

Deep Learning
**Unsupervised Learning and
Generative Models**

Generative modeling and unsupervised learning

Given a dataset X of (**multidimensional**) examples $\mathbf{x} \in X$, which we assume to have been drawn independently from some underlying distribution $p(X)$, a **generative model** can learn to **approximate** this distribution $p(X)$.

Such a model could be used:

- to **generate new samples** that look like they could have been part of the original dataset and/or
- to **infer the likelihood** of an example (useful for OOD detection).

Generative modeling and unsupervised learning

We distinguish **implicit** and **explicit** generative models:

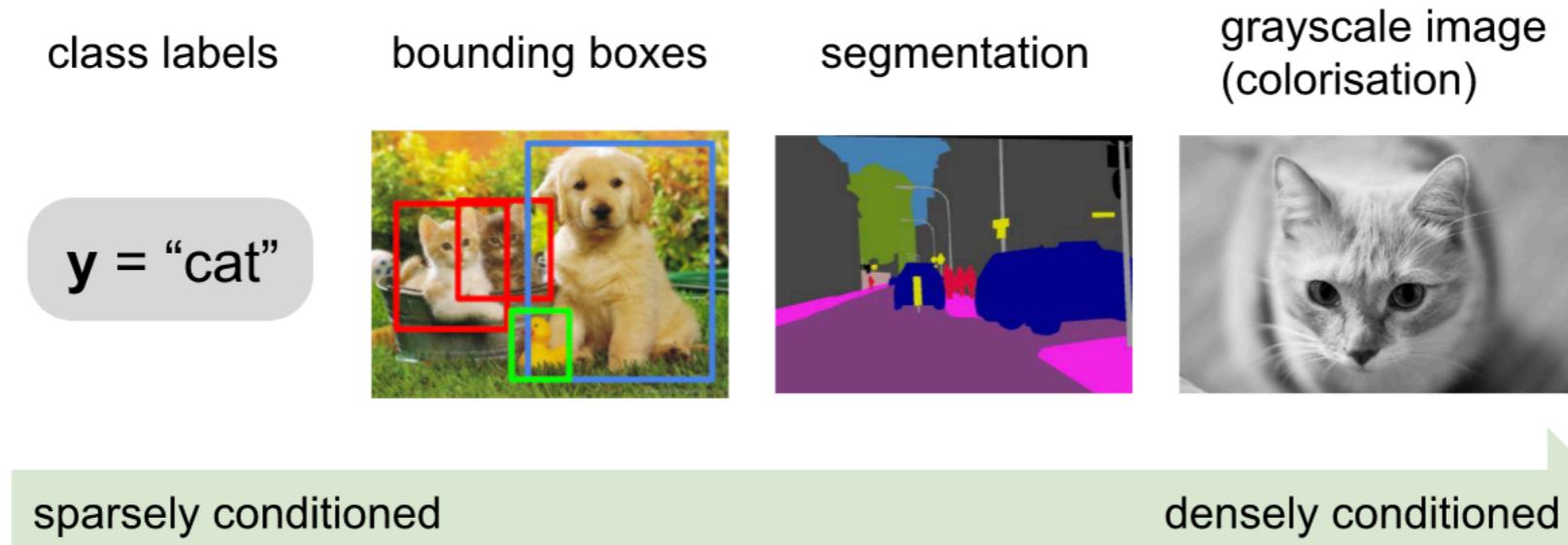
- An **implicit model** can produce new samples $\mathbf{x} \sim p(\mathbf{X})$, but cannot be used to infer the likelihood of an example (i.e. we cannot tractably compute $p(\mathbf{x})$).
- If we have an **explicit model**, we can infer the likelihood of an example.

Conditional Generative Models

Generative models become more practically useful when we can **exert some influence over the samples** we draw from them.

We can do this by providing a **conditioning signal** c , which contains side information about the kind of samples we want to generate. The model is then fit to the conditional distribution $p(X | c)$ instead of $p(X)$.

The generative modeling problem becomes easier if the **conditioning signal c is richer**, because it reduces uncertainty about x .



Generative Models Families

- Likelihood-based models.
 - Autoregressive models.
 - Normalizing Flows (<https://deepgenerativemodels.github.io/notes/flow/>).
 - Variational Autoencoders.
 - Diffusion models.
- Adversarial Models (GAN's).

Likelihood based models

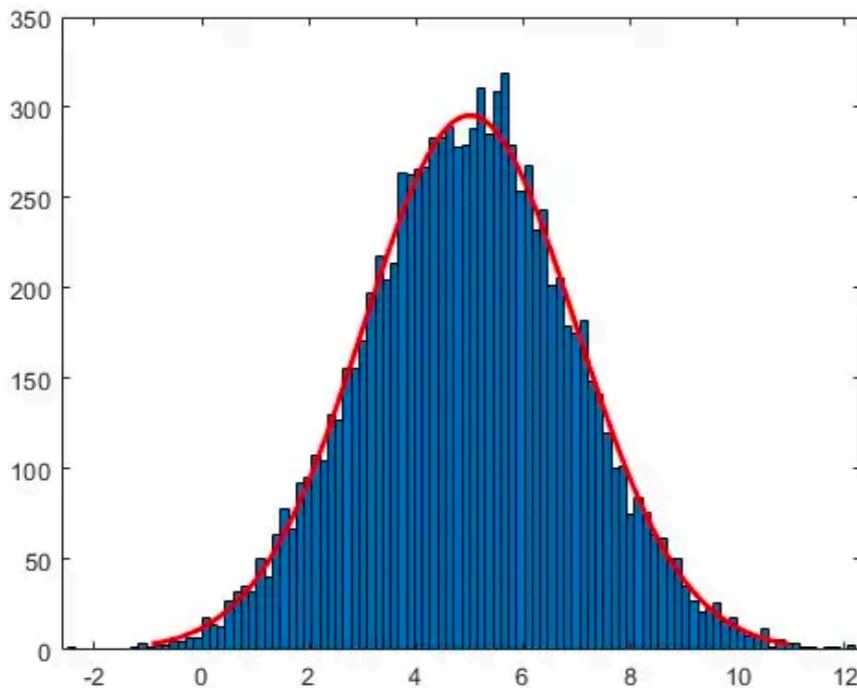
Likelihood-based models directly parameterize $p(X)$.

The parameters θ are then fit by maximizing the (log) likelihood of the data under the model:

$$\mathcal{L}_\theta = \sum_{x \in X} p(x | \theta)$$

$$\theta^* = \arg \max_{\theta} \mathcal{L}_\theta(x)$$

Gaussian data example:



$$\mu, \sigma = \arg \max_{\mu, \sigma} \prod_{i=1}^N p(x_i | \mu, \sigma) = \arg \max_{\mu, \sigma} \prod_{i=1}^N \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}$$

Data likelihood

This is an optimization
problem that can be
solved with SGD.

Likelihood based models

This is typically done in the log-domain because it **simplifies computations and improves numerical stability**.

$$\mathcal{L}_\theta = \sum_{x \in X} \log p(x | \theta) \quad \theta^* = \arg \max_{\theta} \mathcal{L}_\theta(x)$$

Because the model directly parameterizes $p(X)$, we can easily infer the likelihood of any \mathbf{x} , so we get an **explicit model**.

Three popular flavors of likelihood-based models are:

- **autoregressive models**,
- flow-based models, and
- **variational autoencoders**.

Autoregressive models

In an autoregressive model, we assume that the distribution of our **multidimensional examples** $\mathbf{x} \in X$ can be factorized into a product of conditionals, using the **chain rule of probability**:

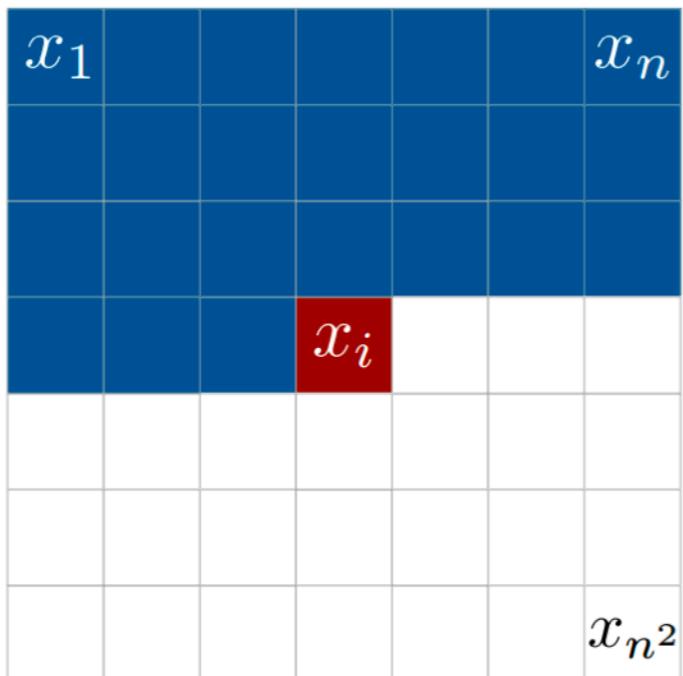
$$p(\mathbf{x}) = \prod_i p(x_i | x_{<i}) = p(x_0)p(x_1 | x_0)p(x_2 | x_0, x_1) \dots p(x_n | x_0, \dots x_{n-1})$$

These conditional distributions, $p(\cdot)$, are typically scalar-valued and much **easier to model**.

Autoregressive models

For time-series, audio signals, NLP..., AR models are a very natural thing to do.

But **we can also do this for other types of structured data** by arbitrarily choosing an order (e.g. raster scan order for images, as in *PixelRNN2* and *PixelCNN3*).



$$p(\mathbf{x}) = \prod_{i=1}^{n^2} p(x_i | x_{<i})$$

Everyone of this conditional probability distributions can be implemented by a NN (f.e. a RNN)

Autoregressive models

Autoregressive models are attractive because they are able to **accurately capture correlations between the different elements** x_i in a sequence, and they allow for **fast inference** (i.e. computing $p(\mathbf{x})$).

Unfortunately they tend to be **slow to sample from** (get $\mathbf{x} \sim p(\mathbf{X})$), because samples need to be drawn sequentially from the conditionals for each position in the sequence.

Autoregressive models

In the case of color images, the conditional probability of i^{th} pixel becomes:

$$p(x_{i,R} | \mathbf{x}_{<i}) p(x_{i,G} | \mathbf{x}_{<i}, x_{i,R}) p(x_{i,B} | \mathbf{x}_{<i}, x_{i,R}, x_{i,G})$$

Thus, each color is conditioned on other colors as well as the previously generated pixels.

To get the appropriate pixel value we use a 256-way softmax layer.

We could use **LSTMs** to model conditional probabilities in a natural way.

Autoregressive models

In order to use **convolutional layers** when modeling conditional probabilities, **masked convolutions** can be used:

1	1	1	1	1
1	1	1	1	1
1	1	0	0	0
0	0	0	0	0
0	0	0	0	0

In this way we can build functions that restrict information flow from ‘future’ pixels into the one we are predicting.

Conditional autoregressive models

$p(\mathbf{x}, C) = \prod_i p(x_i | x_{<i}, C)$ where each $p(x_i | x_{<i}, C)$ is a neural network
that predicts a probability distribution over pixel values.

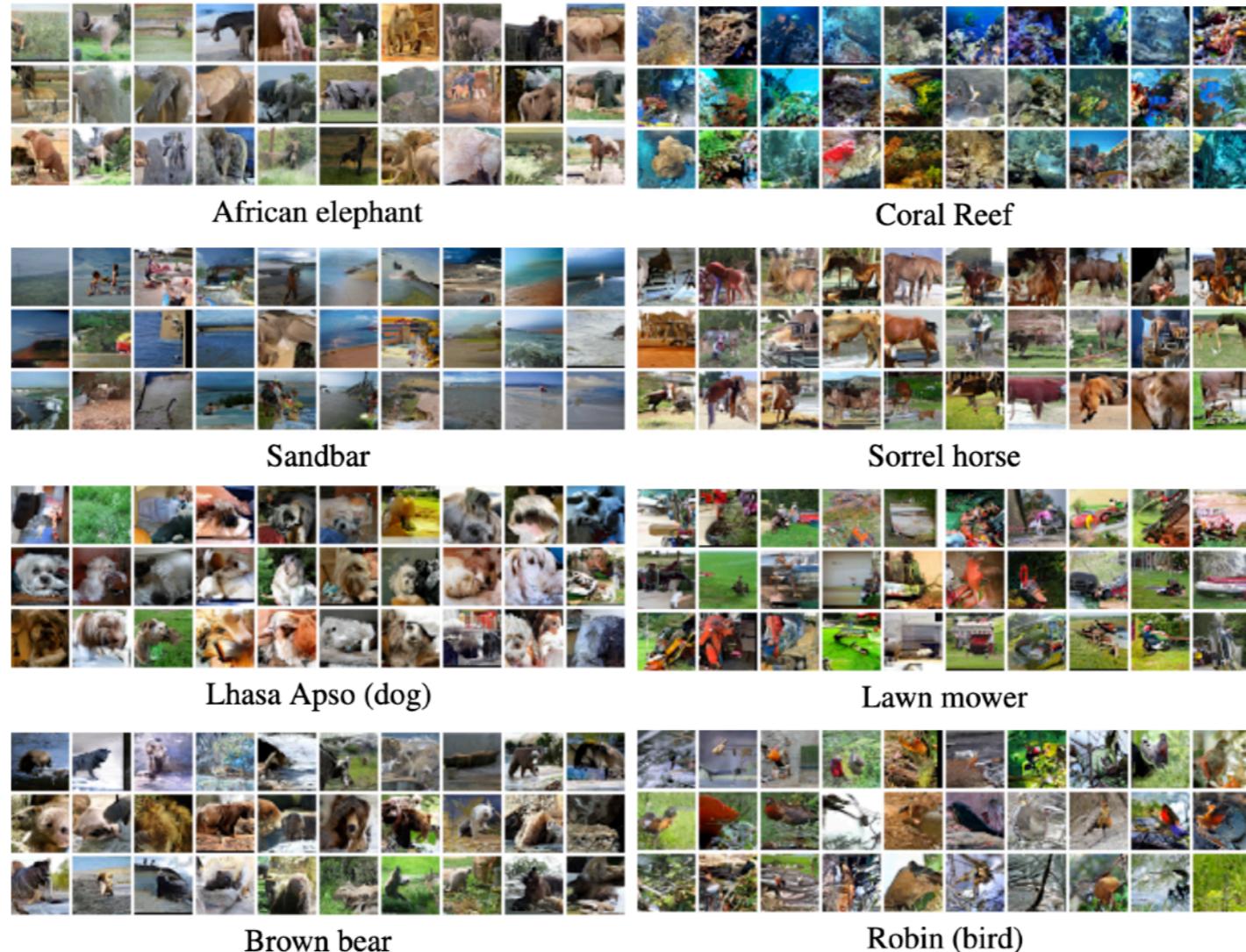


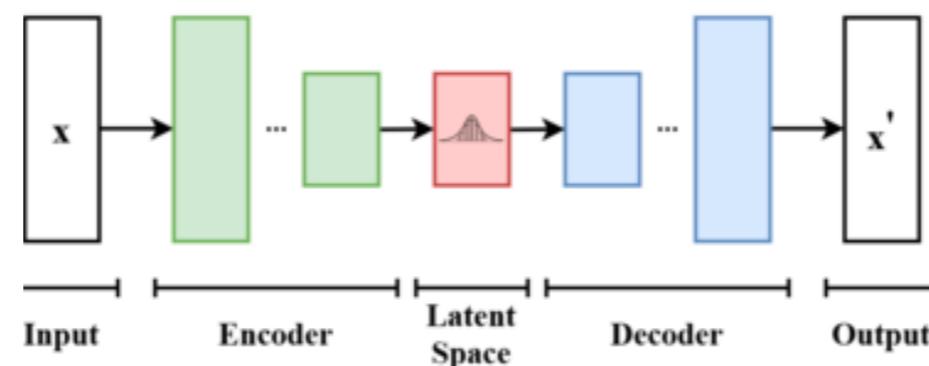
Figure 3: Class-Conditional samples from the Conditional PixelCNN.

Variational Autoencoders

By far the **most popular class** of likelihood-based generative models,

Variational Autoencoders learn two neural networks: an *inference network* or encoder $q(z|x)$ learns to probabilistically map examples x into a **latent space** Z , and a *generative network* or decoder $p(x|z)$ learns the distribution of the data conditioned on a latent representation z .

These are trained to maximise a lower bound on $p(x)$, called the **ELBO** (*Evidence Lower BOund*), because computing $p(x)$ given x (exact inference) is not tractable.



Adversarial models

Generative Adversarial Networks (GANs) take a **very different approach** to capturing the data distribution.

Two networks are trained simultaneously:

- a **generator** G attempts to produce examples according to the data distribution $p(x)$, given latent noise vectors z , while
- a **discriminator** D attempts to tell apart generated examples and real examples. In doing so, the discriminator provides a learning signal for the generator which enables it to better match the data distribution.

In the original formulation, the loss function is as follows:

$$\mathcal{L}(x) = \mathbb{E}_x[\log D(x)] + \mathbb{E}_z[\log(1 - D(G(z)))].$$

Adversarial models

$$\mathcal{L}(x) = \mathbb{E}_x[\log D(x)] + \mathbb{E}_z[\log(1 - D(G(z)))].$$

Probability that the
discriminator is rightly
classifying the real image

Probability that the
discriminator is rightly
classifying the fake image

The generator is trained to **minimize** this loss, whereas the discriminator attempts to **maximize** it.

Discriminator Loss

It penalizes itself for misclassifying a real instance as fake, or a fake instance (created by the generator) as real, by **maximizing** the function.

$$\mathcal{L}(x) = \mathbb{E}_x[\log D(x)] + \mathbb{E}_z[\log(1 - D(G(z)))].$$

Adversarial models

$$\mathcal{L}(x) = \mathbb{E}_x[\log D(x)] + \mathbb{E}_z[\log(1 - D(G(z)))].$$

Probability that the generator is rightly classifying the real image

Probability that the generator is rightly classifying the fake image

Generator Loss

It gets rewarded if it successfully fools the discriminator, and gets penalized otherwise:

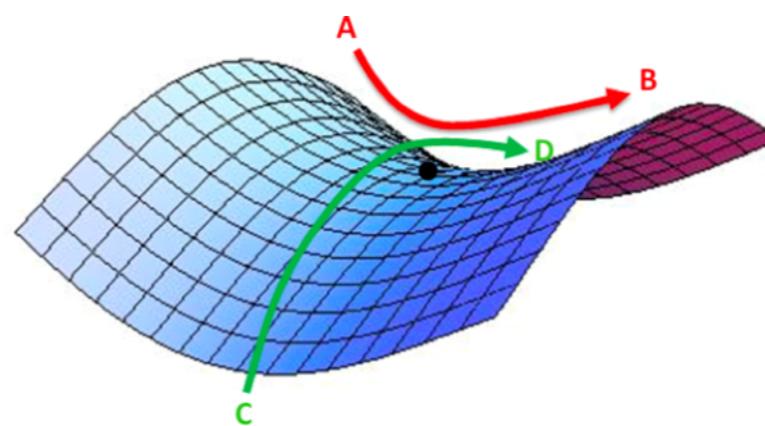
$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D \left(G \left(z^{(i)} \right) \right) \right)$$

The generator is trained to **minimize** this loss, whereas the discriminator attempts to **maximize** it.

Adversarial models

$$\mathcal{L}(x) = \mathbb{E}_x[\log D(x)] + \mathbb{E}_z[\log(1 - D(G(z)))].$$

This means the training procedure is a **two-player minimax game**, rather than an optimization process, as it is for most machine learning models. Balancing this game and keeping training stable has been one of the main challenges for this class of models. Many alternative formulations have been proposed to address this.



Adversarial models

While adversarial and likelihood-based models are both ultimately trying to model $p(\mathbf{x})$, they approach this target from very different angles. As a result, GANs tend to be **better at producing realistic examples, but worse at capturing the full diversity of the data distribution**, compared to likelihood-based models.



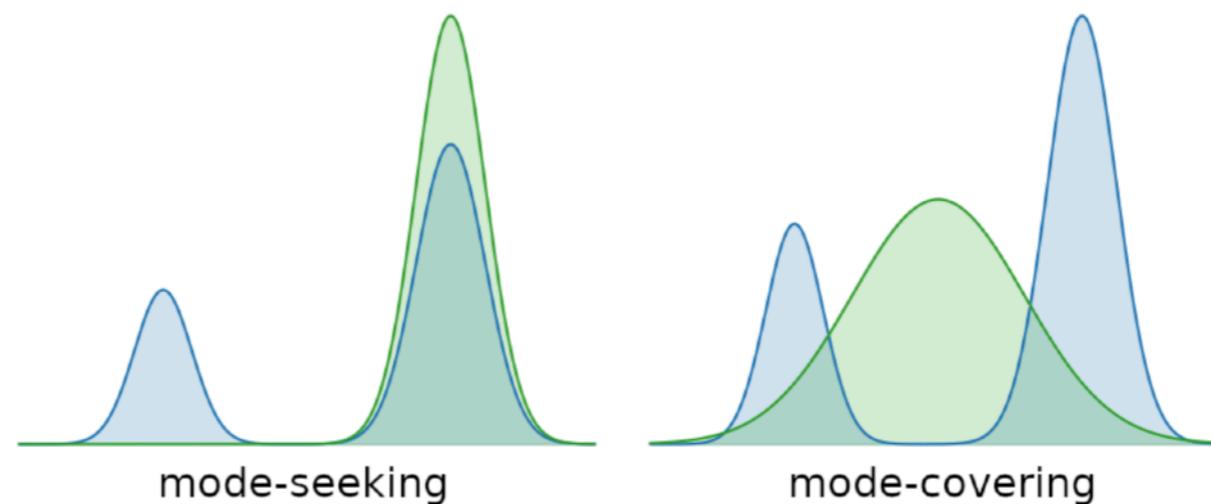
An image generated by a [StyleGAN](#) that looks deceptively like a photograph of a real person.

Mode-covering vs. mode-seeking behaviour

An important consideration when determining which type of generative model is appropriate for a particular application, is **the degree to which it is mode-covering or mode-seeking**. When a model does not have enough capacity to capture all the variability in the data, different compromises can be made.

If all examples should be reasonably likely under the model, it will have to overgeneralise and put probability mass on interpolations of examples that may not be meaningful (**mode-covering**).

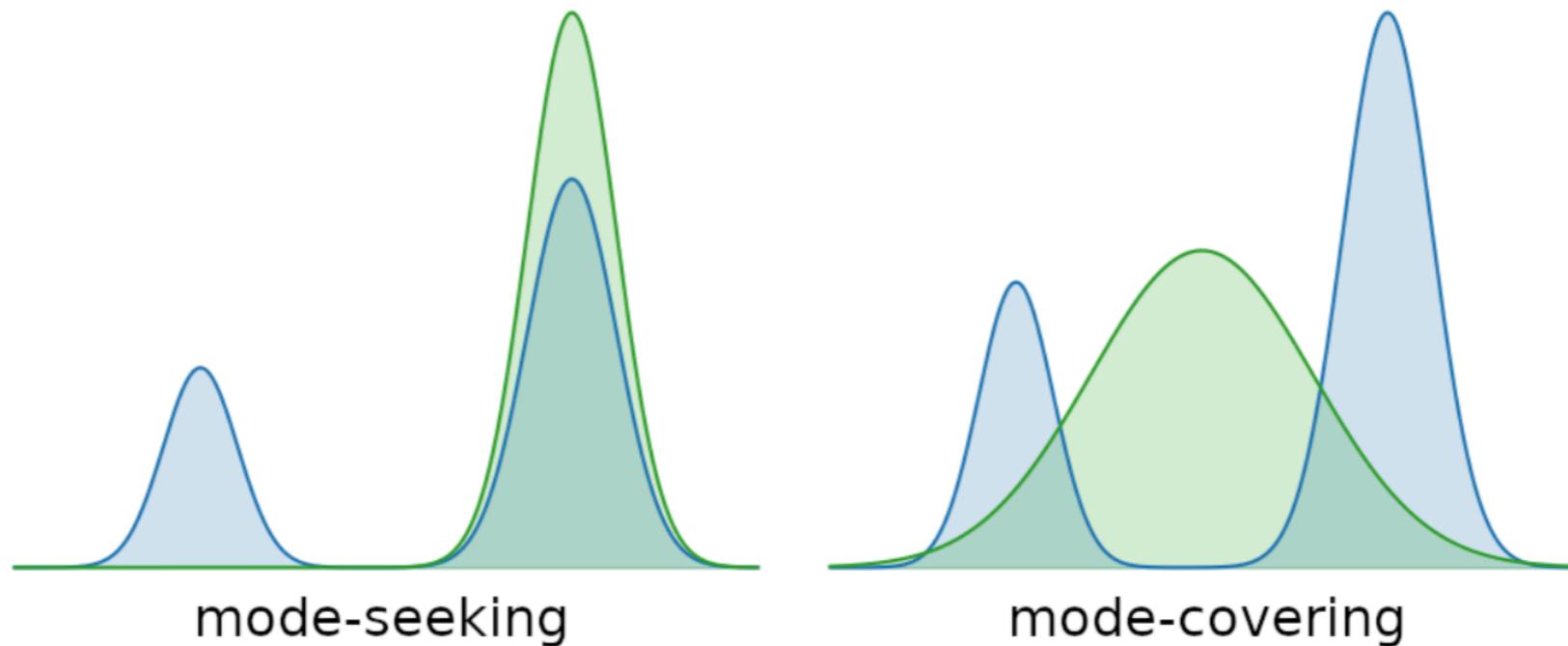
If there is no such requirement, the probability mass can be focused on a subset of examples, but then some parts of the distribution will be ignored by the model (**mode-seeking**).



Mode-covering vs. mode-seeking behaviour

Likelihood-based models are usually mode-covering. This is a consequence of the fact that they are fit by maximising the joint likelihood of the data.

Adversarial models on the other hand are typically mode-seeking.



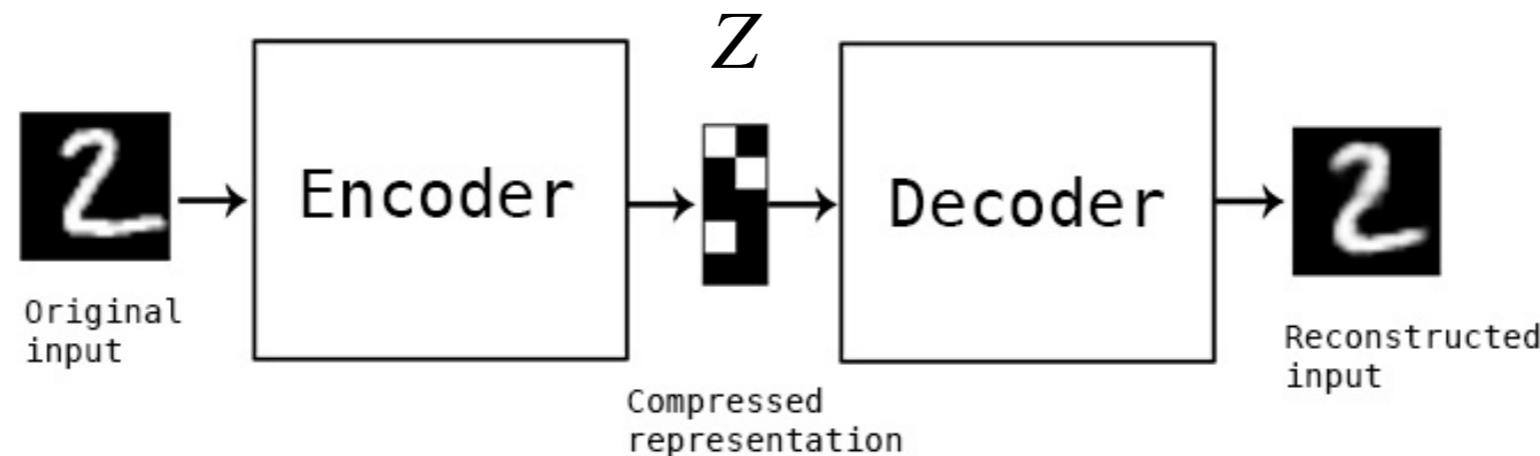
Algorithms

Autoencoders

Suppose we have only a set of unlabeled training examples x_1, x_2, x_3, \dots , where $x_i \in \mathbb{R}^n$.

An autoencoder neural network is an **unsupervised** learning algorithm that applies backpropagation and uses a loss function that is optimal when setting the target values to be equal to the inputs, $y_i = x_i$.

To build an autoencoder, you need three things: an **encoding function**, a **decoding function**, and a **distance function** between the amount of information loss between the compressed representation of your data and the decompressed representation.



Simple Autoencoder for MNIST

```
1 # Source: Adapted from https://blog.keras.io/building-autoencoders-in-keras.html
2
3 from tensorflow.keras.layers import Input, Dense
4 from tensorflow.keras.models import Model
5
6 # this is the size of our encoded representations
7 encoding_dim = 32 # 32 floats -> compression of factor 24.5,
8 # assuming the input is 784 floats
9
10 input_img = Input(shape=(784,))
11
12 # encoded representation of the input
13 encoding_layer = Dense(encoding_dim, activation='relu')
14 encoded = encoding_layer(input_img)
15
16 # lossy reconstruction of the input
17 decoding_layer = Dense(784,
18 activation='sigmoid')
19 decoded = decoding_layer(encoded)
20
21 # this model maps an input to its reconstruction
22 autoencoder = Model(input_img, decoded)
```

Input



Output

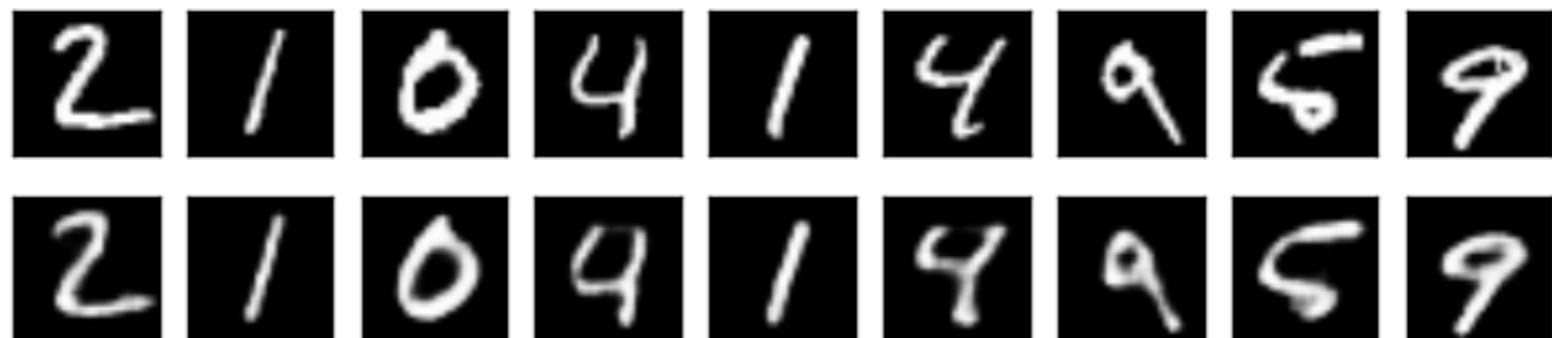


Adding depth and sparsity constraints

```
1 from tensorflow.keras import regularizers
2 from tensorflow.keras import optimizers
3 from tensorflow.keras.regularizers import l2, activity_l1
4 from tensorflow.keras.layers import Input, Dense
5 from tensorflow.keras.models import Model
6
7 input_img = Input(shape=(784,))
8 encoded = Dense(128, activation='relu')(input_img)
9 encoded = Dense(64, activation='relu')(encoded)
10 encoded = Dense(32, activation='relu')(encoded) Z
11 decoded = Dense(64, activation='relu')(encoded)
12 decoded = Dense(128, activation='relu')(decoded)
13 decoded = Dense(784, activation='sigmoid')(decoded)
14
15 autoencoder = Model(input_img, decoded)
16 autoencoder.compile(optimizer='adadelta',
17                      loss='binary_crossentropy',
18                      activity_regularizer=regularizers.l1(10e-5))
```

Convolutional Autoencoders

```
5 input_img = Input(shape=(28, 28, 1))
6
7 x = Conv2D(16, 3, 3, activation='relu', border_mode='same')(input_img)
8 x = MaxPooling2D((2, 2), border_mode='same')(x)
9 x = Conv2D(8, 3, 3, activation='relu', border_mode='same')(x)
10 x = MaxPooling2D((2, 2), border_mode='same')(x)
11 x = Conv2D(8, 3, 3, activation='relu', border_mode='same')(x)
12 encoded = MaxPooling2D((2, 2), border_mode='same')(x) Z
13
14 # at this point the representation is (4, 4, 8) i.e. 128-dimensional
15
16 x = Conv2D(8, 3, 3, activation='relu', border_mode='same')(encoded)
17 x = UpSampling2D((2, 2))(x)
18 x = Conv2D(8, 3, 3, activation='relu', border_mode='same')(x)
19 x = UpSampling2D((2, 2))(x)
20 x = Conv2D(16, 3, 3, activation='relu')(x)
21 x = UpSampling2D((2, 2))(x)
22 decoded = Conv2D(1, 3, 3, activation='sigmoid', border_mode='same')(x)
23
24 # at this point the representation is (28, 28, 1) i.e. 784-dimensional
25
26 autoencoder = Model(input_img, decoded)
27 autoencoder.compile(optimizer='adadelta', loss='binary_crossentropy')
```



Denoising Autoencoders

Noise as a **regularization** strategy



Input



Output



Variational Autoencoders

A **variational autoencoder** is an autoencoder that adds **probabilistic constraints** on the representations (Z) being learned.

When using probabilistic models, the compressed representation is called **latent variable model**.

So, instead of learning a function this model is **learning a probabilistic distribution function that models your data**.

Variational Autoencoders

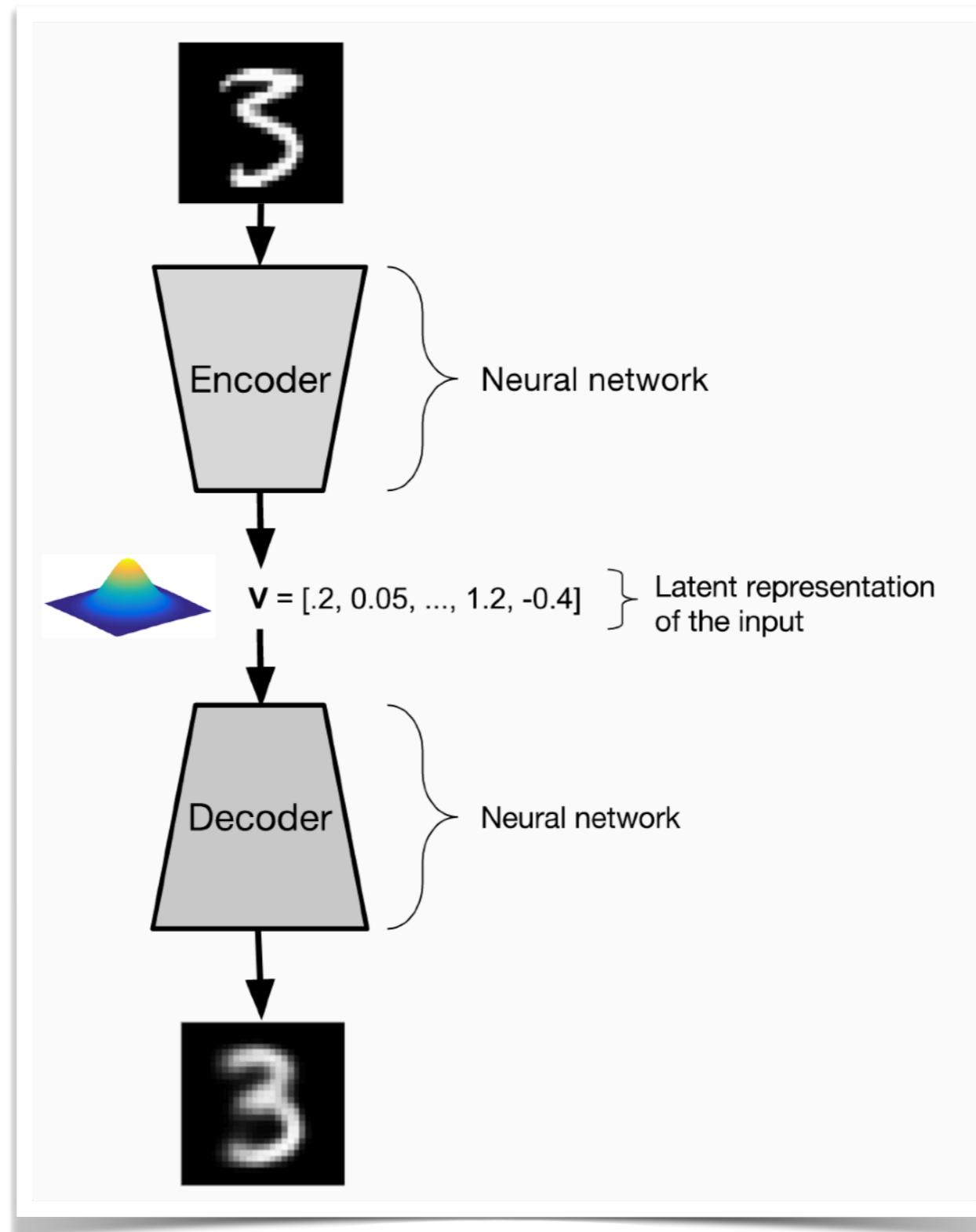
Why?

Standard autoencoders are not suited to work as a generative model. **If you pick a random value for your decoder you won't get necessarily a good reconstruction:** the value can far away from any previous value the network has seen before! That's why attaching a probabilistic model to the compressed representation is a good idea!

For sake of simplicity, let's use a **standard normal distribution (V)** to represent the distribution of inputs the decoder will receive.

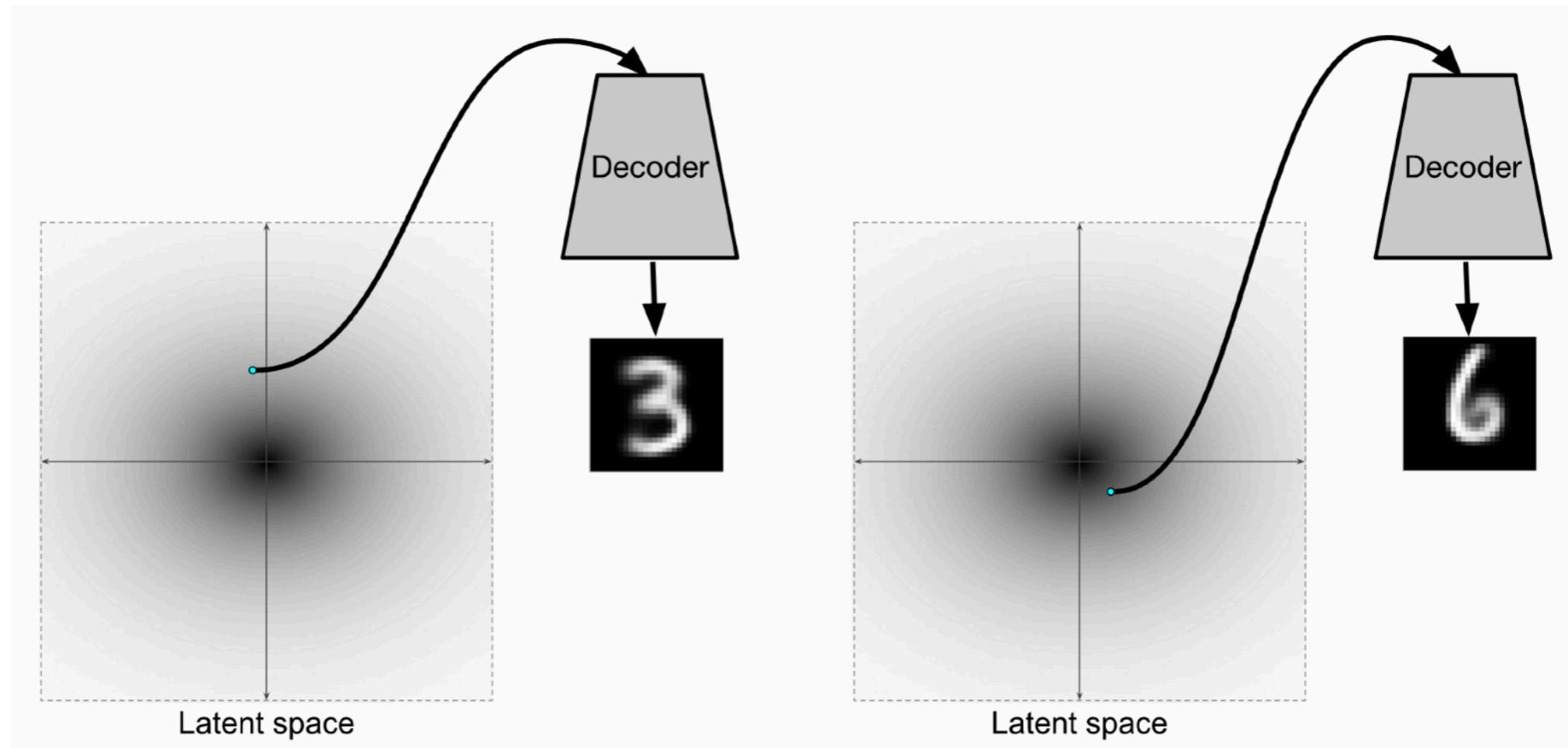
Variational Autoencoders

Architecture of a variational autoencoder (VAE)



Variational Autoencoders

We want the decoder to take any point taken from a standard normal distribution to return a reasonable element of our dataset.

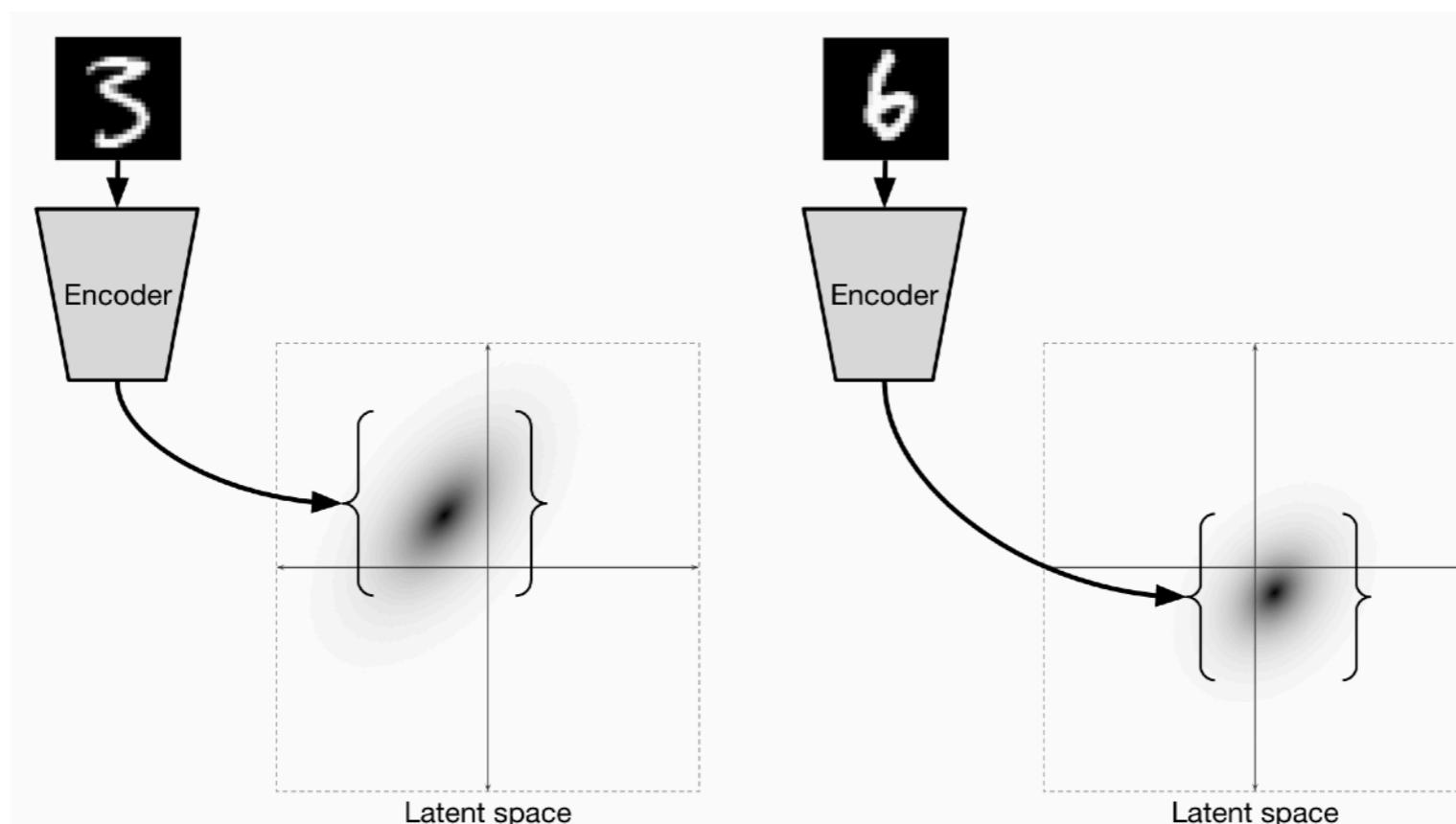


Variational Autoencoders

Let's consider the **encoder role** in this architecture.

In a traditional autoencoder, the encoder model takes a sample from data and returns a single point in the latent space, which is then passed to the decoder.

In VAE the **encoder instead produces (the parameters of) a probability distribution** in the latent space:

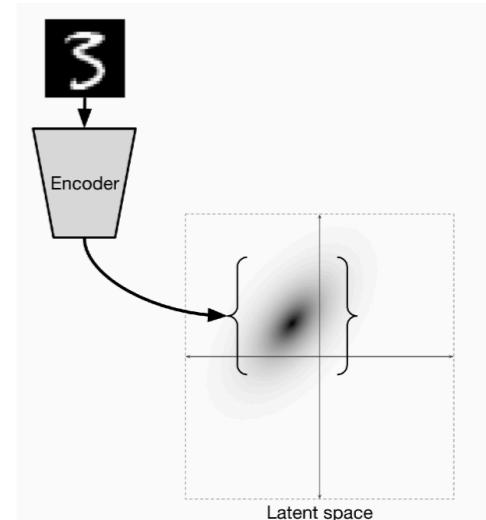


Variational Autoencoders

First, let's implement the encoder net, which takes input X and outputs two things: $\mu(X)$ and $\sigma(X)$, the parameters of the Gaussian. Our encoder will be a neural net with one hidden layer.

Our latent variables are two dimensional, so that we could easily visualize it.

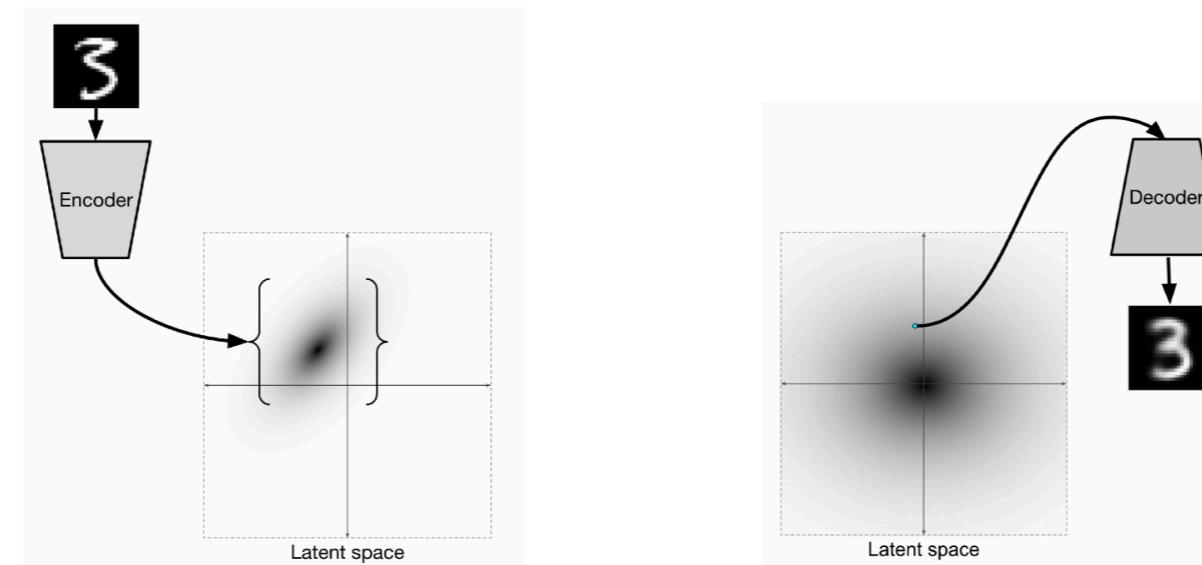
```
18 # encoder
19 inputs = Input(shape=(784,))
20 h_q = Dense(512, activation='relu')(inputs)
21 mu = Dense(2, activation='linear')(h_q)
22 log_sigma = Dense(2, activation='linear')(h_q)
    Diagonal covariance matrix
```



Variational Autoencoders

Up to now we have an encoder that takes images and produce (the parameters of) a pdf in the latent space.

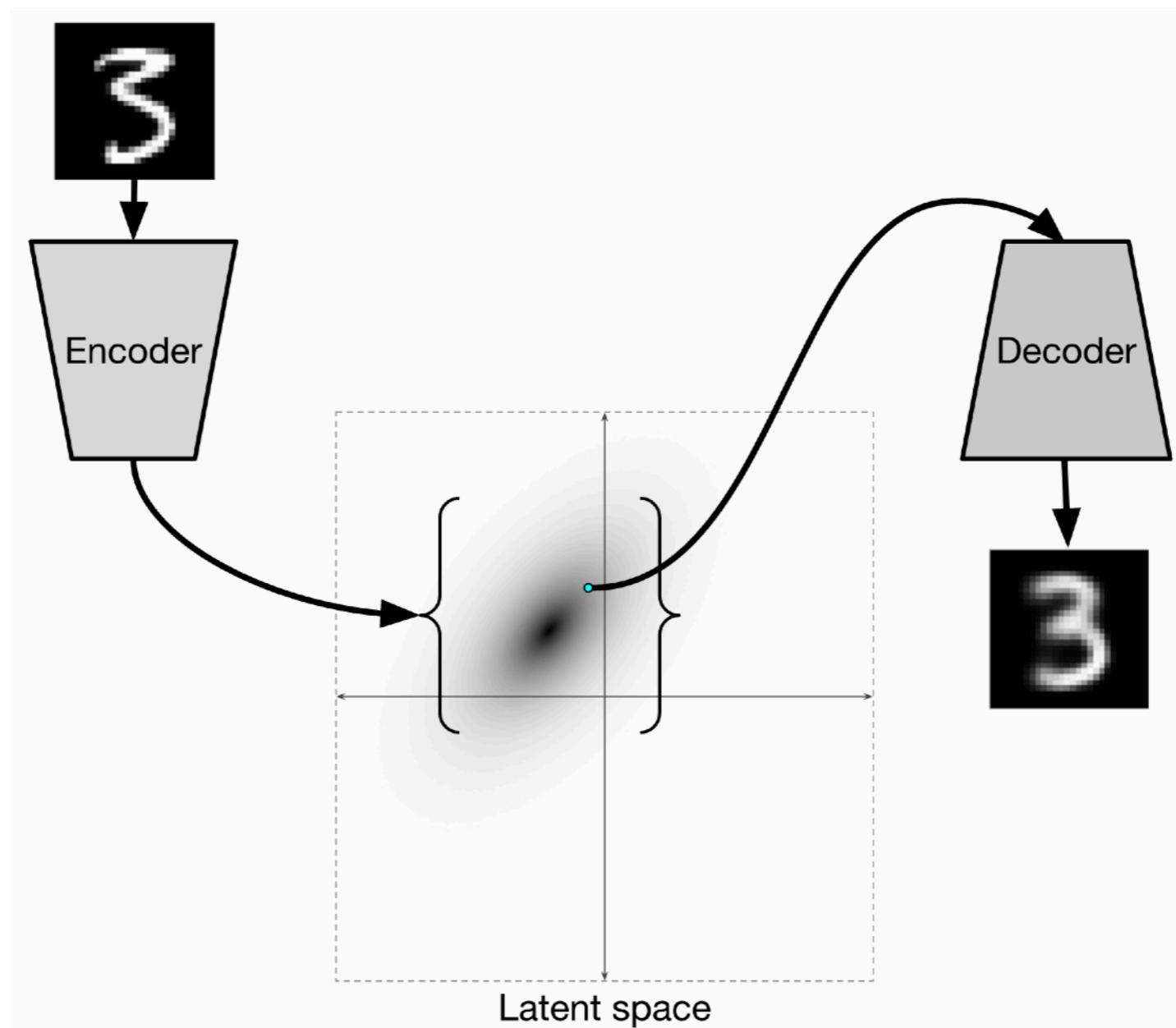
The decoder takes points in the latent space and return reconstructions.



How do we connect both models?

By sampling from the produced distribution!

Variational Autoencoders



Variational Autoencoders

Sampling from the produced distribution

To this end we will implement a random variate **reparameterization**: the substitution of a random variable by a deterministic transformation of a simpler random variable.

There are several methods by which non-uniform random numbers, or random variates, can be generated. The most popular methods are the **one-liners**, which give us the simple tools to generate random variates in one line of code, following the classic paper by Luc Devroye (*Luc Devroye, Random variate generation in one line of code, Proceedings of the 28th conference on Winter simulation, 1996*).

Variational Autoencoders

In the case of a Gaussian, we can use the following algorithm:

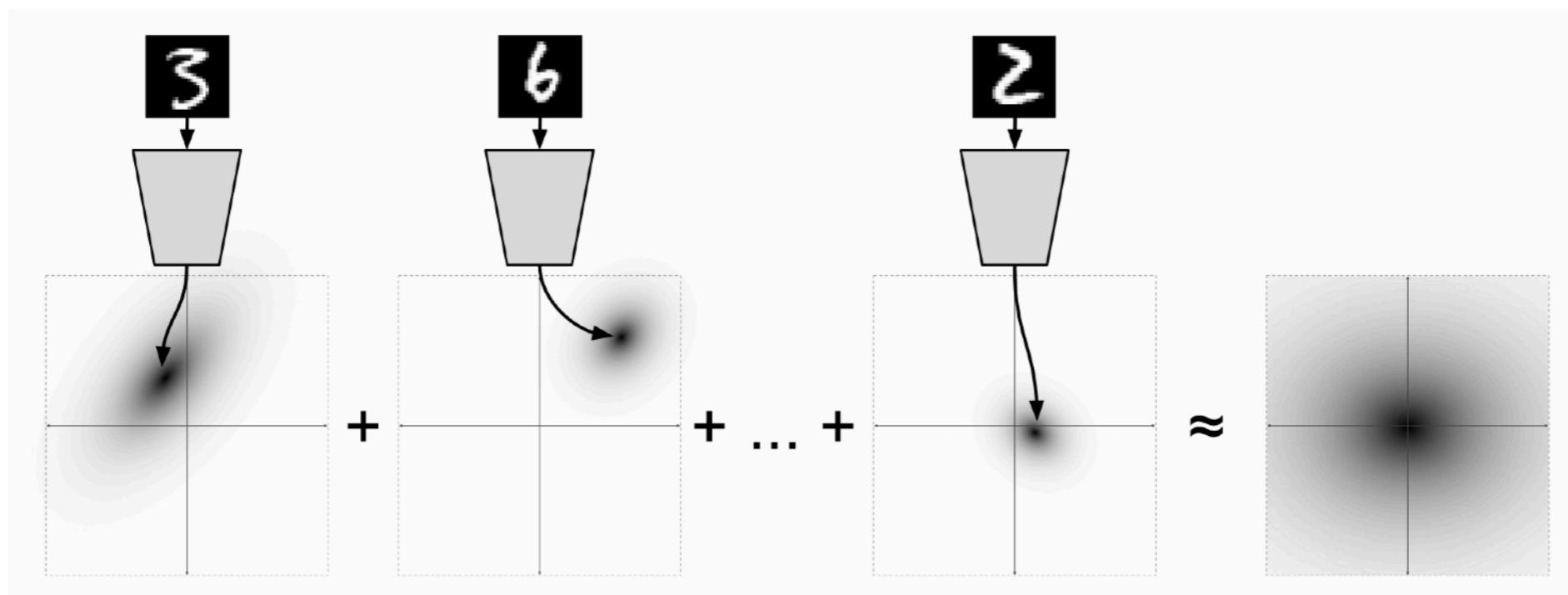
- Generate $\epsilon \sim \mathcal{N}(0; 1)$.
- Compute a sample from $\mathcal{N}(\mu; RR^T)$ as $\mu + R\epsilon$.

```
1 def sample_z(args):
2     mu, log_sigma = args
3     eps = K.random_normal(shape=(m, n_z), mean=0., std=1.)
4     return mu + K.exp(log_sigma / 2) * eps
5
6 # Sample z
7 z = Lambda(sample_z)([mu, log_sigma])
```

Variational Autoencoders

From this model, we can do three things: **reconstruct inputs, encode inputs into latent variables, and generate data from latent variable.**

In order to be coherent with our previous definitions, we must assure that points sampled from the latent space fit a standard normal distribution, but the encoder is producing non standard normal distributions. So, we must add a constraint for getting something like this:



Variational Autoencoders

In order to impose this constraint we can add a term to the loss function:

VAE's loss function comprises a **Reconstruction** error, and a **KL-divergence** error used to model the networks' objectives: $L_{MSE} + L_{KL}$.

The final loss is a weighted sum of both losses.

The reconstruction loss in VAE is similar to the Loss we used in the traditional Autoencoder i.e. **Mean-Squared-Error**

$$L_{MSE}(\theta, \phi) = \frac{1}{N} \sum_{i=1}^N (x_i - D_\phi(E_\theta(x_i)))^2$$

where θ and ϕ are the parameters of the encoder and decoder, respectively.

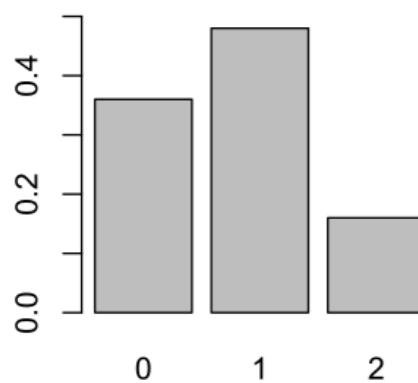
Variational Autoencoders

The Kullback–Leibler divergence is a measure of how one probability distribution diverges from a second expected probability distribution.

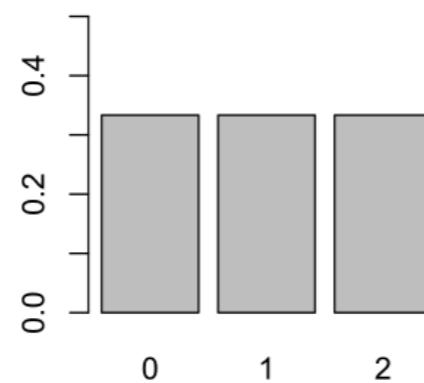
For discrete probability distributions P and Q , the Kullback–Leibler divergence from Q to P is defined to be

$$D_{\text{KL}}(Q||P) = \sum_i Q(i) \log \frac{Q(i)}{P(i)}$$

Distribution P
Binomial with $p = 0.4$, $N = 2$



Distribution Q
Uniform with $p = 1/3$



x	0	1	2
Distribution $P(x)$	$\frac{9}{25}$	$\frac{12}{25}$	$\frac{4}{25}$
Distribution $Q(x)$	$\frac{1}{3}$	$\frac{1}{3}$	$\frac{1}{3}$

$$D_{\text{KL}}(P \parallel Q) = \sum_{x \in \mathcal{X}} P(x) \ln \left(\frac{P(x)}{Q(x)} \right)$$

$$= \frac{9}{25} \ln \left(\frac{9/25}{1/3} \right) + \frac{12}{25} \ln \left(\frac{12/25}{1/3} \right) + \frac{4}{25} \ln \left(\frac{4/25}{1/3} \right)$$

$$= \frac{1}{25} (32 \ln(2) + 55 \ln(3) - 50 \ln(5)) \approx 0.0852996$$

$$D_{\text{KL}}(Q \parallel P) = \sum_{x \in \mathcal{X}} Q(x) \ln \left(\frac{Q(x)}{P(x)} \right)$$

$$= \frac{1}{3} \ln \left(\frac{1/3}{9/25} \right) + \frac{1}{3} \ln \left(\frac{1/3}{12/25} \right) + \frac{1}{3} \ln \left(\frac{1/3}{4/25} \right)$$

$$= \frac{1}{3} (-4 \ln(2) - 6 \ln(3) + 6 \ln(5)) \approx 0.097455$$

A relative entropy of 0 indicates that the two distributions in question have identical quantities of information.

Variational Autoencoders

We train our VAE to minimize the KL divergence between the encoder's distribution $Q_\phi(Z|X)$ and a predefined $P(Z)$.

In VAE, $P(Z)$ is defined to follow a standard or unit $P(Z) = \mathcal{N}(0,1)$. If the encoder outputs encoding Z far from a standard normal distribution, KL-divergence loss will penalize it more.

Trick: While calculating the KL-divergence we choose to map the variance to the logarithm of the variance. By taking the **logarithm of the variance**, we force the network to have the output range of the natural numbers rather than just positive values (**variances** would only have positive values). This allows for smoother representations for the latent space.

In this special case, the KL divergence has the closed form:

$$L_{KL}(\theta, \phi) = -0.5 \sum_{i=1}^N 1 + \log \sigma_i^2 - \mu_i^2 - e^{\log \sigma_i^2}$$

```
# D_KL(Q(z|x) || P(z|x)); calculate in closed form as both dist. are Gaussian
kl_loss = 0.5 * tf.reduce_sum(tf.exp(z_logvar) + z_mu**2 - 1. - z_logvar, 1)
# VAE loss
vae_loss = tf.reduce_mean(recon_loss + kl_loss)
```

How do we train a Variational Autoencoder?

We have a sampling operation (stochastic node) in the computational graph!

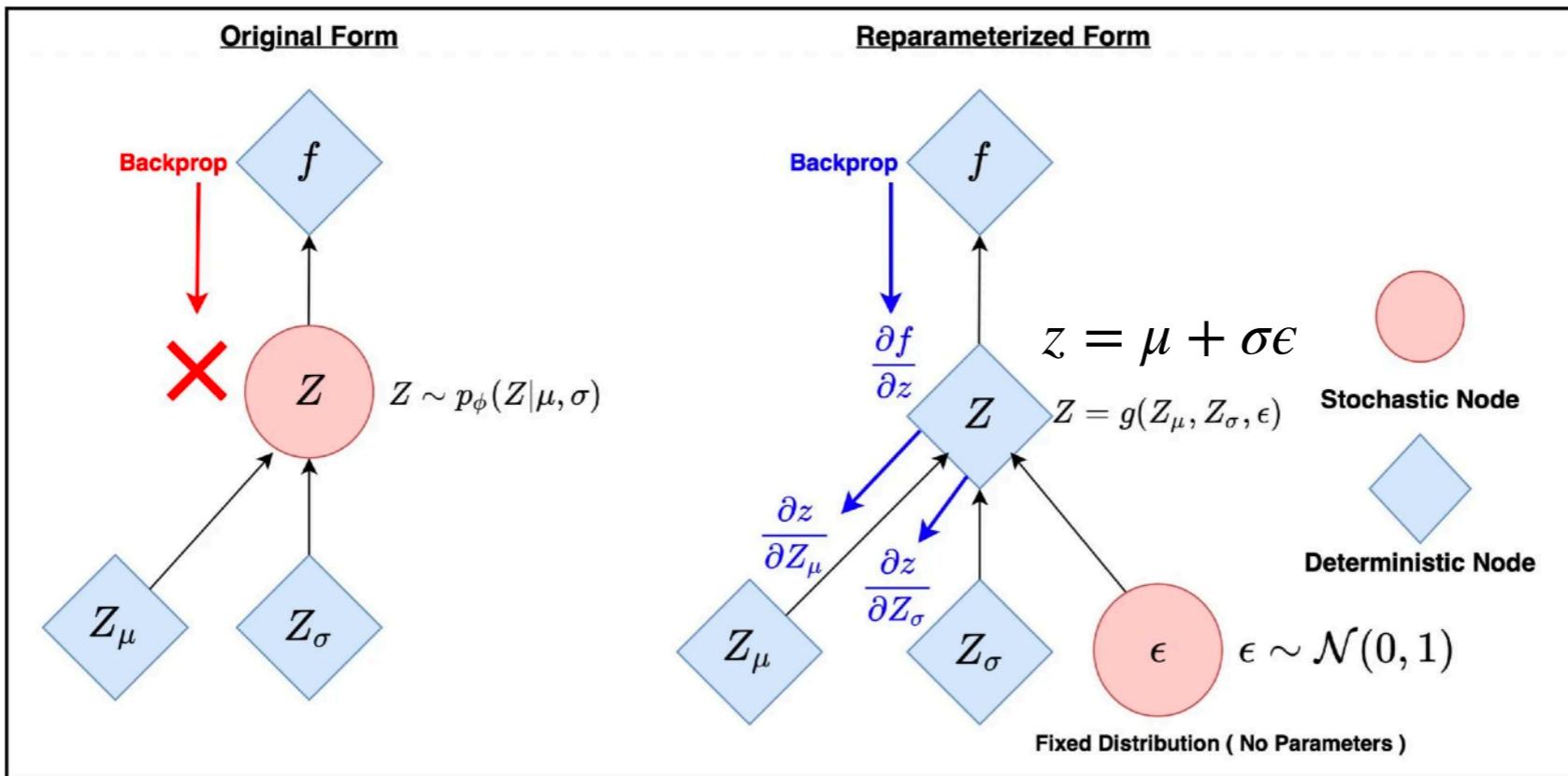
In fact this is not a problem!

By using the one-liner method for sampling **we have expressed the latent distribution in a way that its parameters are factored out of the parameters of the random variable** so that **backpropagation can be used to find the optimal parameters of the latent distribution.**

For this reason this method is called **reparametrization trick**.

By using this trick we can train end-to-end a VAE with backpropagation.

Reparametrization trick

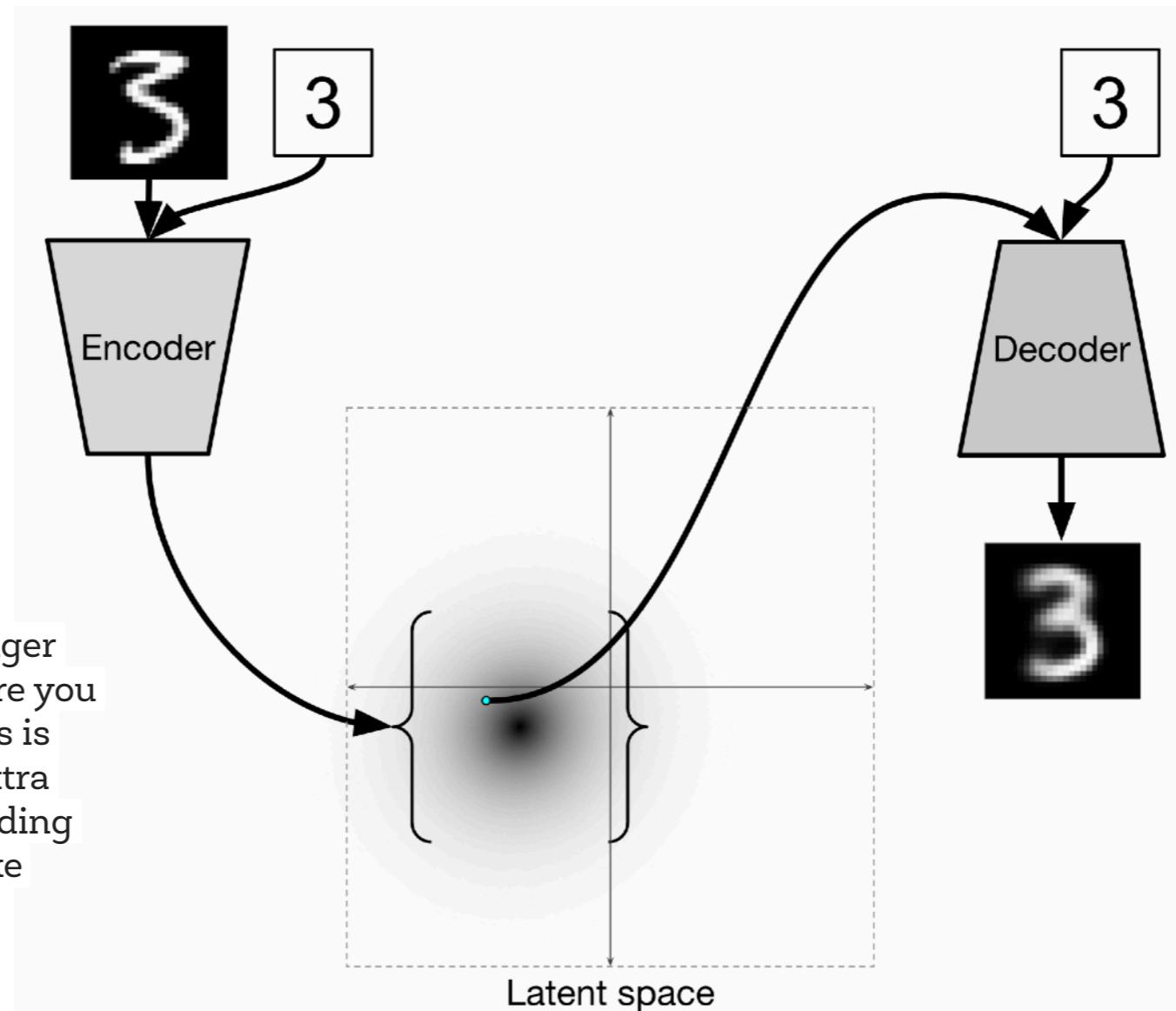


```
1 def sample_z(args):
2     mu, log_sigma = args
3     eps = K.random_normal(shape=(m, n_z), mean=0., std=1.)
4     return mu + K.exp(log_sigma / 2) * eps
5
6 # Sample z
7 z = Lambda(sample_z)([mu, log_sigma])
```

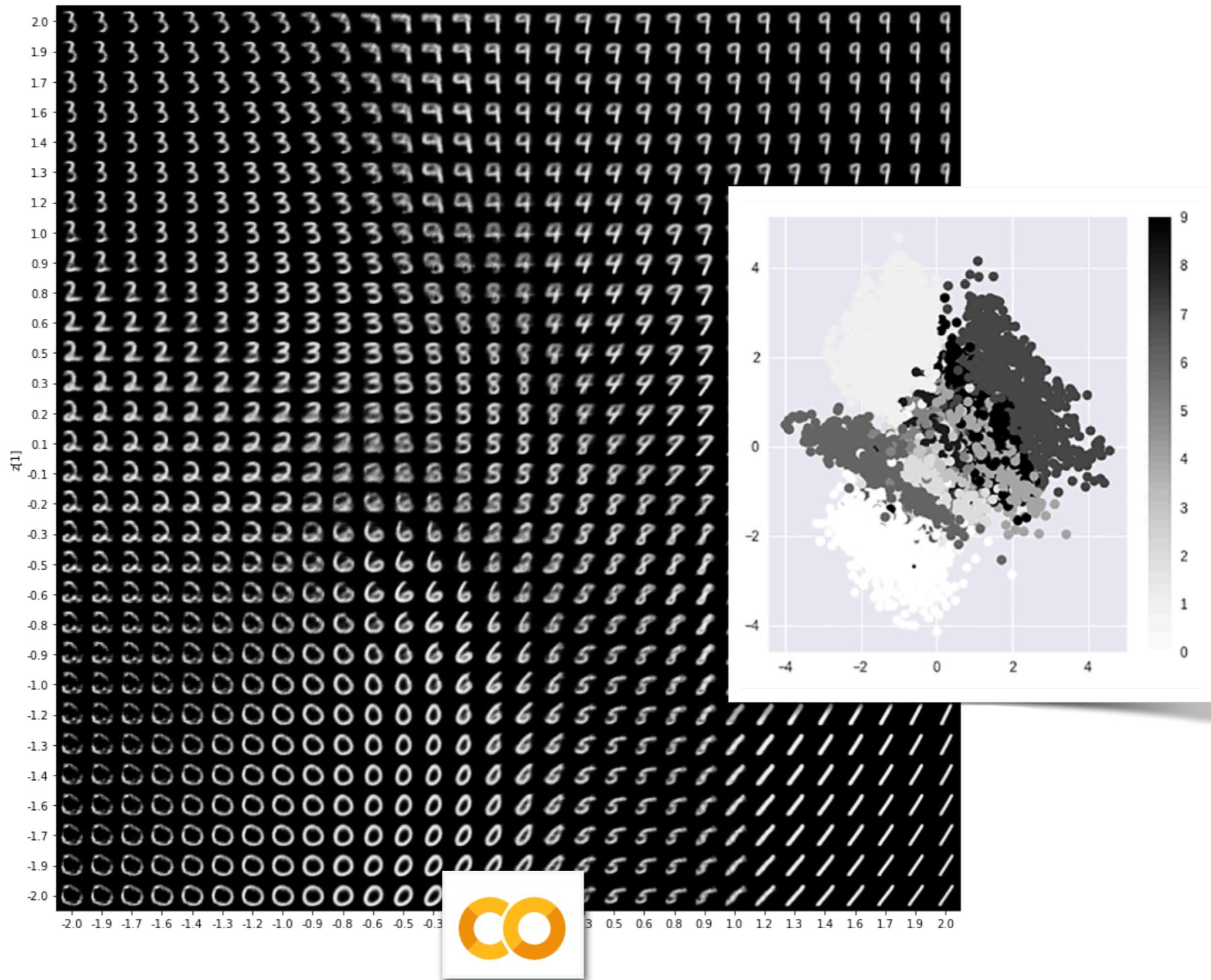
Conditional Variational Autoencoder

What about producing specific number instances on demand?

We can do this by adding an extra input (as a one-hot encoding) to both the encoder and the decoder:



Variational Autoencoder in Keras (d=2)



Is it enough?

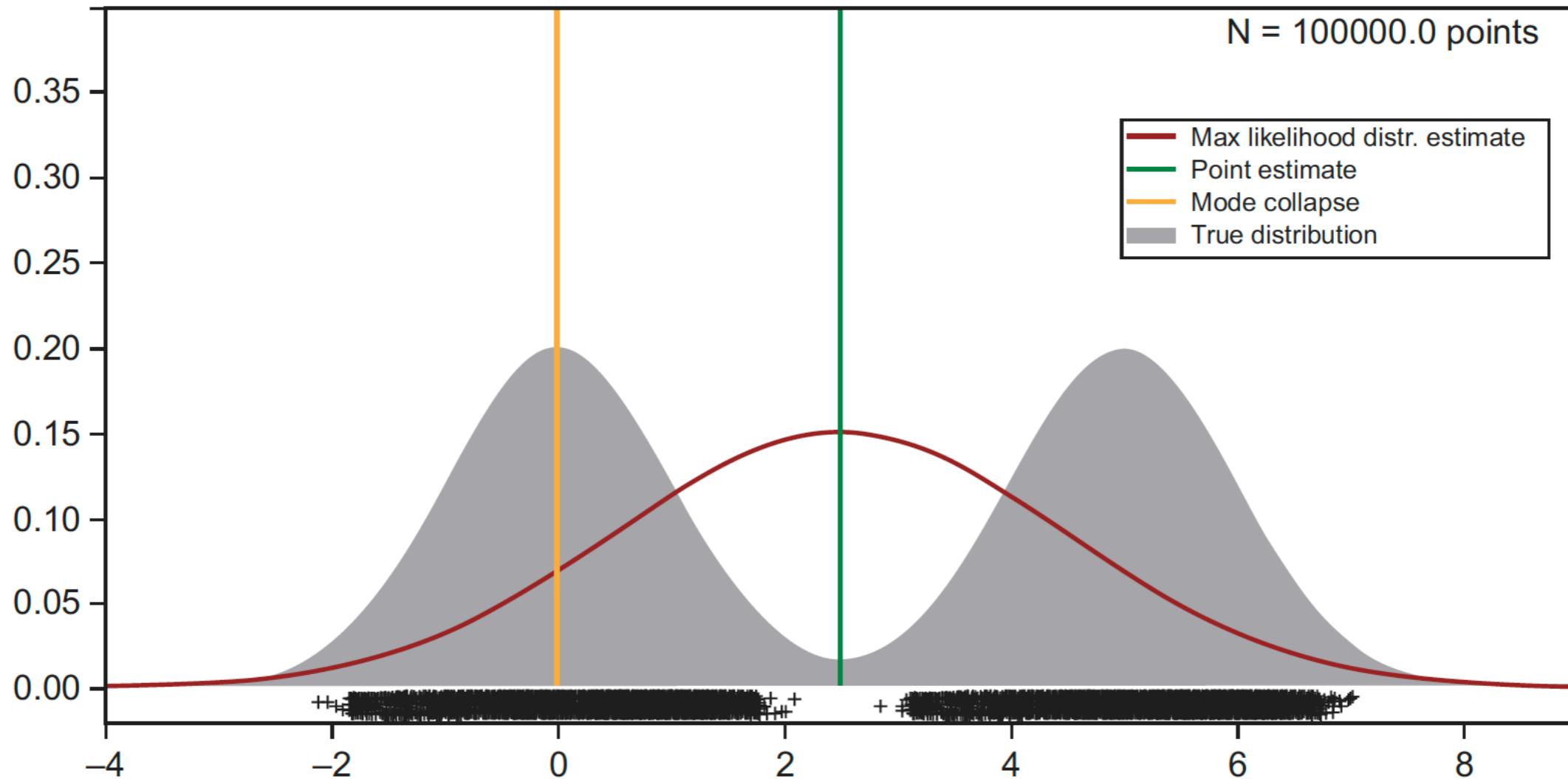
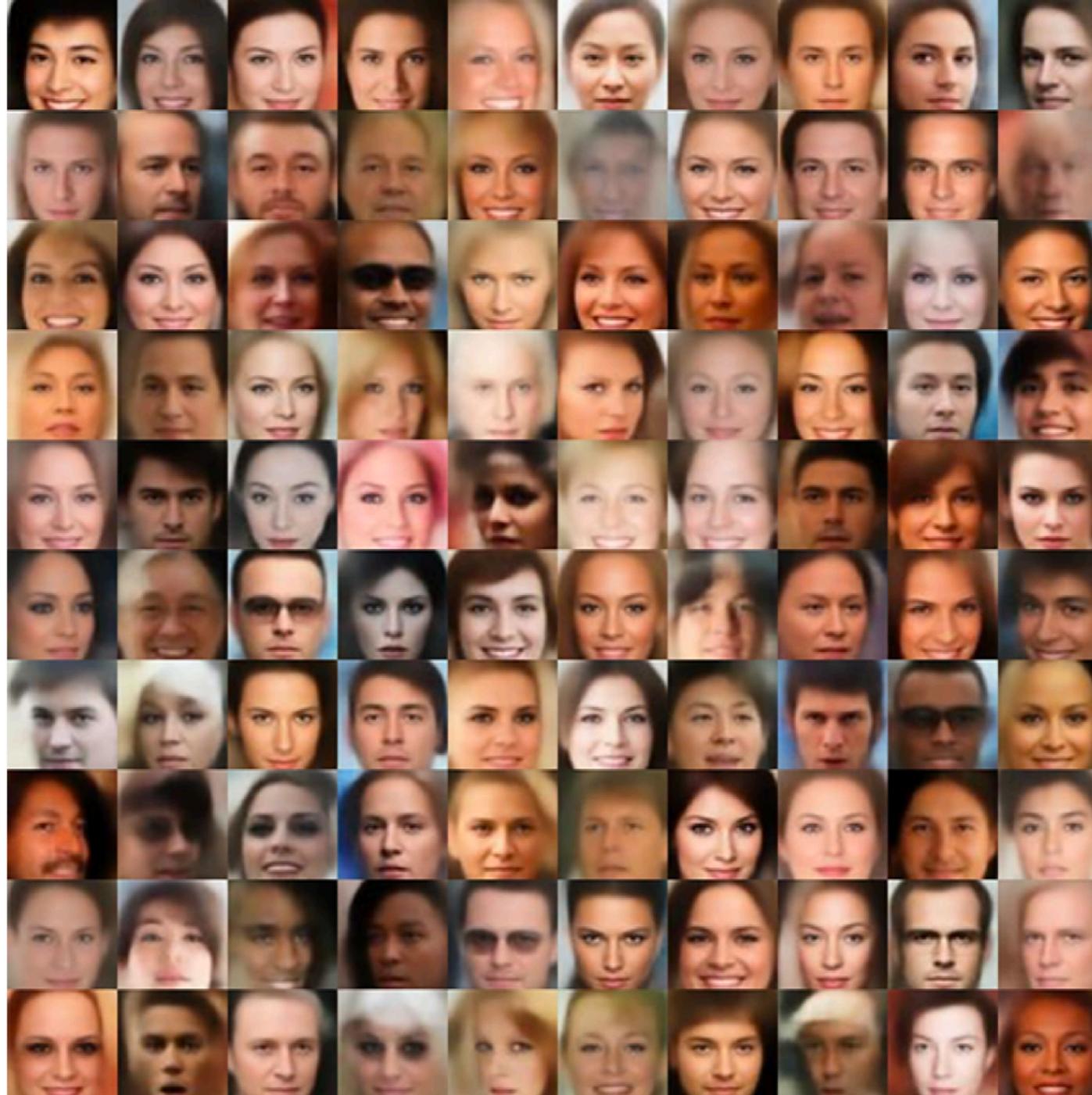


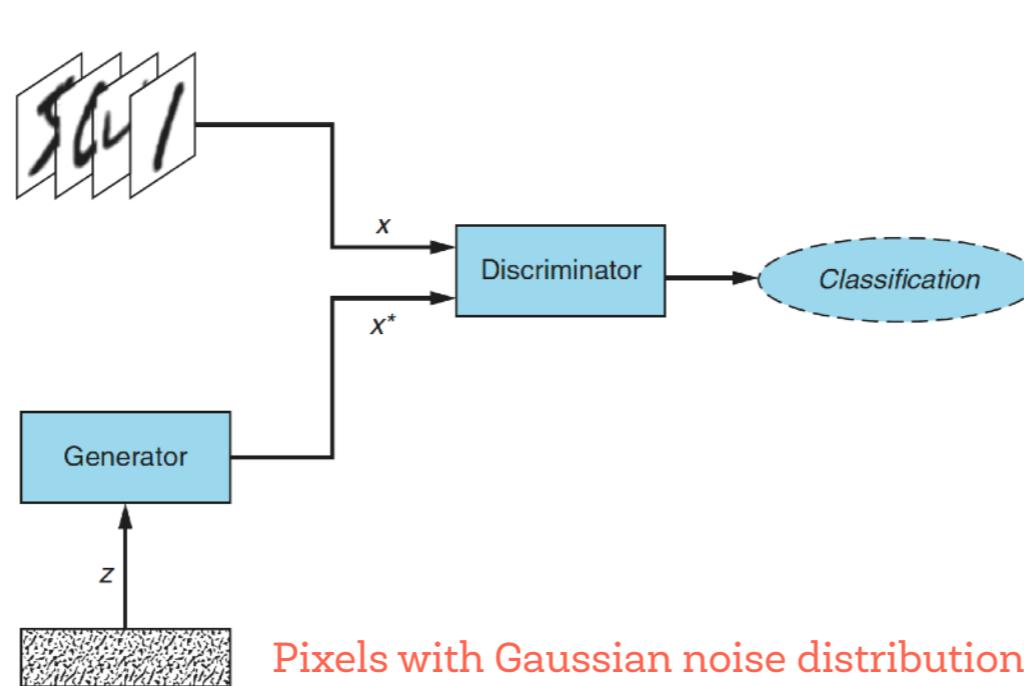
Figure 2.8 Maximum likelihood, point estimates, and true distributions. The gray (theoretical) distribution is bimodal rather than having a single mode. But because we have assumed this, our model is catastrophically wrong. Alternatively, we can get mode collapse, which is worth keeping in mind for chapter 5. This is especially true when we are using flavors of the KL, such as the VAE or early GANs.

Is it enough?



The VAE picks the safe path and makes the background blurry by choosing a “safe” pixel value, which minimizes the loss, but does not provide good images.

Generative Adversarial Networks



GAN training algorithm

For each training iteration **do**

- 1 Train the Discriminator:
 - a Take a random mini-batch of real examples: x .
 - b Take a mini-batch of random noise vectors z and generate a mini-batch of fake examples: $G(z) = x^*$.
 - c Compute the classification losses for $D(x)$ and $D(x^*)$, and backpropagate the total error to update $\theta^{(D)}$ to *minimize* the classification loss.
- 2 Train the Generator:
 - a Take a mini-batch of random noise vectors z and generate a mini-batch of fake examples: $G(z) = x^*$.
 - b Compute the classification loss for $D(x^*)$, and backpropagate the loss to update $\theta^{(G)}$ to *maximize* the classification loss.

End for

GAN problems

Many GAN models suffer the following major problems:

- **Non-convergence**: the model parameters oscillate, destabilize and never converge,
- **Mode collapse**: the generator collapses which produces limited varieties of samples,
- **Diminished gradient**: the discriminator gets too successful that the generator gradient vanishes and learns nothing,
- Unbalance between the generator and discriminator causing **overfitting**, &
- Highly sensitive to the **hyperparameter** selections.

Solutions:

- Advanced loss functions.
- Advanced training strategies.
- Etc.

GAN problems

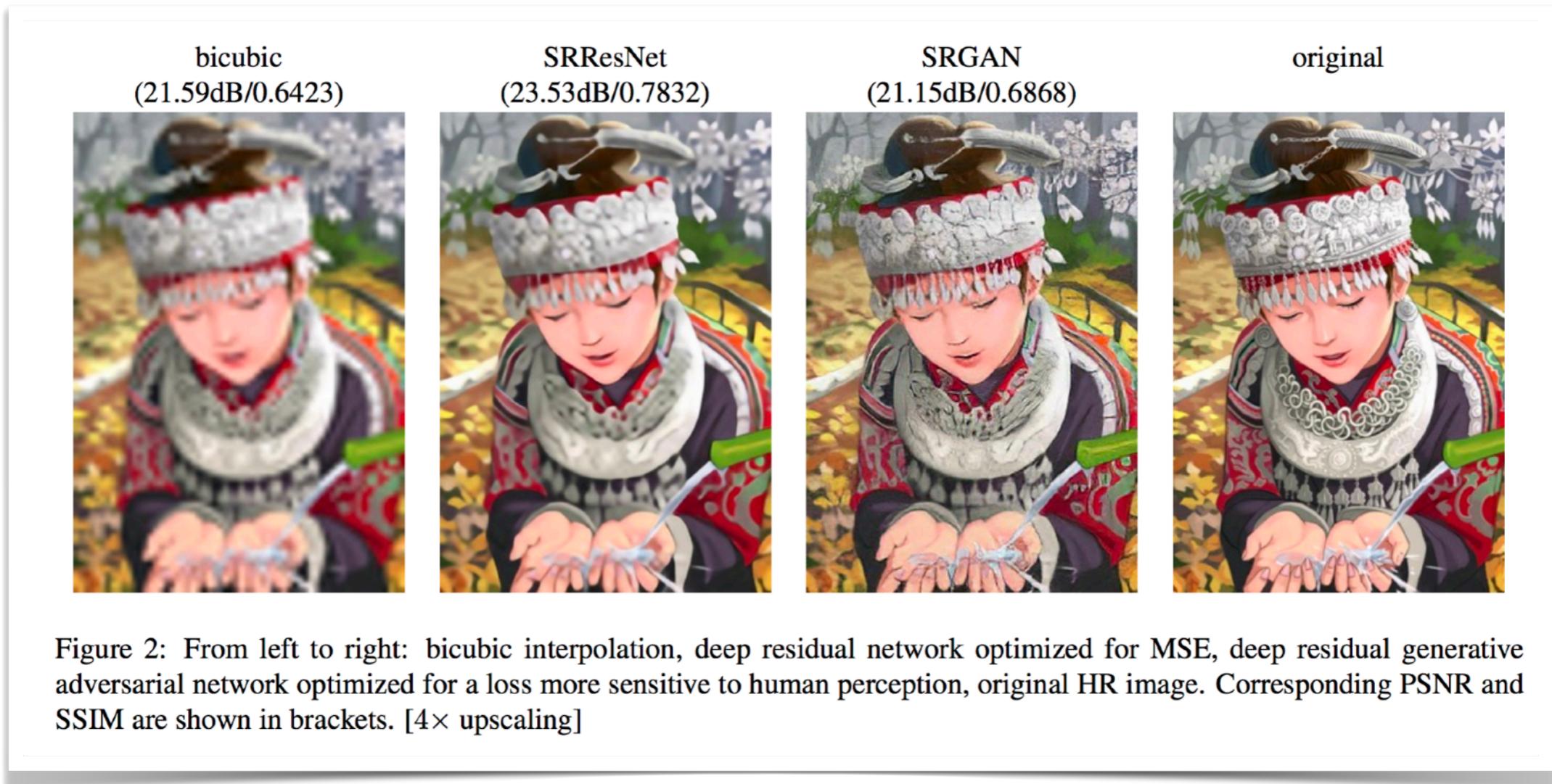
The objective evaluation of GAN generator models remains an open problem.

As such, GAN generator models are evaluated based on the **quality of the images generated**, often in the context of the target problem domain.

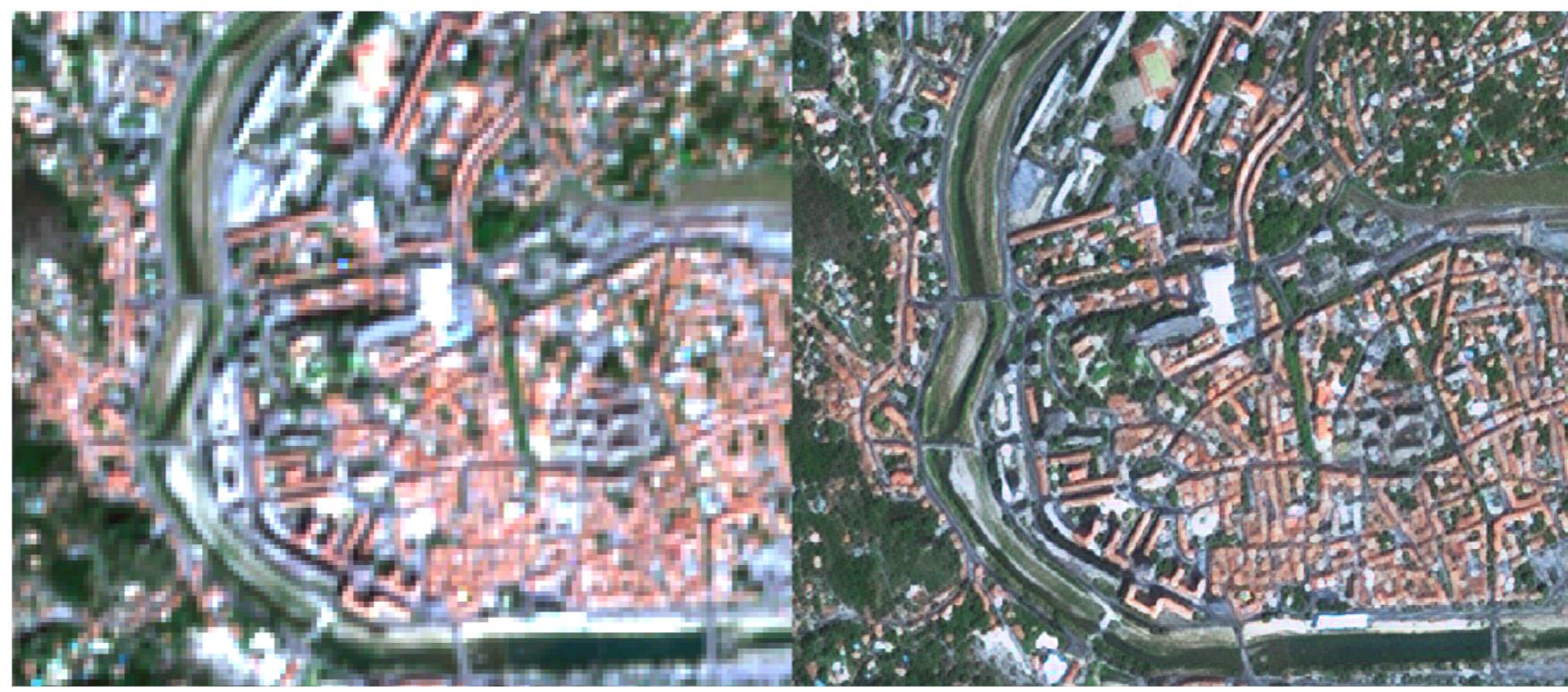
Twenty-four quantitative techniques for evaluating GAN generator models are listed below.

1. Average Log-likelihood
2. Coverage Metric
3. Inception Score (IS)
4. Modified Inception Score (m-IS)
5. Mode Score
6. AM Score
7. Frechet Inception Distance (FID)
8. Maximum Mean Discrepancy (MMD)
9. The Wasserstein Critic
10. Birthday Paradox Test
11. Classifier Two-sample Tests (C2ST)
12. Classification Performance
13. Boundary Distortion
14. Number of Statistically-Different Bins (NDB)
15. Image Retrieval Performance
16. Generative Adversarial Metric (GAM)
17. Tournament Win Rate and Skill Rating
18. Normalized Relative Discriminative Score (NRDS)
19. Adversarial Accuracy and Adversarial Divergence
20. Geometry Score
21. Reconstruction Error
22. Image Quality Measures (SSIM, PSNR and Sharpness Difference)
23. Low-level Image Statistics
24. Precision, Recall and F1 Score

Applications



Applications



Enhancement of 10m Sentinel-2 images at 1.5m

Applications

Hugging Face Models Datasets Spaces Resources Solutions Pricing Log In Sign Up

Spaces: akhaliq/AnimeGANv2 like 363 Running

App Files and versions

AnimeGANv2

Gradio Demo for AnimeGanv2 Face Portrait. To use it, simply upload your image, or click one of the examples to load them. Read more at the links below. Please use a cropped portrait picture for best results similar to the examples below.

img  Output  0.23s

version

version 1 (▲ stylization, ▼ robustness)

version 2 (▲ robustness, ▼ stylization)

Screenshot

Clear Submit

