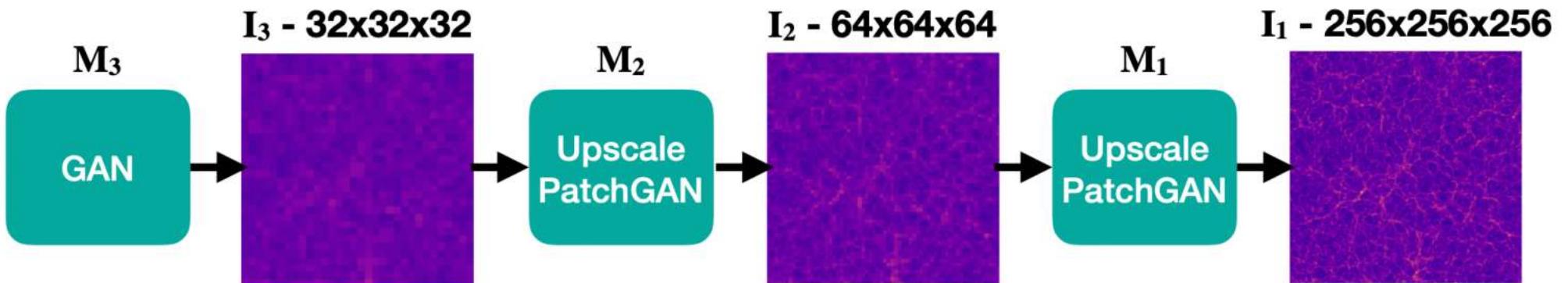
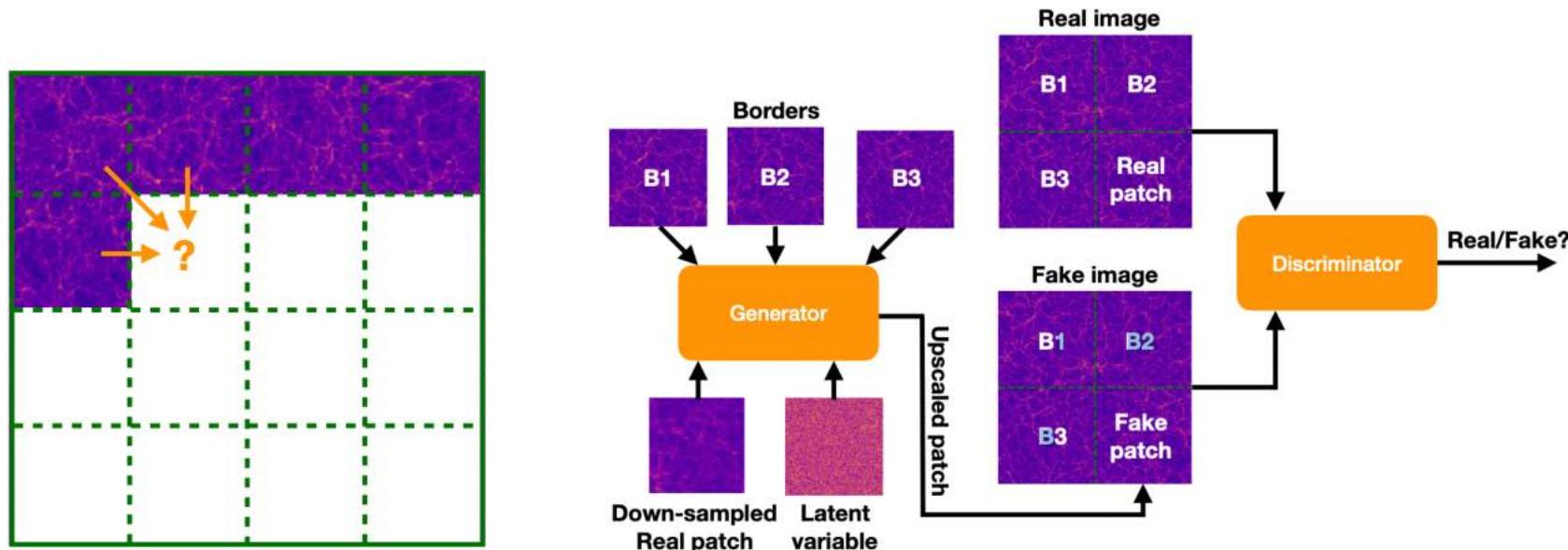
### GENERATOR LOSS

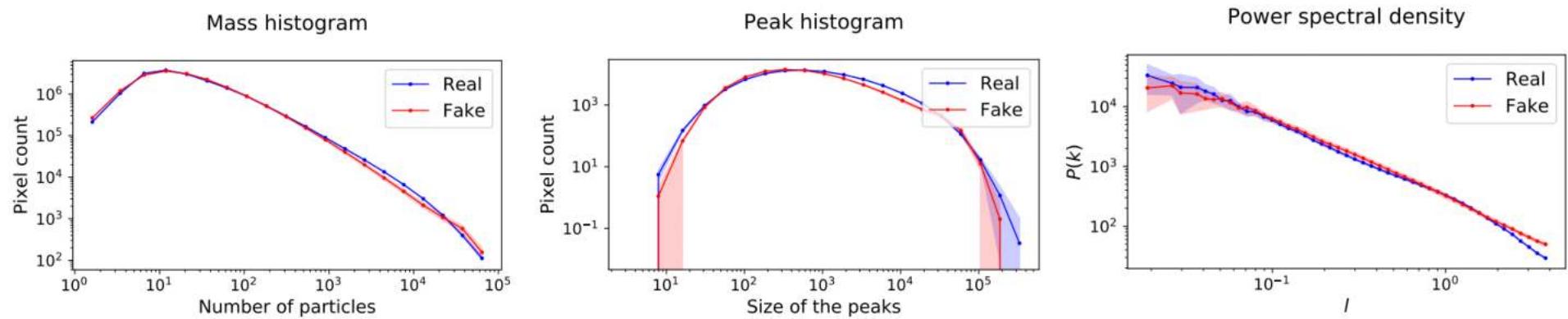
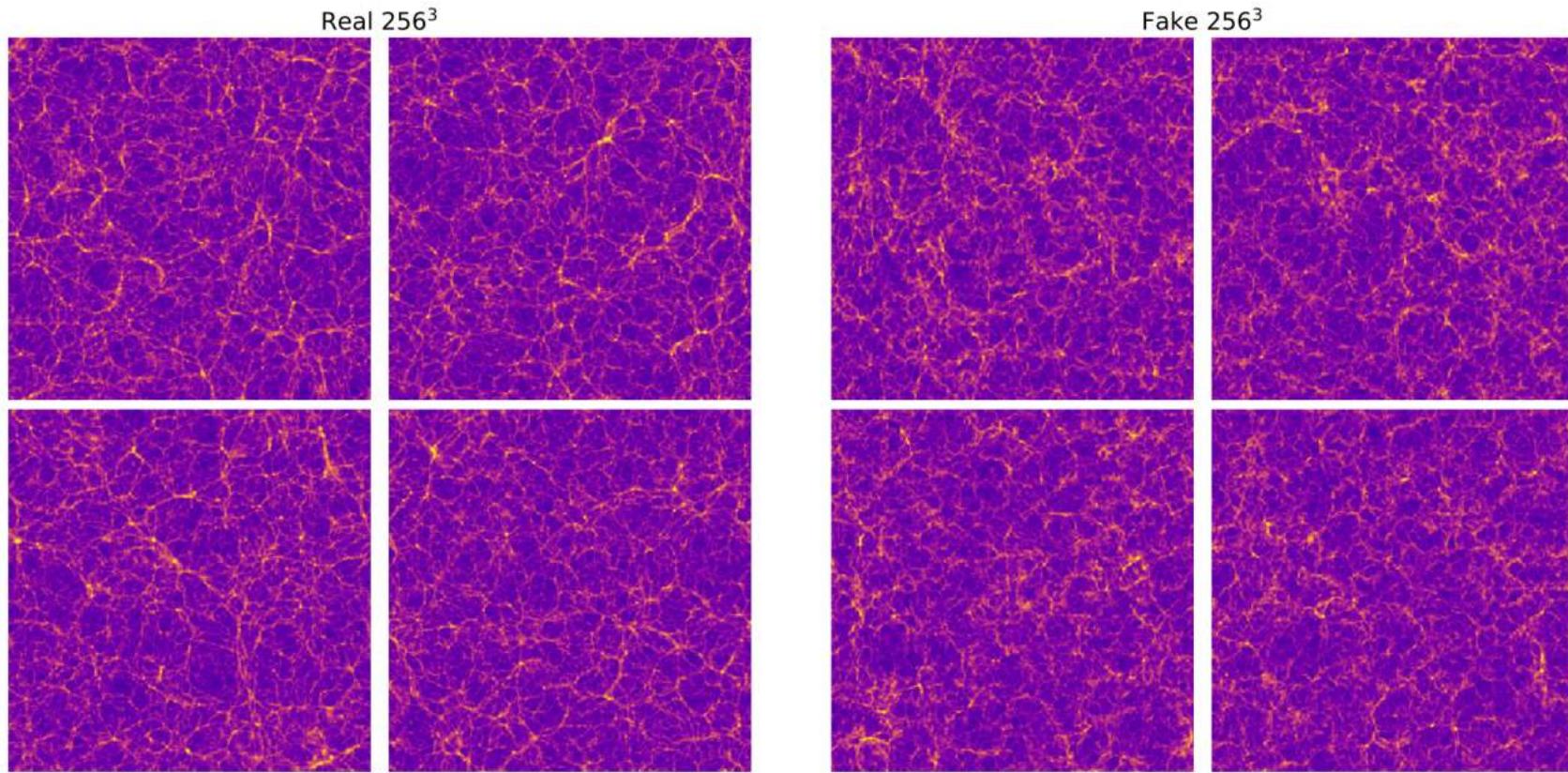
# GANs have remained state-of-the-art for data generation (until diffusion...)

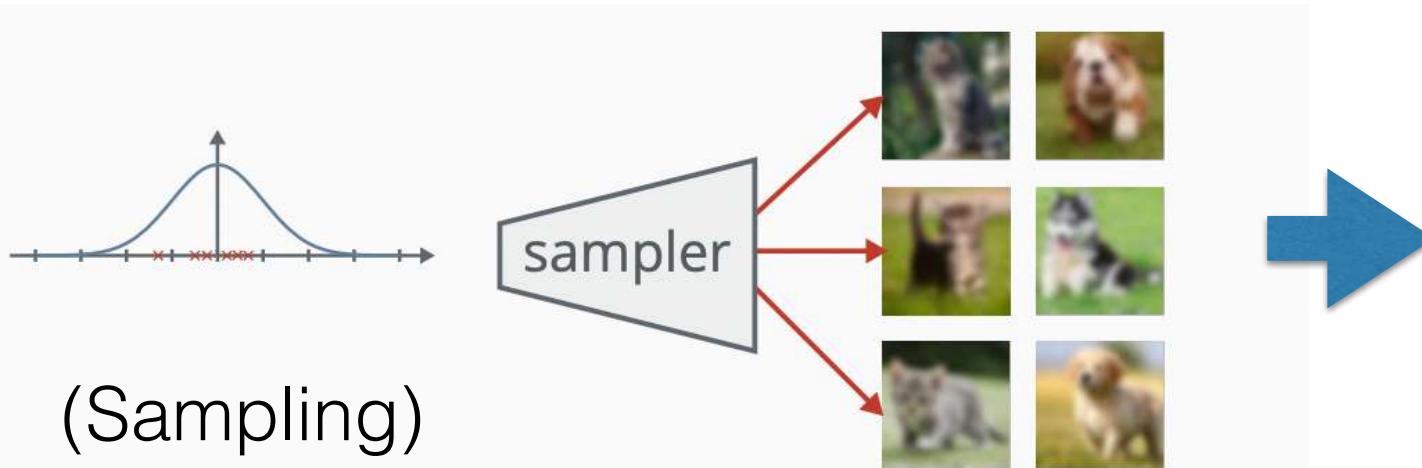


Karras+19

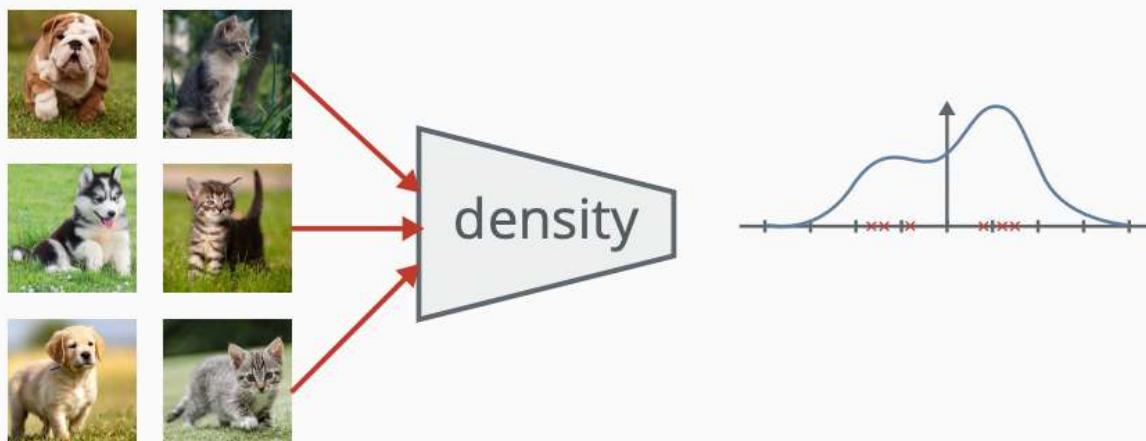
# N-body emulation by Deep Generative Modelling







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  $p$*

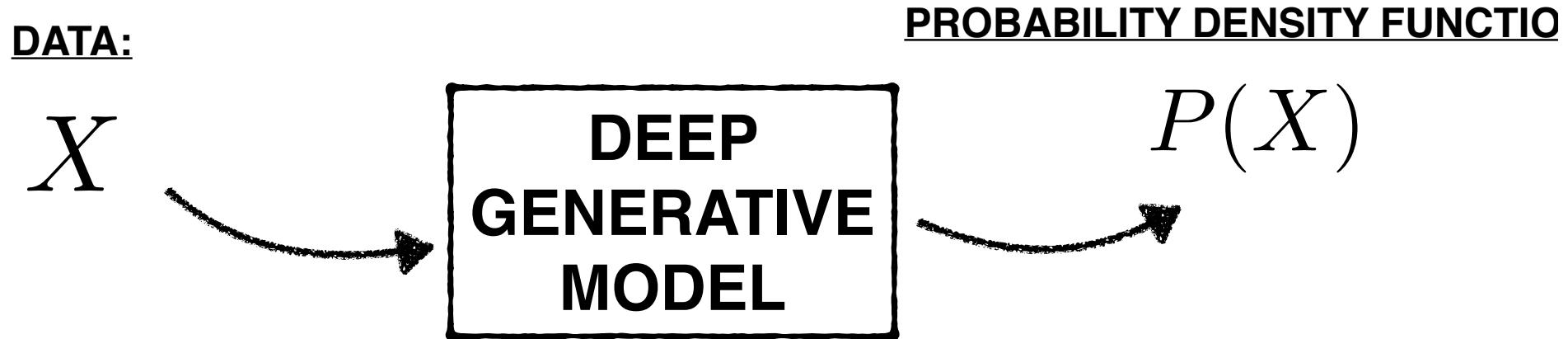


Then an object will be anomalous if:

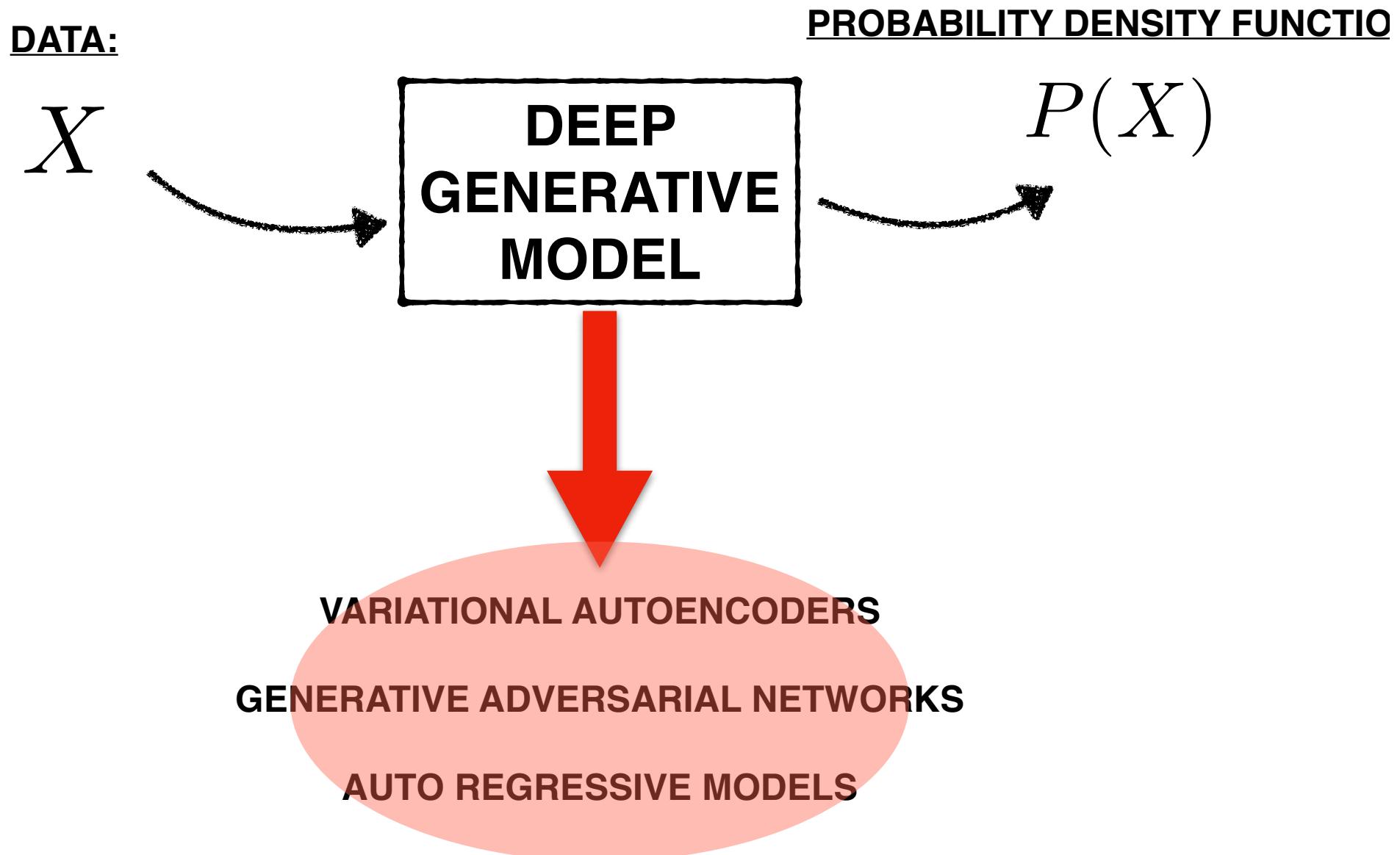
$$p(x_i) < \epsilon$$

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

## GENERATIVE MODELS DO PRECISELY THAT:



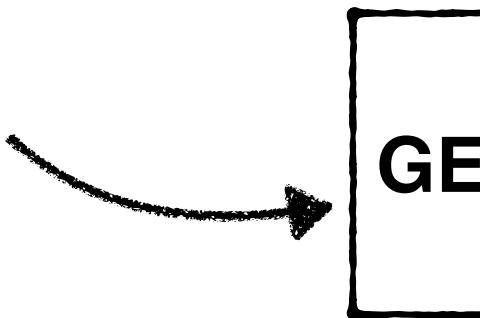
## GENERATIVE MODELS DO PRECISELY THAT:



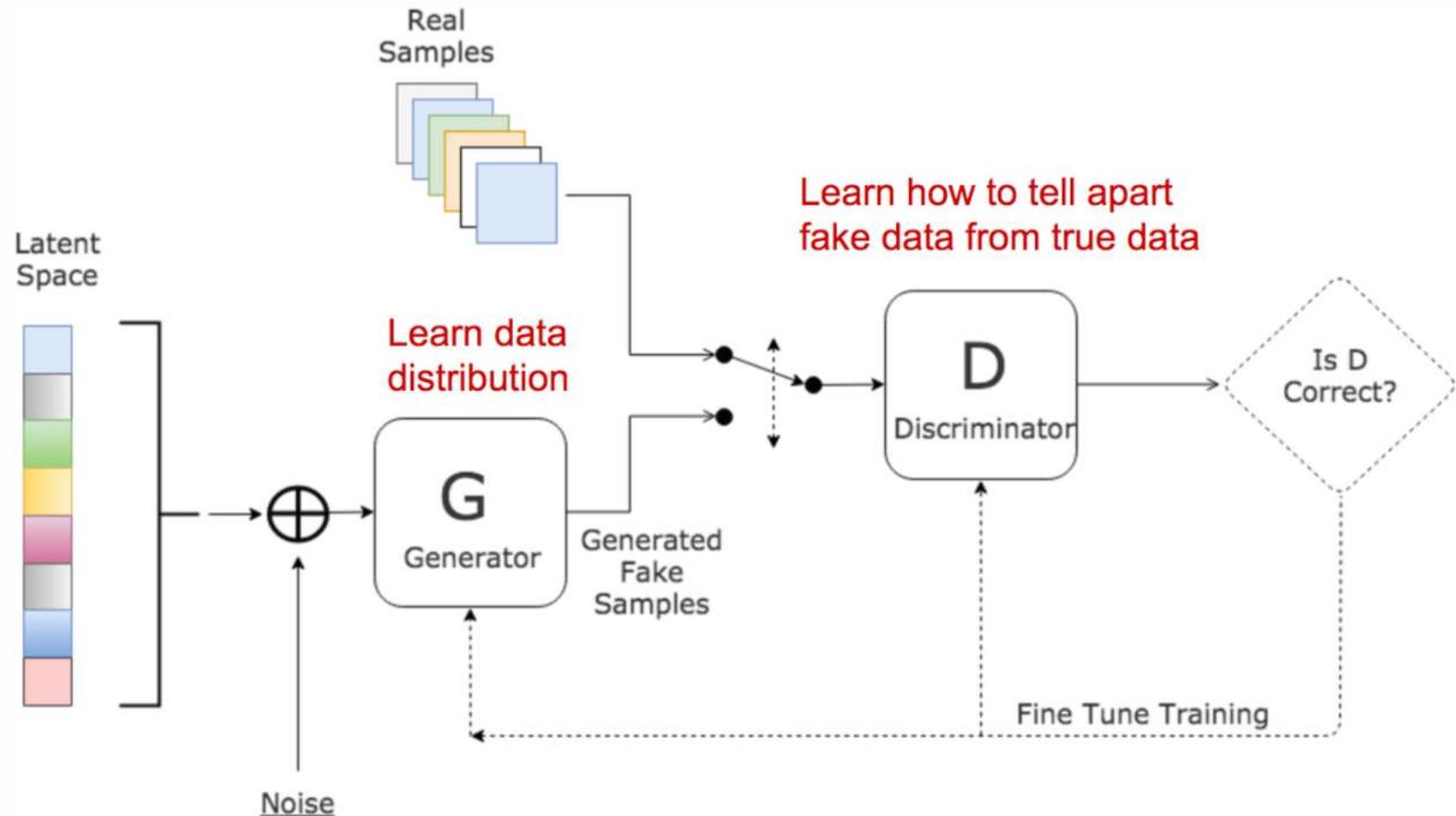
# GENERATIVE MODELS DO PRECISELY THAT:

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**

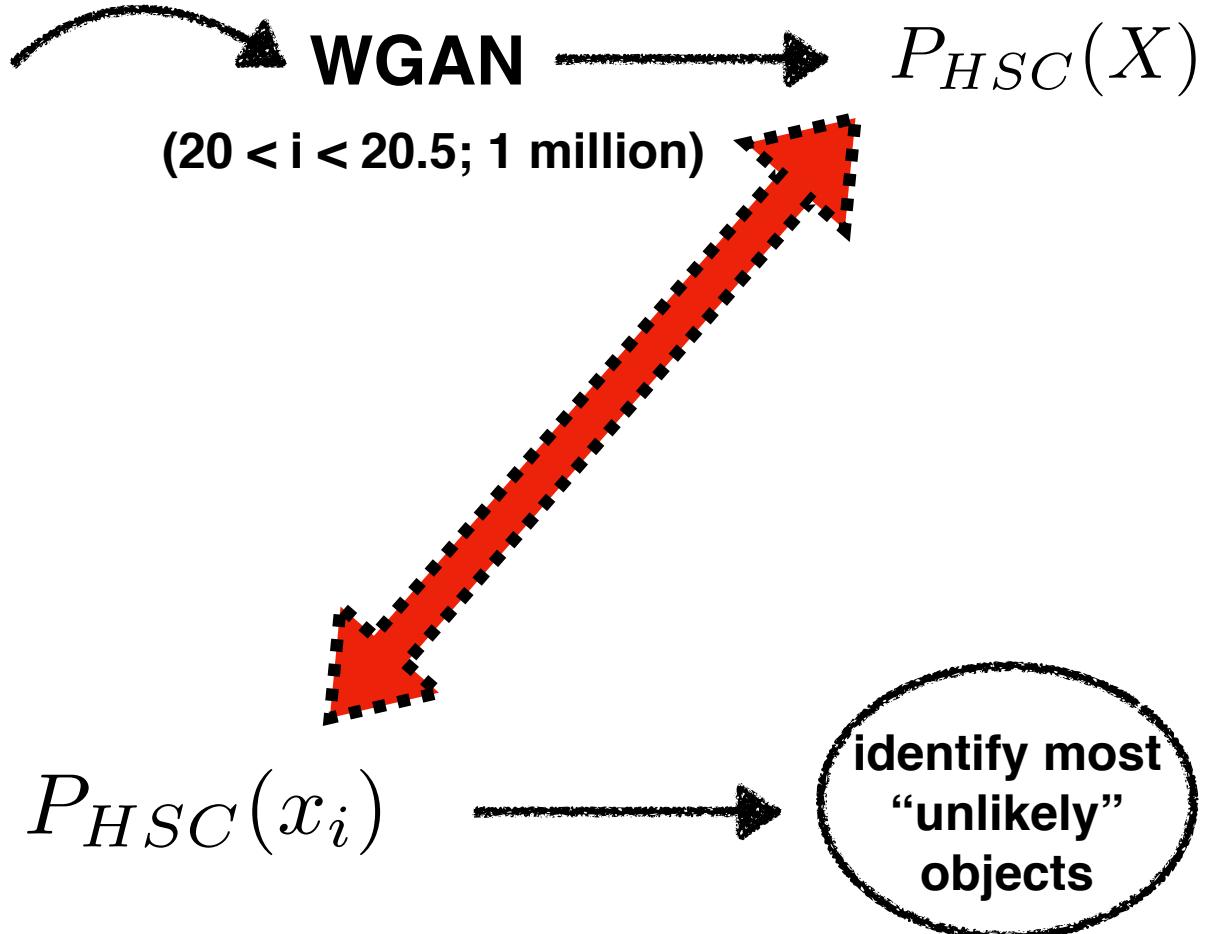


WGAN →  $P_{HSC}(X)$   
 $(20 < i < 20.5; 1 \text{ million})$

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

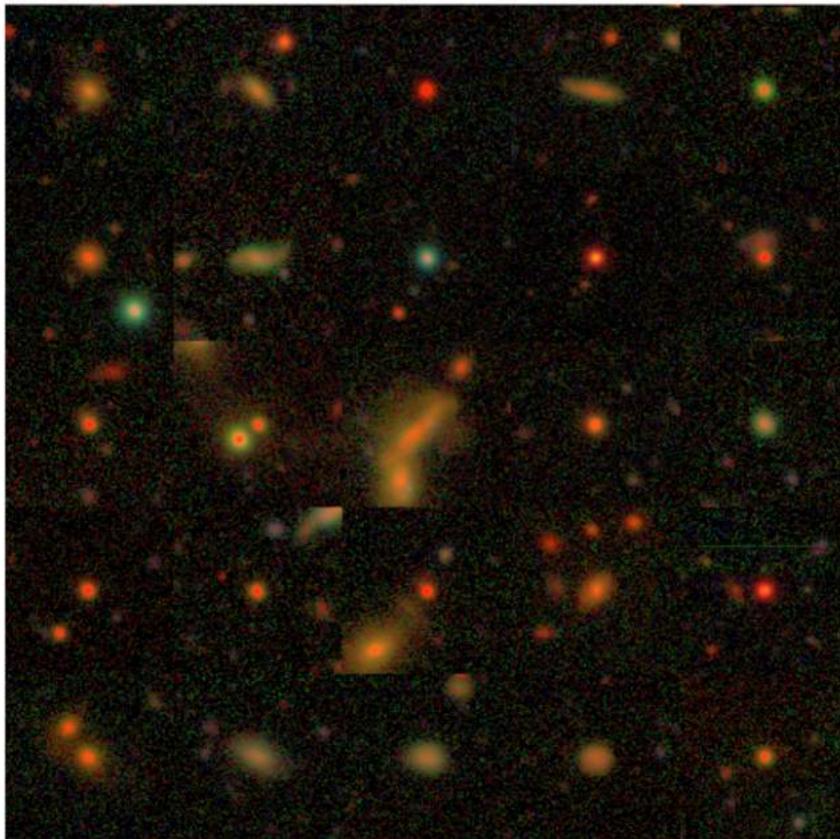


$x_i$  →  
**(individual detection)**

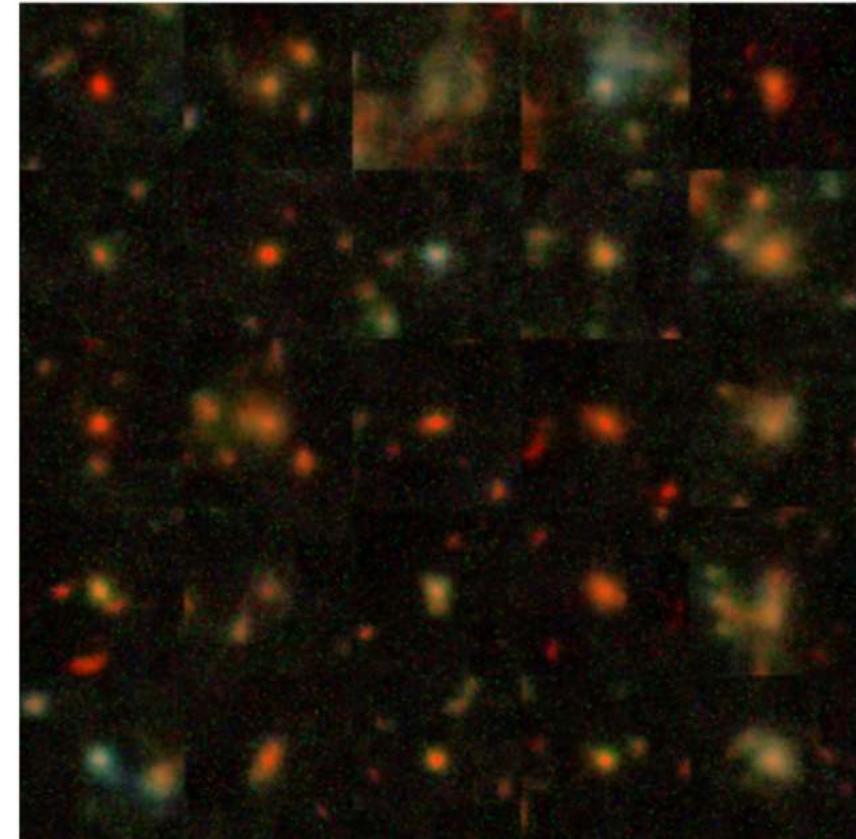


**Storey-Fisher, MHC+21**  
**Margalef-Bentabol, MHC+20**

**REAL**



**WGAN GENERATED**



# COMPUTE ANOMALY SCORE BASED ON WGAN RECONSTRUCTION ERROR

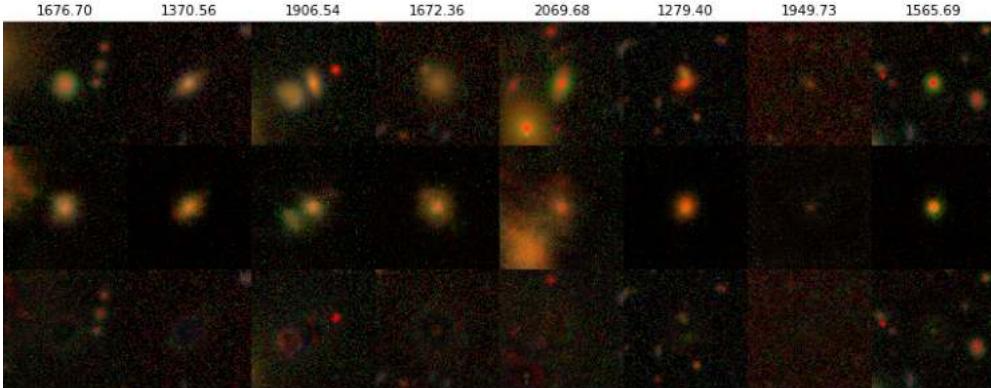
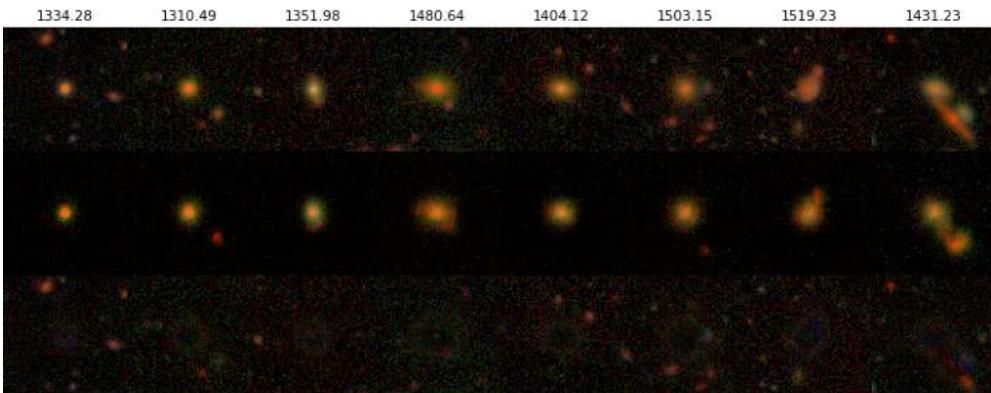
ANOMALY SCORE:

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

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

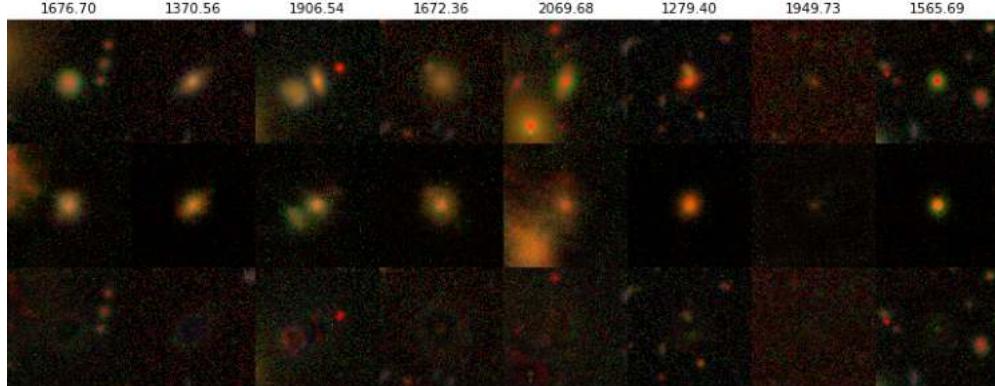
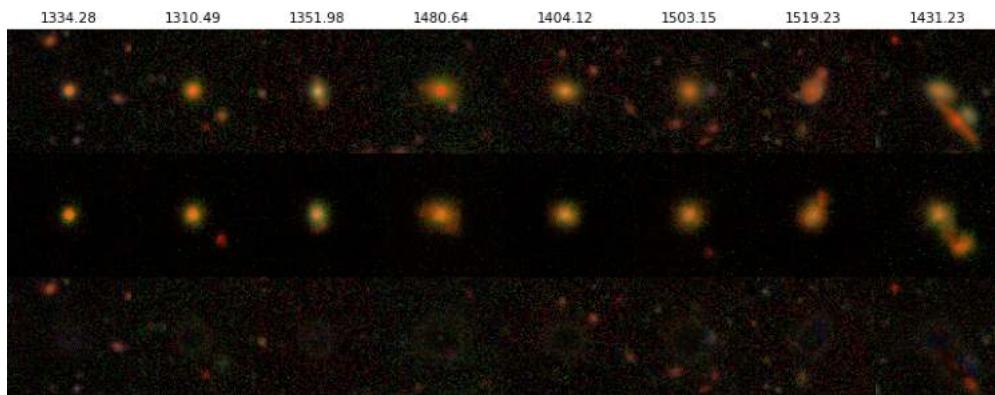
ANOMALY SCORE:

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

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

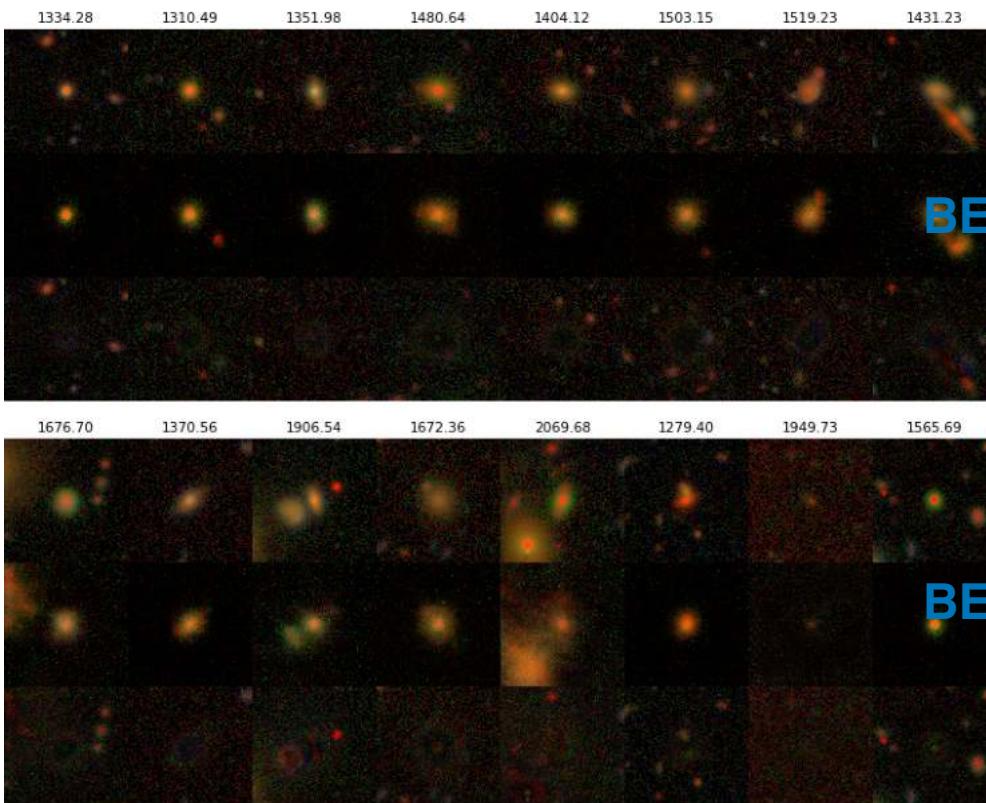
ANOMALY SCORE:

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

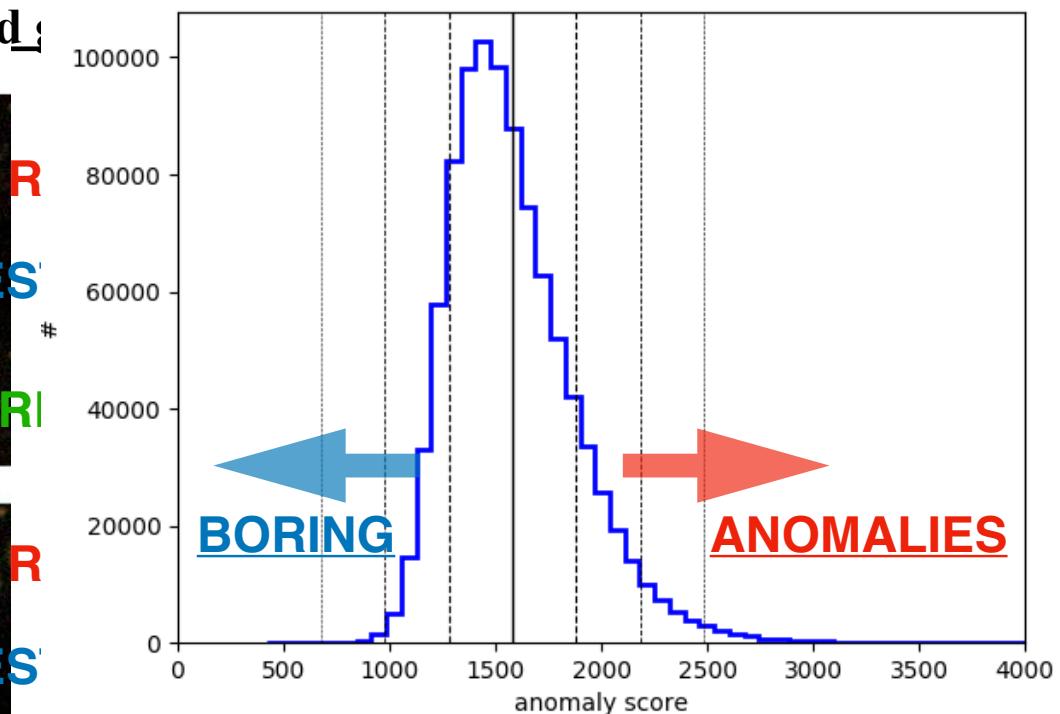
Hyperparameter

RMS  
and

RMS btw critic features



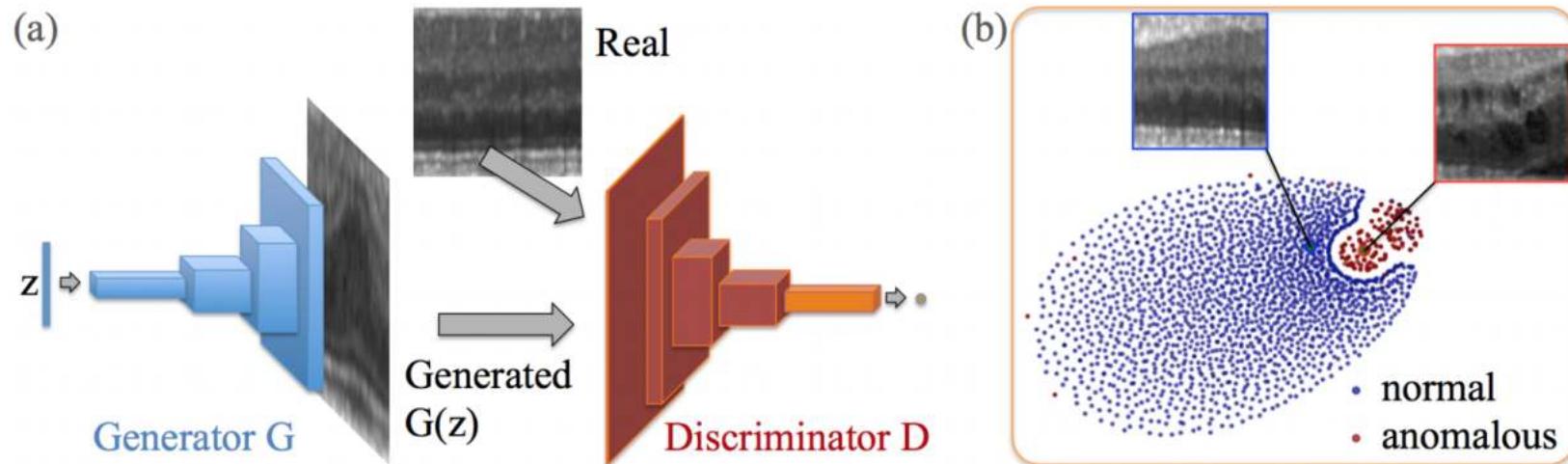
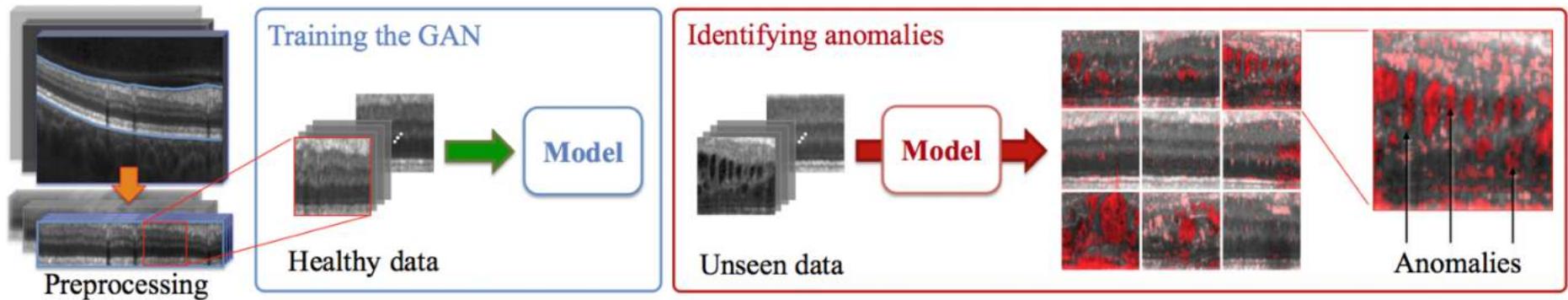
**RESIDUAL**



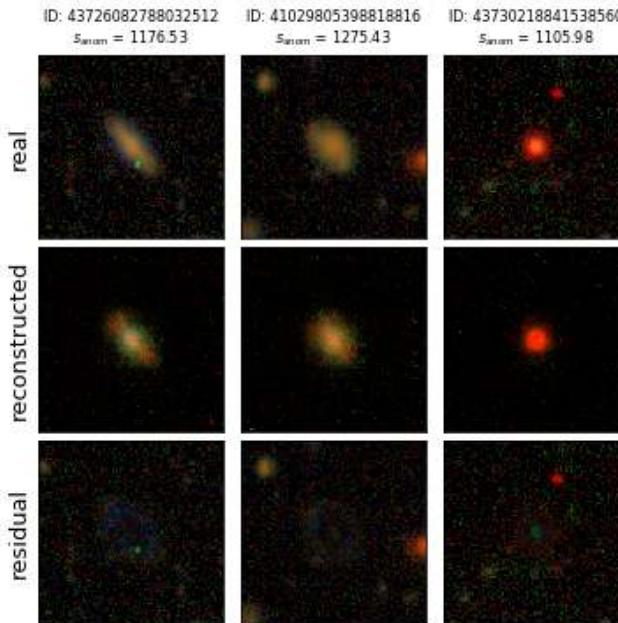
Storey-Fisher+21

Margalef-Bentabol, MHC+20

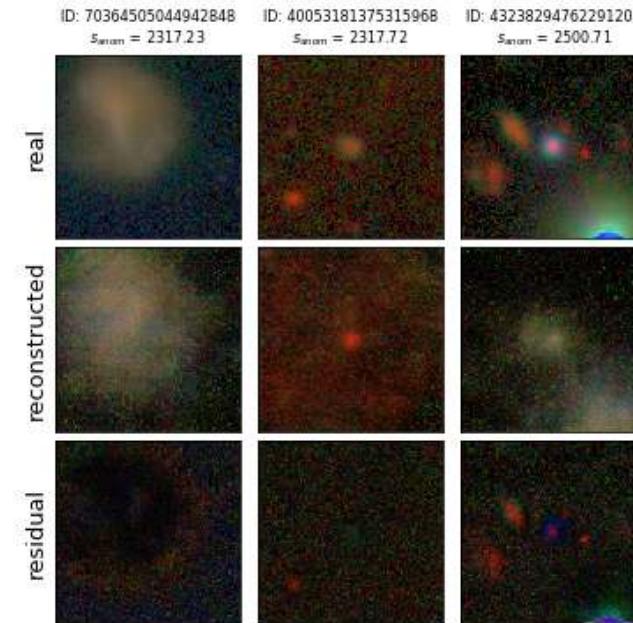
# This is again inspired by biomedical applications



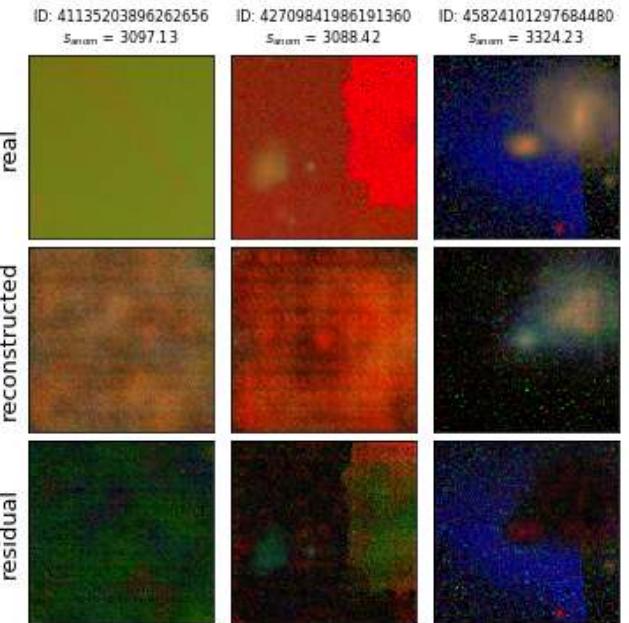
## AVERAGE



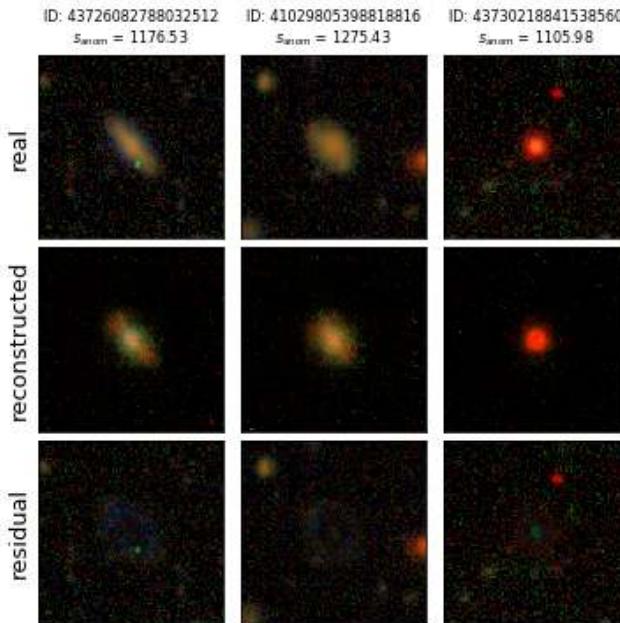
## 1-SIGMA



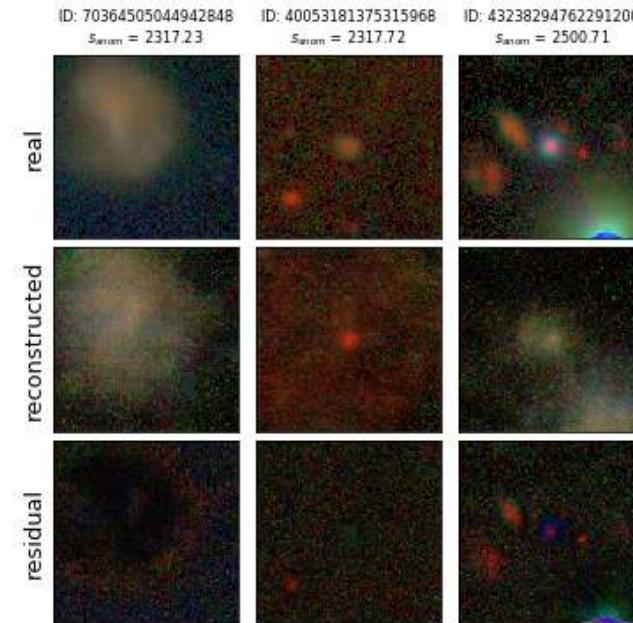
## >3-SIGMA



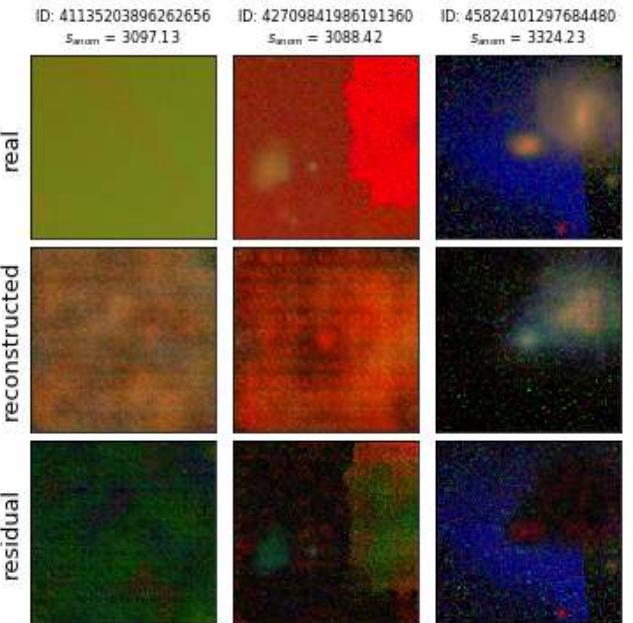
## AVERAGE



## 1-SIGMA

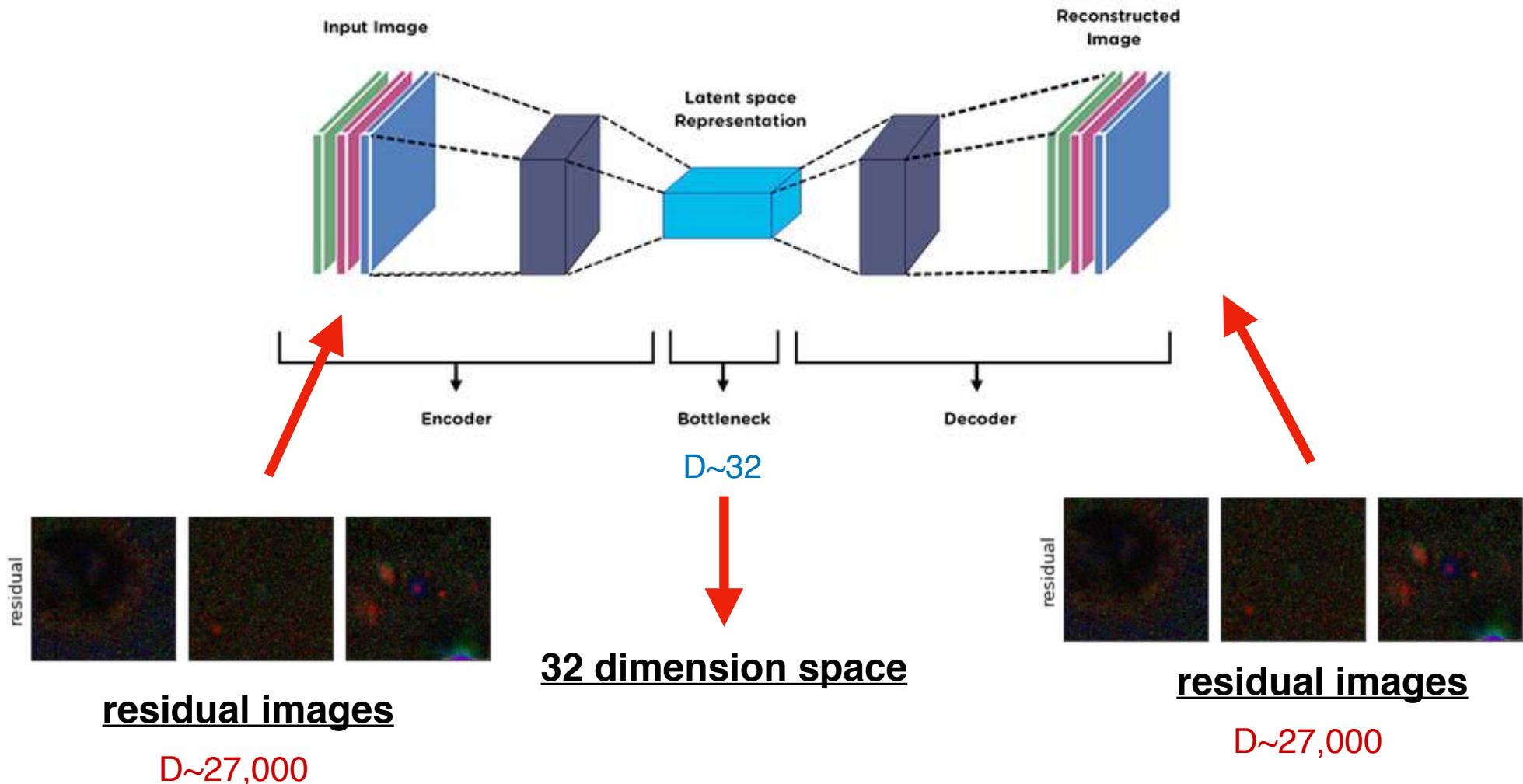


## >3-SIGMA



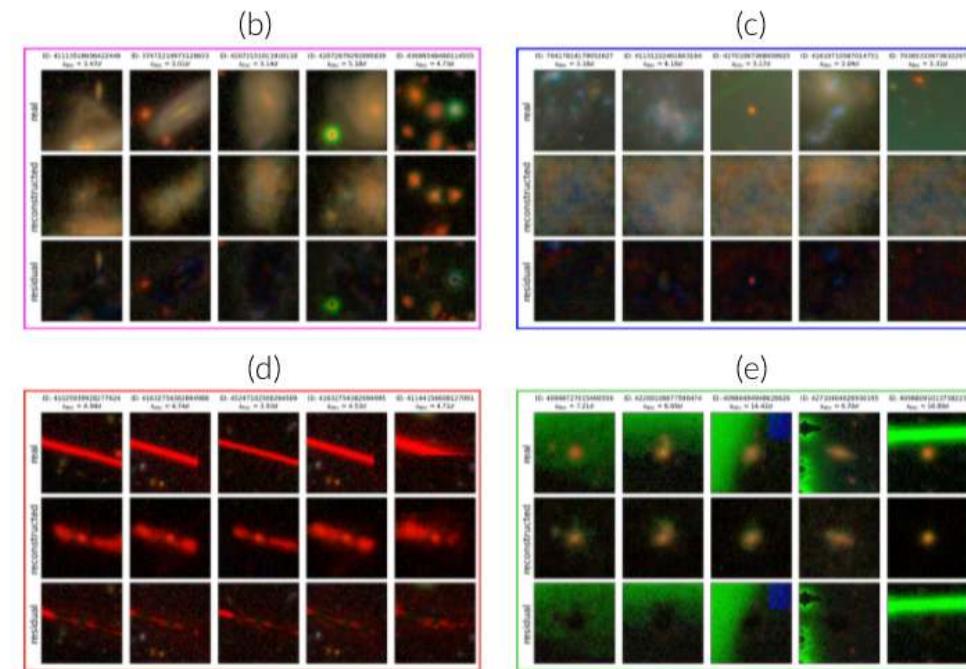
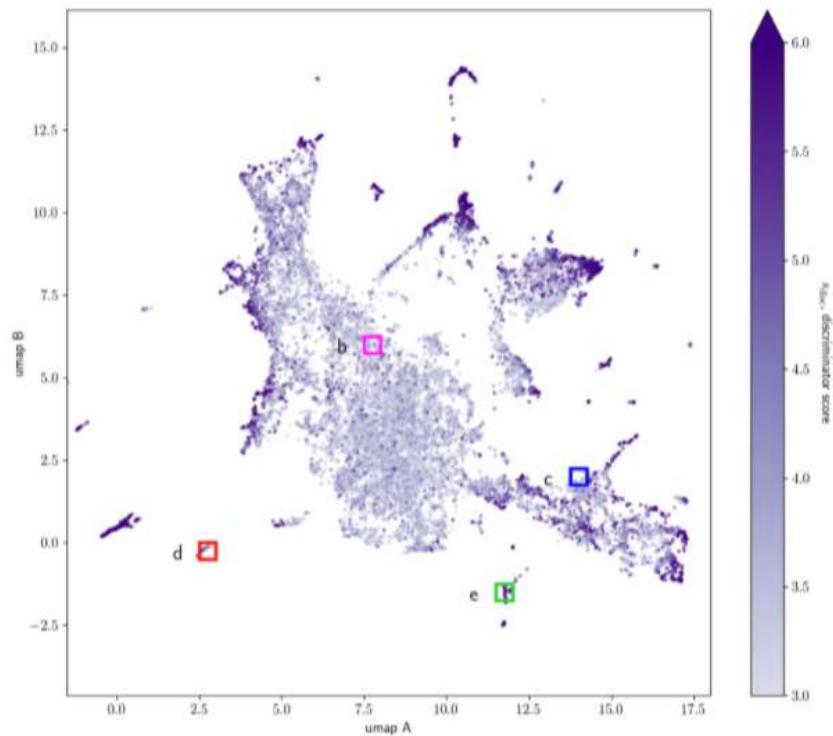
**HOW DO WE FILTER OUT “INTERESTING”  
ANOMALIES?**

# HUNTING “INTERESTING” OUTLIERS

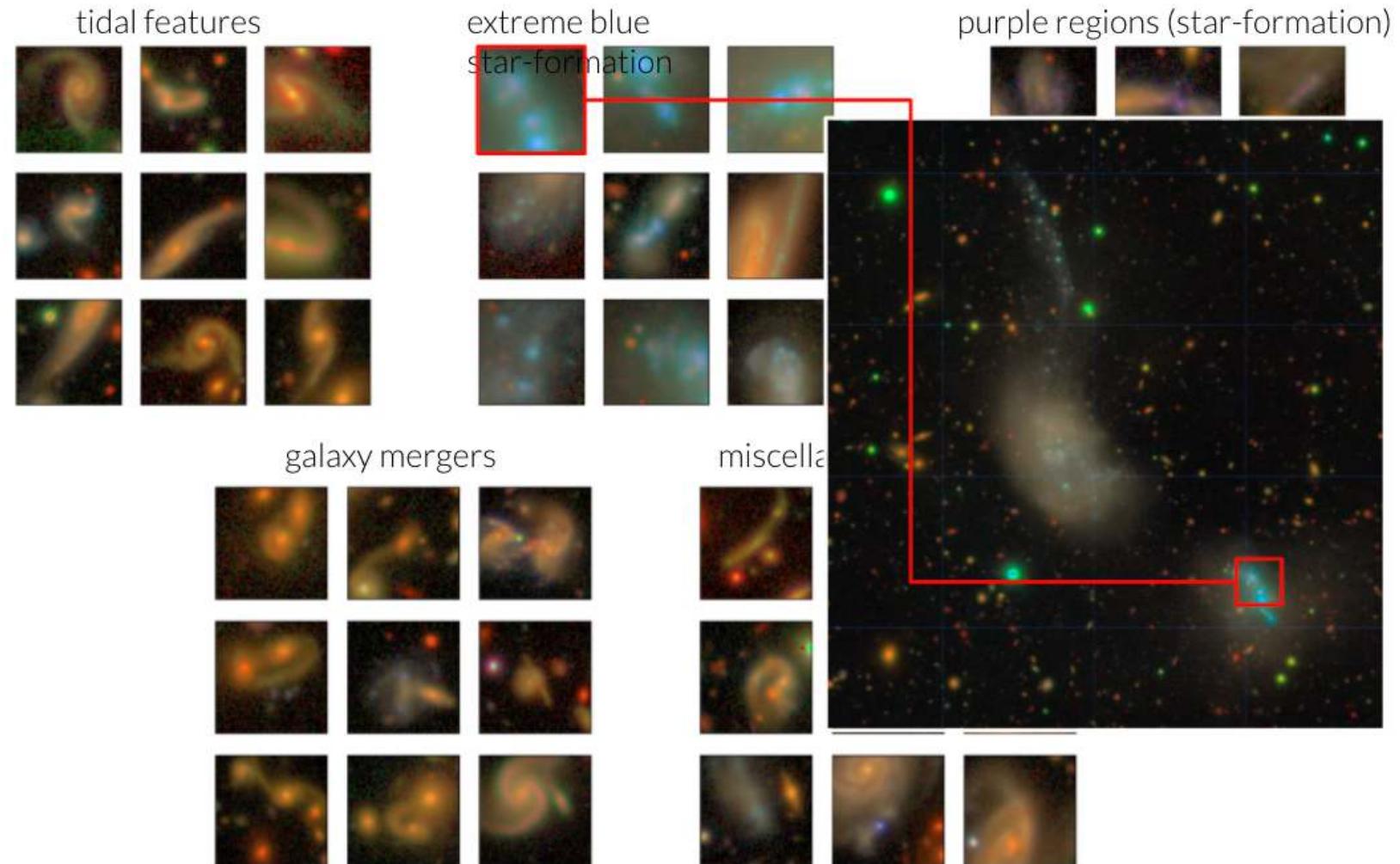


AUTO-ENCODERS TO REDUCE THE DIMENSIONALITY  
OF RESIDUAL IMAGES

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



# Detected interesting anomalies



### 3. Density Estimation for likelihood evaluation: simulation based inference



# Normalizing Flows

Based on change of variables:

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

unknown pdf

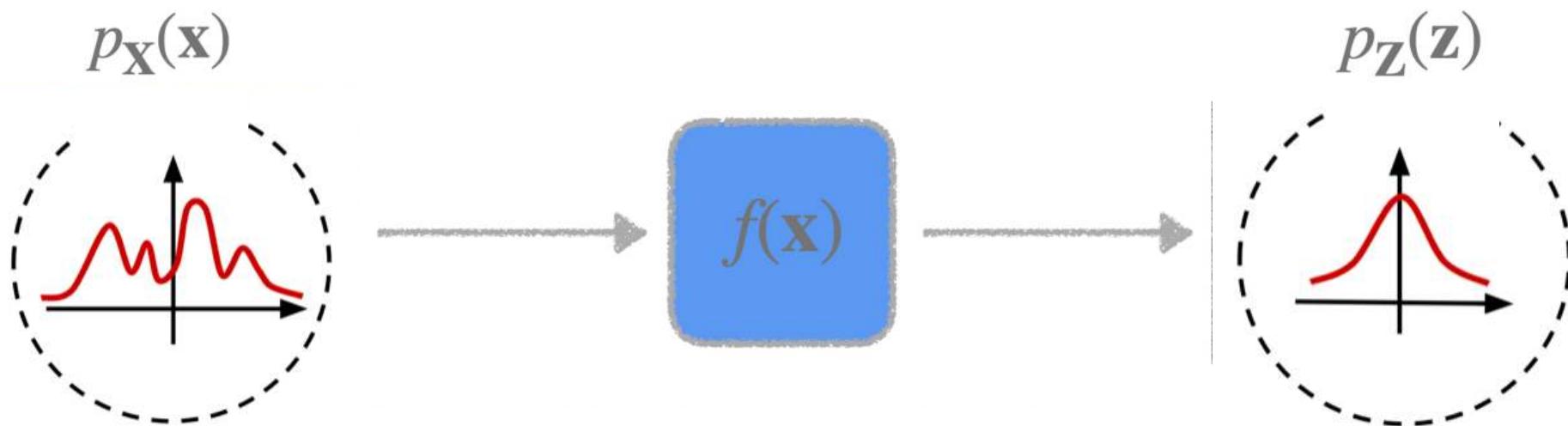
Invertible differentiable function

Determinant of the Jacobian of  $f(x)$

tractable pdf  
(typically Gaussian)

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

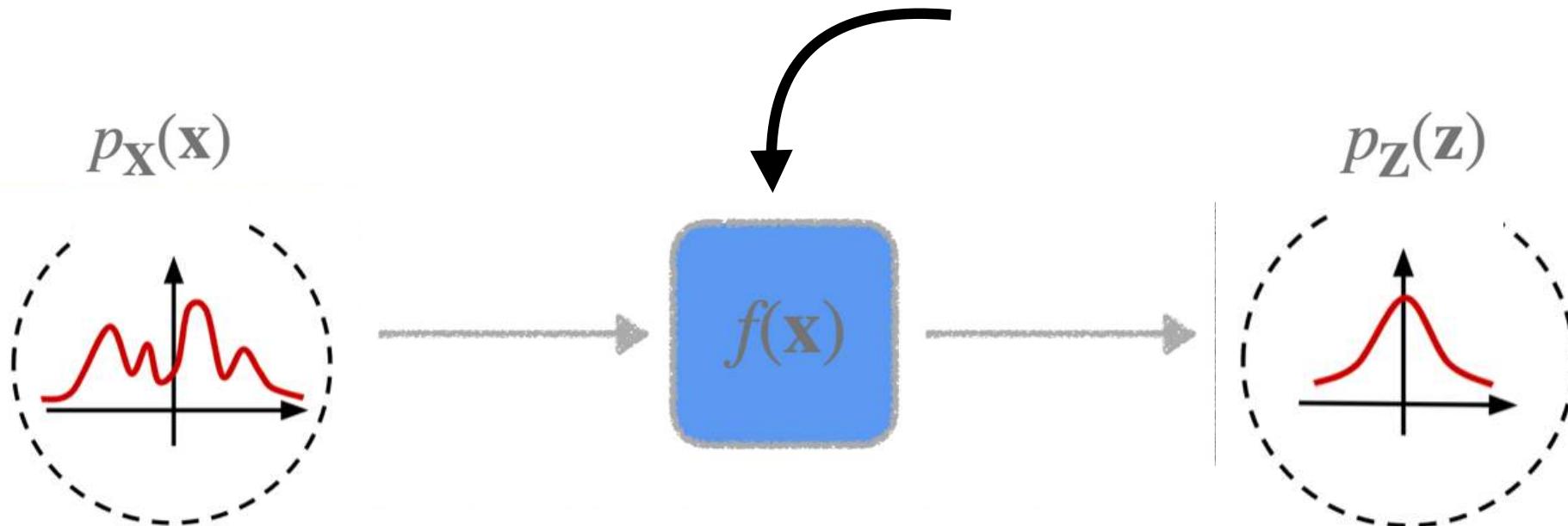
Can represent  $p_X(x)$  in terms of  $f(x)$  and  $p_Z(z)$   
(notice there is no dimensionality reduction as opposed to VAEs )



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

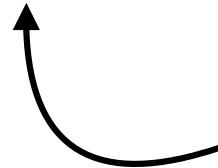
Can represent  $p_X(x)$  in terms of  $f(x)$  and  $p_Z(x)$   
(notice there is no dimensionality reduction as opposed to VAEs )

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



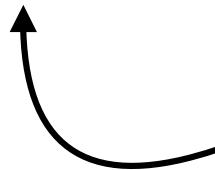
Then  $f(x)$  becomes:

$$f(x; W)$$



NN learnable weights

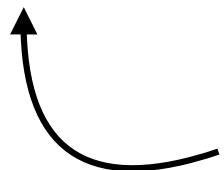
Then  $f(x)$  becomes:  $f(x; W)$



NN learnable weights

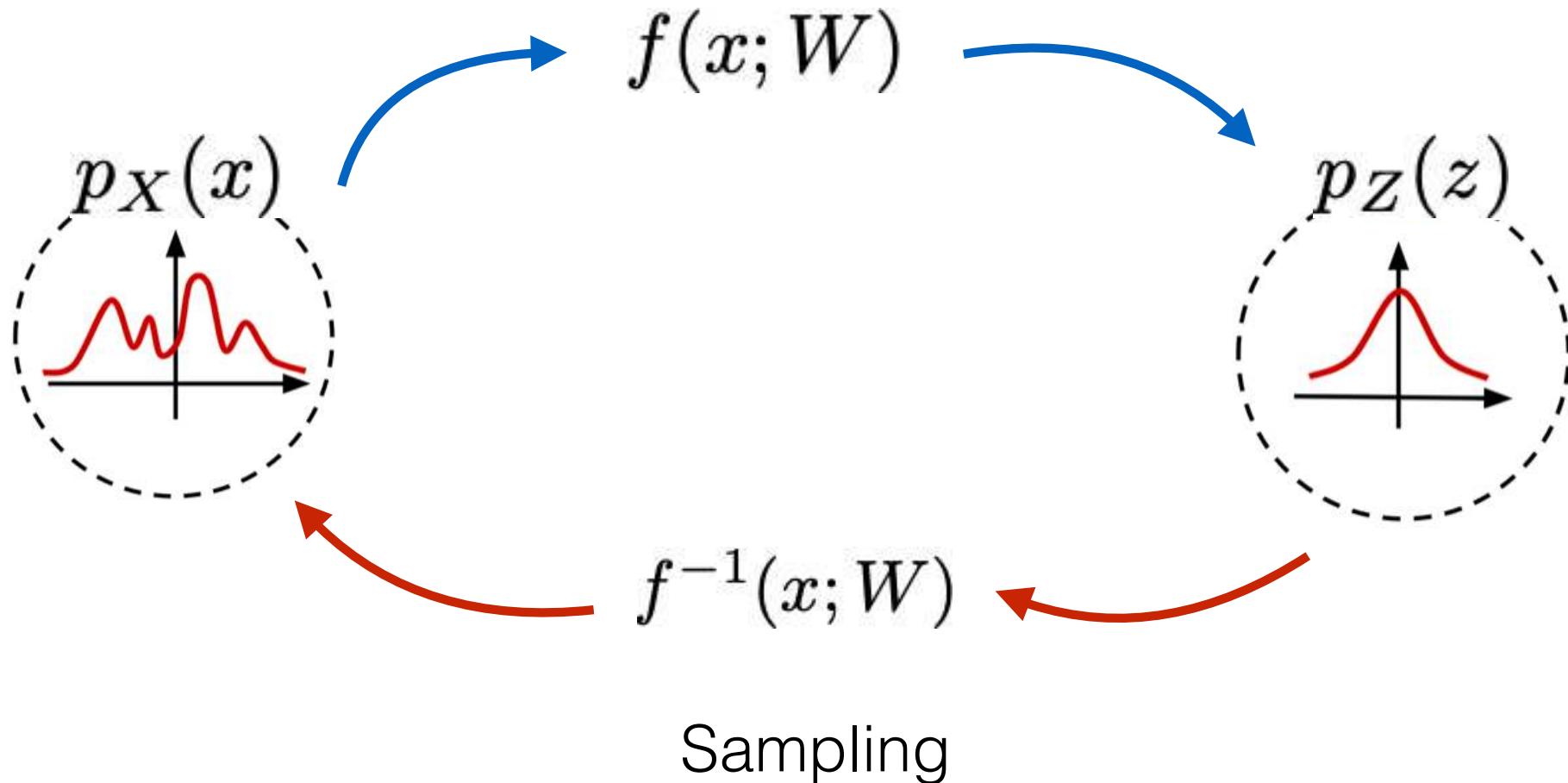
And the loss function: log likelihood

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

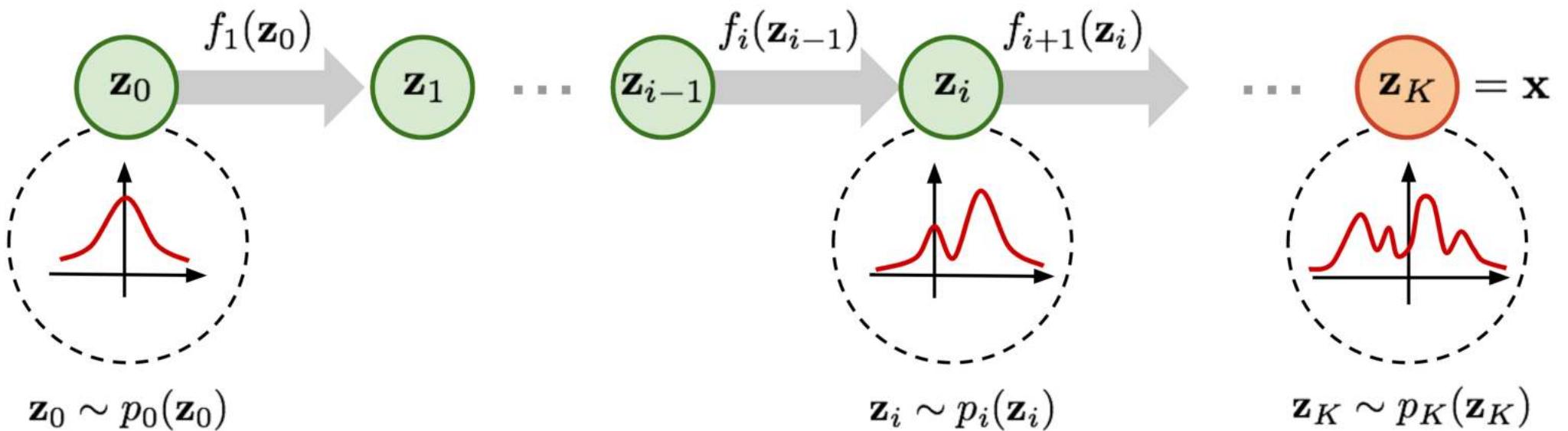


This is a gaussian

## Density 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

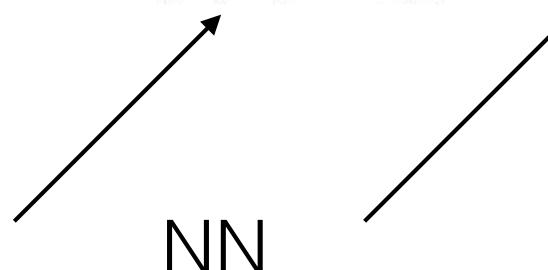
# What functions satisfy these conditions?

- 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:

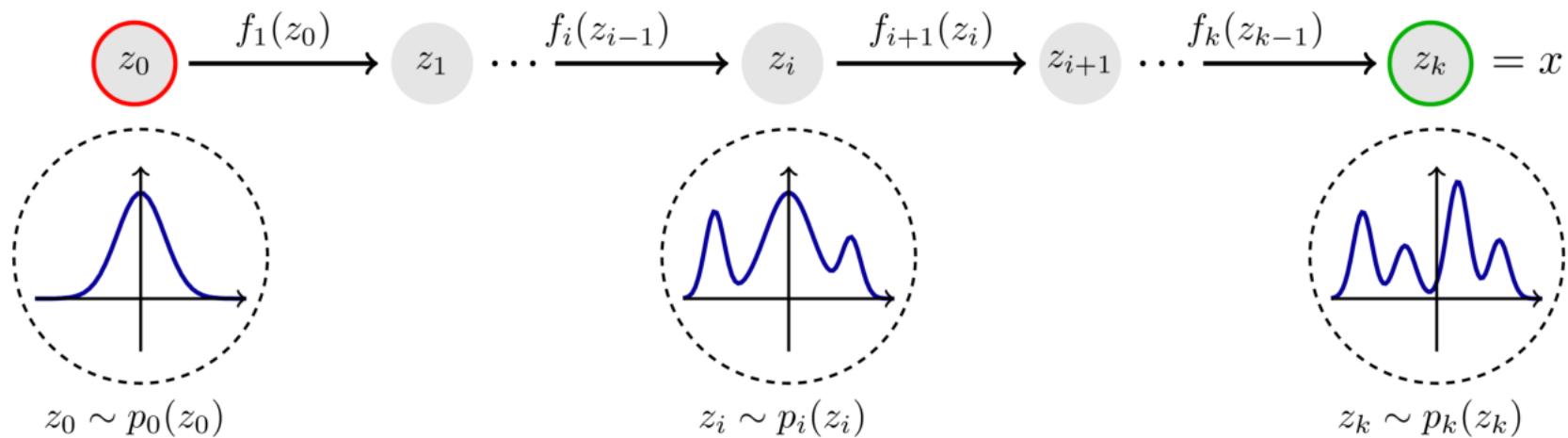
$$\begin{cases} \mathbf{y}_{1:d} &= \mathbf{x}_{1:d} \\ \mathbf{y}_{d+1:D} &= \mathbf{x}_{d+1:D} \odot \exp(s(\mathbf{x}_{1:d})) + t(\mathbf{x}_{1:d}) \end{cases} \Leftrightarrow \begin{cases} \mathbf{x}_{1:d} &= \mathbf{y}_{1:d} \\ \mathbf{x}_{d+1:D} &= (\mathbf{y}_{d+1:D} - t(\mathbf{y}_{1:d})) \odot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

- Jacobian:

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

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

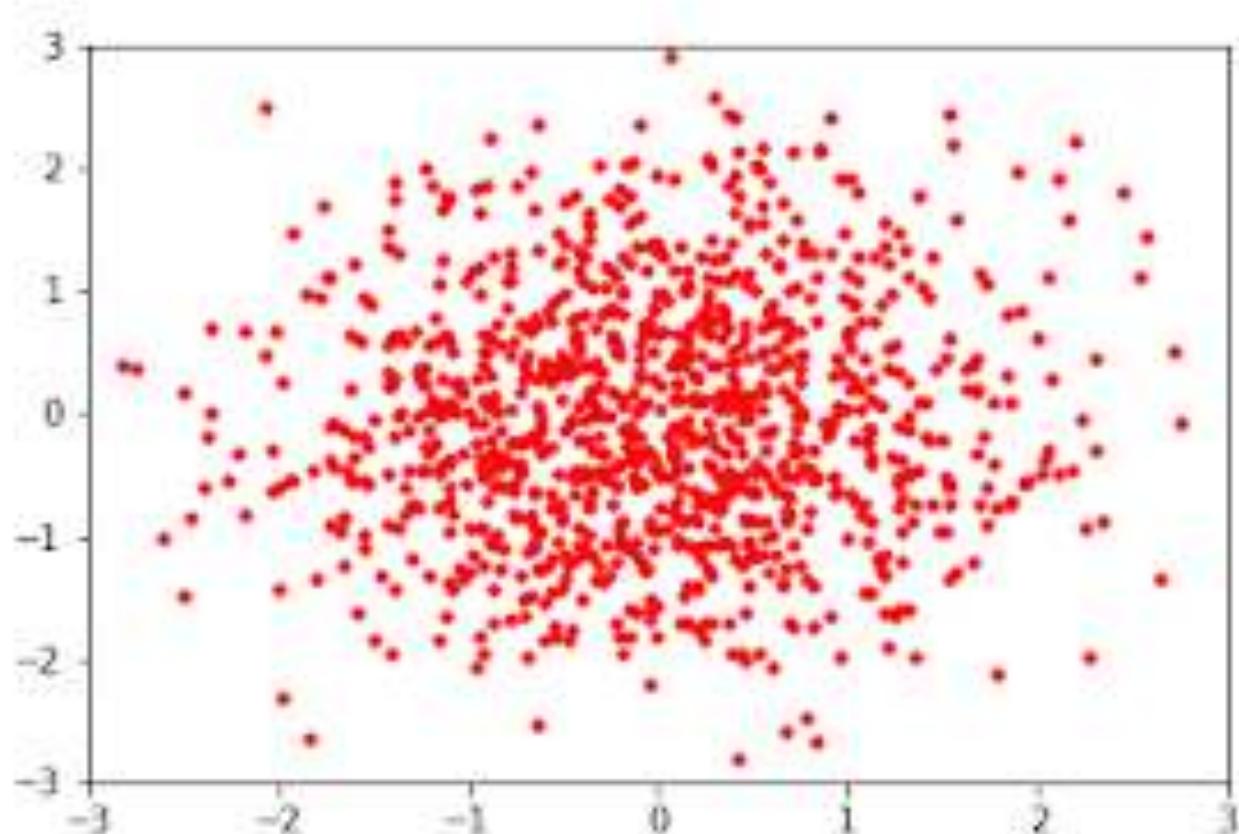
**normalizing flows** are generative models that are easy to evaluate and flexibly expressive



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

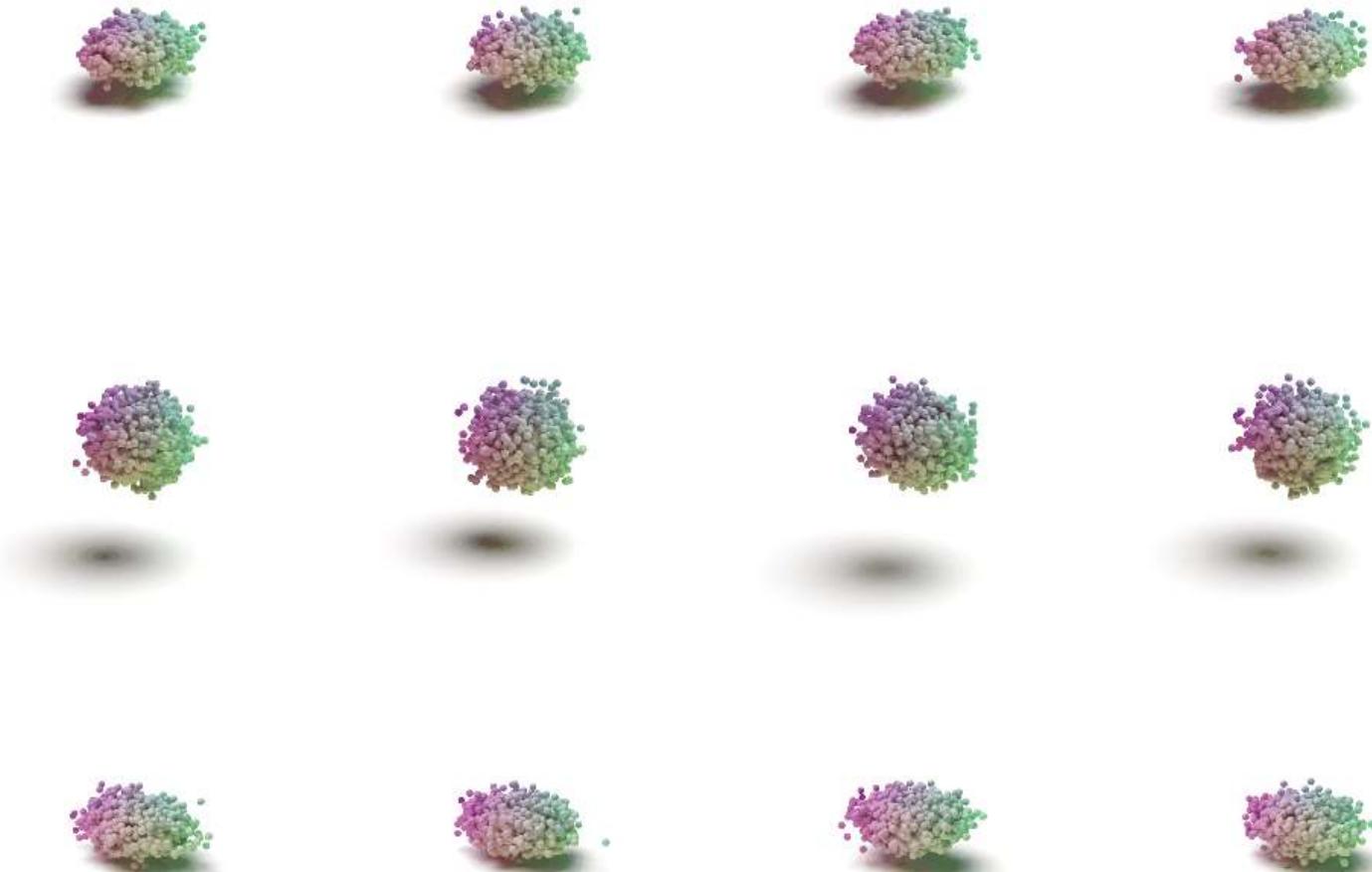
$f = f_1 \circ f_2 \dots \circ f_{k-1} \circ f_k$  is also invertible and differentiable

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



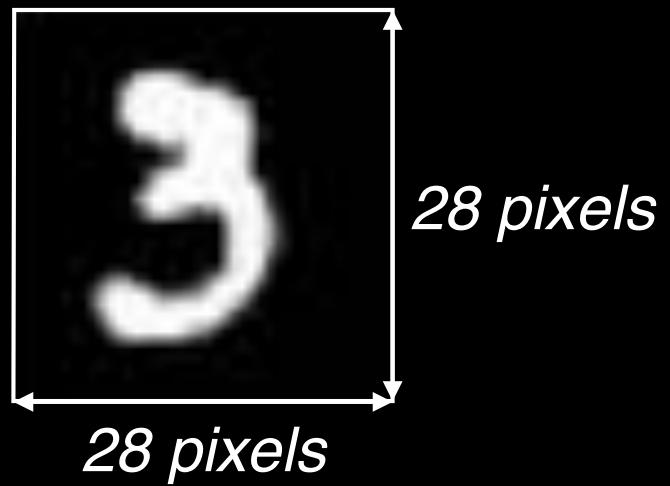
*credit: Eric Jang*

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

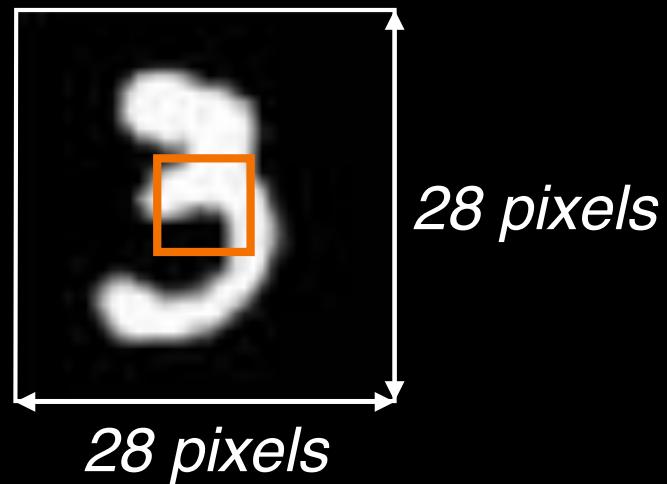


sample  $q_\phi(\text{pixels})$

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 5 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**



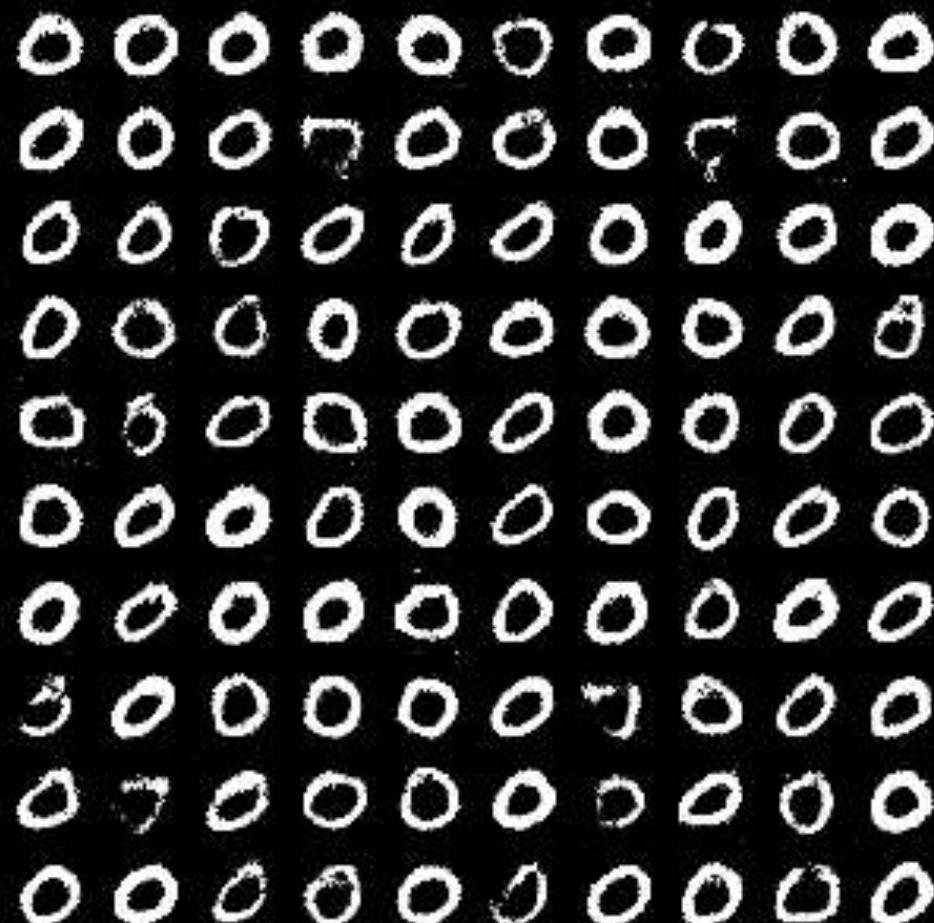
$$p(\text{pixels} \mid \text{central pixels}) \approx q_{\phi}(\text{pixels} \mid \text{central pixels})$$

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



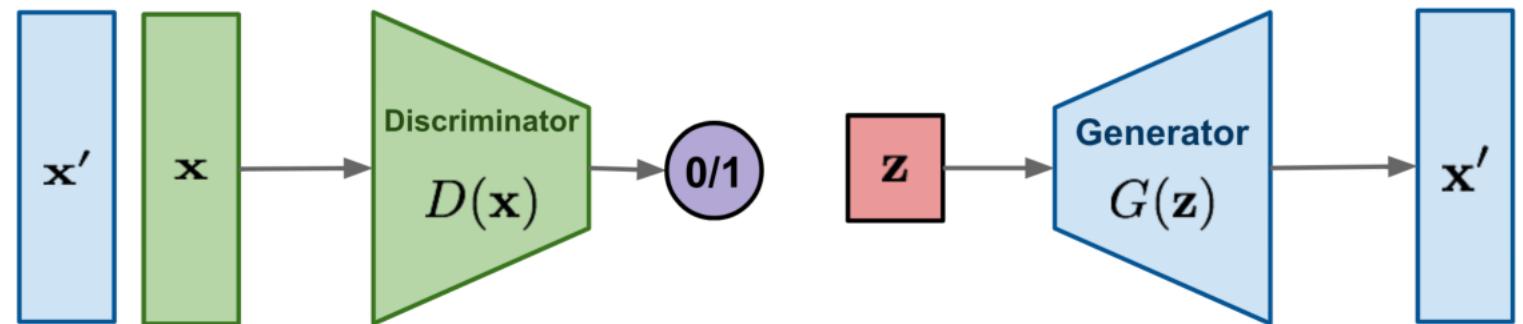
$$p(\text{pixels} \mid \text{central pixels}) \approx q_\phi$$

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

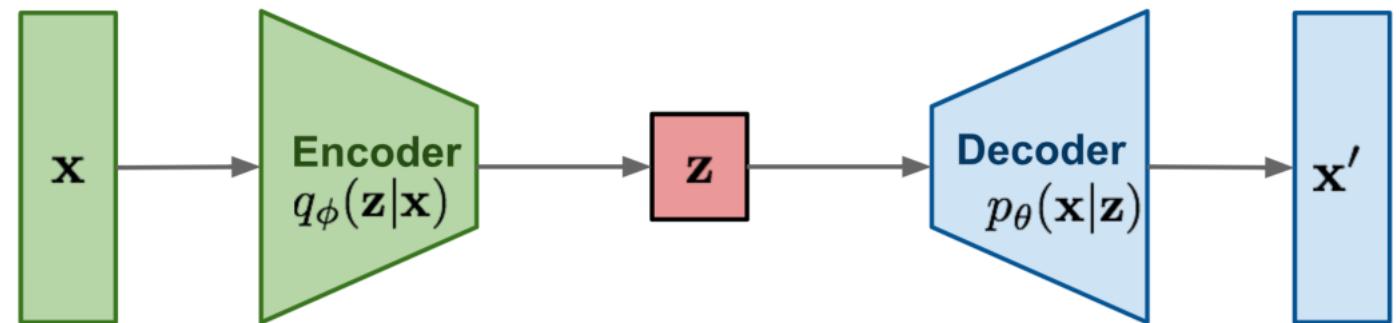
# GANs AND VAEs ARE VERY POWERFUL BUT DO NOT PROVIDE AN EXPLICIT LIKELIHOOD

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

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



**VAE:** maximize ELBO.



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

