# Convolutional NNs and Generative Models

# Friday 08h00-09h00

Géraldine Conti, Matthew Vowels, Mykhailo Vladymyrov
Bern Winter School on Machine Learning 2024, Mürren

1

# Convolutional Neural Networks

- Specialized Neural Network for data arranged on a **continuous** grid
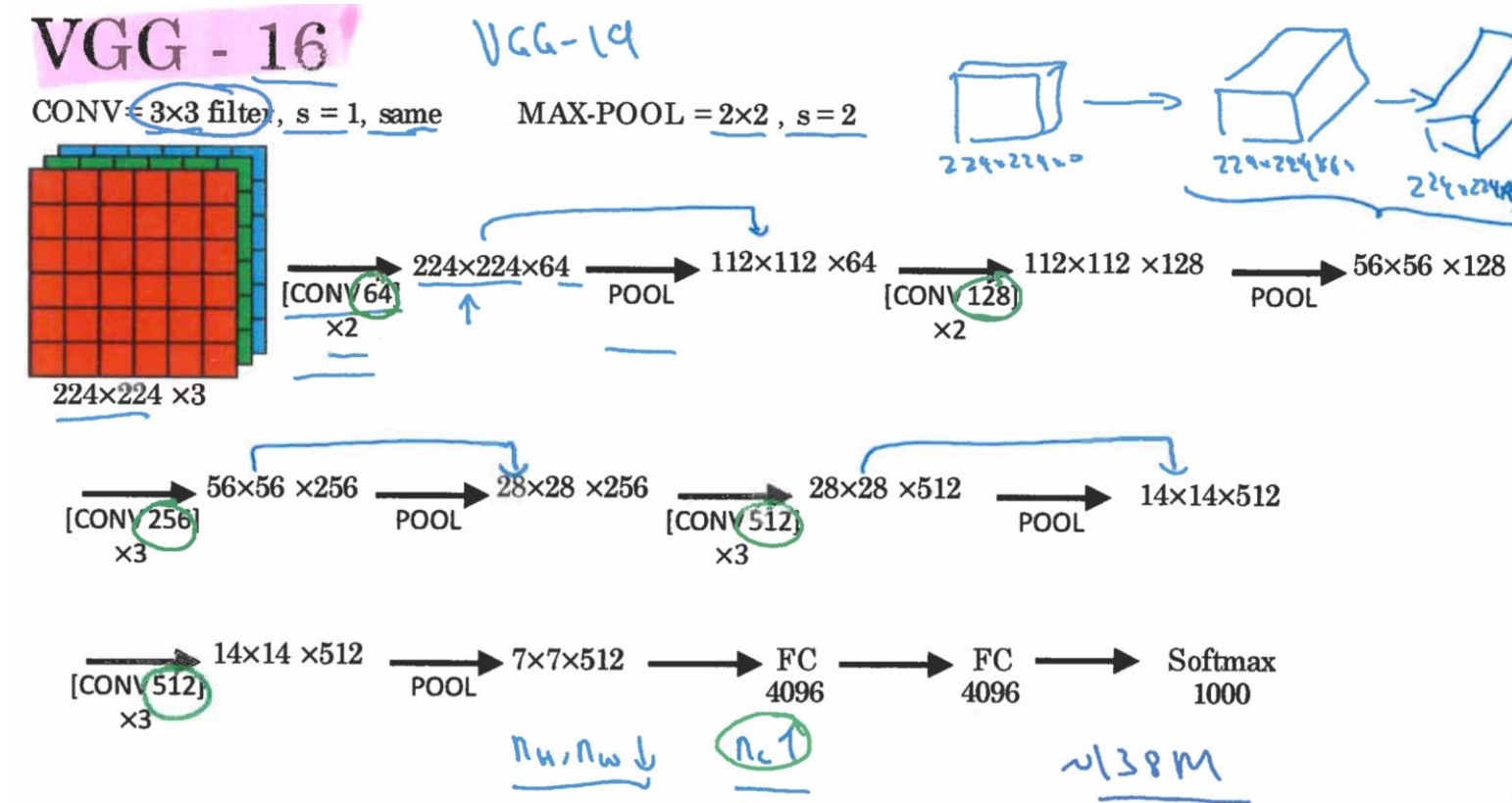  - Images, signal waveform samples, …
  - DNA sequences
  - …

# Types of layer in CNN

- **Convolution** (CONV)
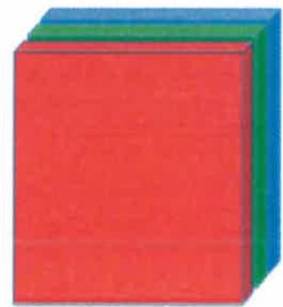
- **Pooling** (POOL)

- **Fully connected** (FC)

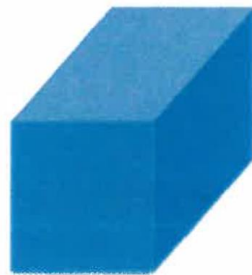- *Usually multiple CONV layers followed by a POOL layer, and FC layers in the last few layers*

# Convolution Layer (CONV)

- **Convolution** transforms an input volume into an output volume of different size, also called feature map

- **Filter kernels** are used to detect features (for example, edge detection in 1st hidden layer)



$n \times n$ image     $f \times f$ filter     out $= n - f + 1$

# Convolution Layer (CONV)

| | | | | Input Volume (+pad 1) (7x7x3) | | | | | | | | Filter W0 (3x3x3) | | | | | | | | Filter W1 (3x3x3) | | | | | | | | Output Volume (3x3x2) |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

**Input Volume (+pad 1) (7x7x3)**   **Filter W0 (3x3x3)**   **Filter W1 (3x3x3)**   **Output Volume (3x3x2)**

x[:,:,0]
```
0  0  0  0  0  0  0
0  1  1  0  1  2  0
0  0  1  0  0  0  0
0  0  1  0  2  1  0
0  0  1  2  0  2  0
0  2  0  2  0  1  0
0  0  0  0  0  0  0
```

w0[:,:,0]
```
 1   1  -1
 0   0  -1
 1   0   0
```
**-1**

w1[:,:,0]
```
 0   1  -1
 1  -1   0
 1   0   1
```

o[:,:,0]
```
-3   2   0
 1   3   4
-2   8   1
```

w0[:,:,1]
```
 1   0   0
 1  -1  -1
 1   1   0
```
**-2**

w1[:,:,1]
```
 0  -1   1
 0   1   1
 0   1   1
```

o[:,:,1]
```
 8   7   3
 5  11   5
 0   3   4
```

x[:,:,1]
```
0  0  0  0  0  0  0
0  1  2  1  0  2  0
0  1  2  1  1  0  0
0  0  0  0  1  0  0
0  2  1  2  1  0  0
0  1  2  0  0  2  0
0  0  0  0  0  0  0
```

w0[:,:,2]
```
 0   0   1
-1   1  -1
-1  -1   0
```
**-1**

w1[:,:,2]
```
 1   1   1
 0   1   0
 1   1   0
```

**Bias b0 (1x1x1)**   **Bias b1 (1x1x1)**

b0[:,:,0]
```
1
```
**1**

b1[:,:,0]
```
0
```

x[:,:,2]
```
0  0  0  0  0  0  0
0  1  1  0  0  1  0
0  1  0  2  0  1  0
0  0  0  1  1  2  0
0  0  0  2  1  0  0
0  1  0  1  1  0  0
0  0  0  0  0  0  0
```

toggle movement

**Hyperparameters:**
- Filter size
- Number of filter sets (channels)
- Padding
- Stride

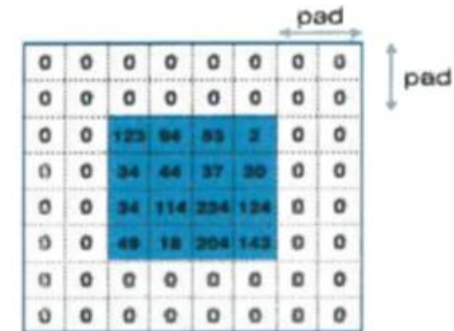N.B. this example the **stride** = 2, **padding** = 1, **channels = 2**
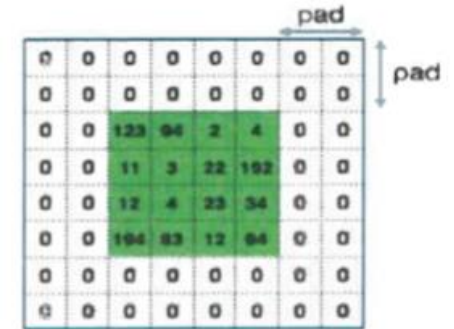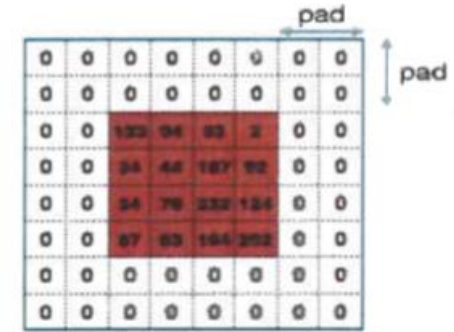
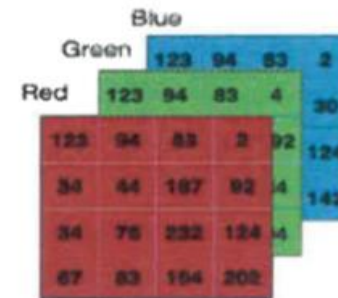https://cs231n.github.io/convolutional-networks/

# Padding

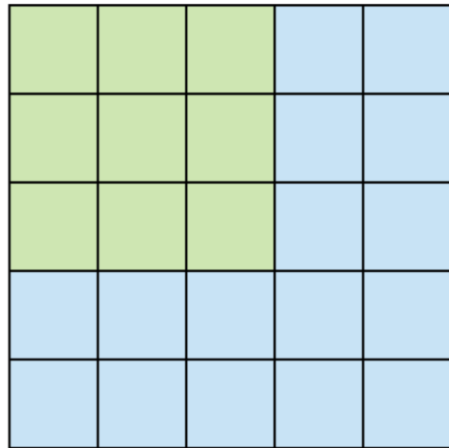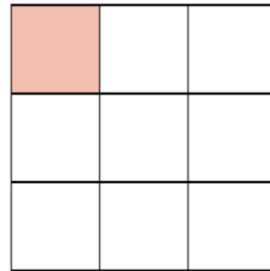*Adds zeros around the border of an image*



**Use cases :**

- Keeps more information at the border of an image

- Allows to *use a CONV layer without shrinking* the height and width of the volumes (important for deeper networks)
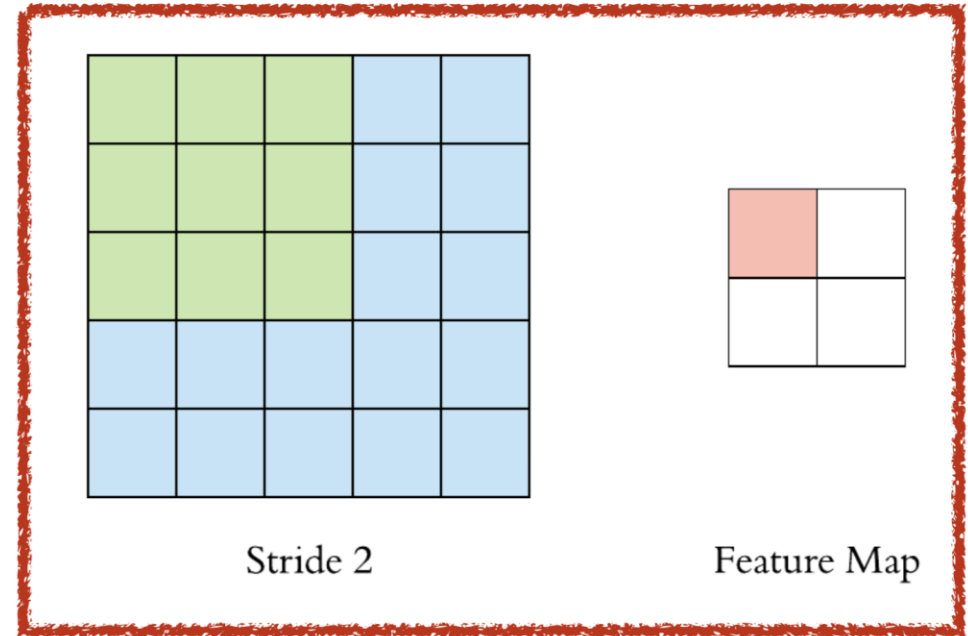
# Strided convolutions
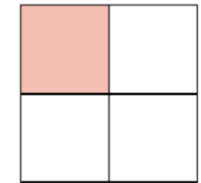
*By how much you move the filter*
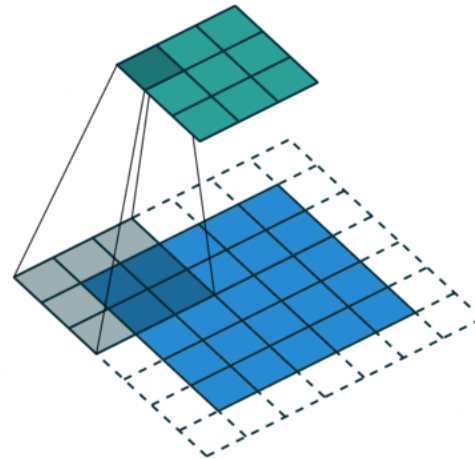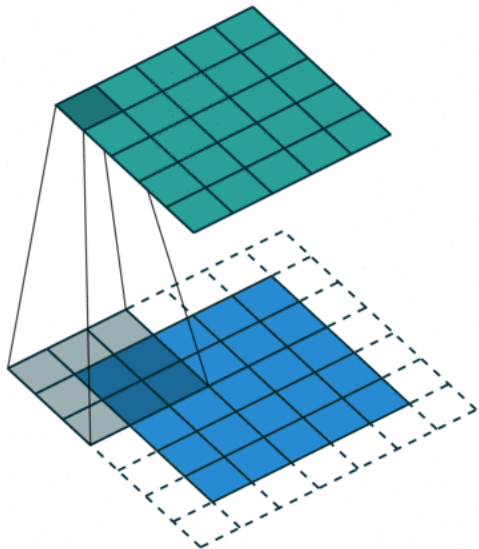


Stride 1      Feature Map      Stride 2      Feature Map
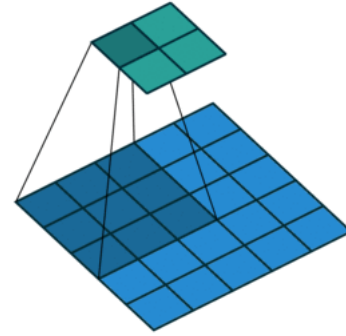
**increasing stride from 1 to 2**

# Illustration



https://hannibunny.github.io/mlbook/neuralnetworks/convolutionDemos.html

# Pooling Layer (POOL)

- reduces the spatial dimension ($n_H$ and $n_W$) to decrease computational power



Max Pool

Max-Pool with a 2 by 2 filter and stride 2.

*Get the max value*

Average Pool

Average Pool with a 2 by 2 filter and stride 2.

*Get the average value*

# Object Detection



Image **classification**

Classification with localization

Detection

"Car"

"Car"

1 object

multiple objects

12

# Yolo (*YOu Look only Once*)

- define a grid in the image

- apply the training to each cell (need ground-truth bounding boxes)

- For each 'anchor box' and 3 classes we have:



- Allows for overlapping objects

# YOLO prediction visualisation



- Filter the boxes using :
  1) score thresholding
  2) non-max suppression

Legend:
- car
- road sign
- tree
- traffic light
- sky
- background

# Score Thresholding

- Throw away boxes that have detected a class with a score less than the threshold (*0.6 for example*)

# Non-max suppression

- ensures that an object is detected only once
  - Choice based on the $p_C$ value : *keep the largest $p_C$ output* and *discard any remaining box with IoU>0.5*



Before non-max suppression

After non-max suppression

Non-Max Suppression

16

# Intersection Over Union (IOU)

- performance metric on how similar two boxes are with each other
- (higher the better!)

**Outcome of the algorithm**



**True bounding box**

$$IoU = \frac{\text{Size of} \ \boxed{//}}{\text{Size of} \ \boxed{//}}$$

# Intersection Over Union (IOU)



Yellow = intersection I
Green = union U

IoU  = I/U

Can express this as:

TP / (TP + FN + FP)

[see also, Jaccard Index]

# Generative Models

Holy grail of Deep Learning these days



DALL-E : 'A photograph of a cow on the moon.'

# Generative Models Examples

**Music Generator**

Random seed = 5 → Sequence Model →



DeepBach

**Text Captioning**

Random seed = 10 → Sequence Model →

*Neuron is a cell that can transmit electrical signals. They can also refer to units in ANN model................*

**Image Generator**

Random seed = 15 → Sequence Model →



GauGAN

# RNN model



$$a^{<t>} = g_1(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$y^{<t>} = g_2(W_{ya}a^{<t>} + b_y)$$

$W_{ax}, W_{aa}, W_{ya}, b_a$ and $b_y$ are weights that are shared temporally and $g_1, g_2$ activation functions

# Applications of RNN

| Type of RNN | Illustration | Example |
|:---:|:---:|:---:|
| One-to-one $T_x = T_y = 1$ | $a^{<0>} \rightarrow \Box \rightarrow \hat{y}$ with input $x$ | Traditional neural network |

# Applications of RNN

| Type of RNN | Illustration | Example |
|---|---|---|
| One-to-many $T_x = 1, T_y > 1$ |  | Music generation |

# Applications of RNN

| Type of RNN | Illustration | Example |
|---|---|---|
| Many-to-one $T_x > 1, T_y = 1$ |  | Sentiment classification |

# Applications of RNN

| Type of RNN | Illustration | Example |
|---|---|---|
| Many-to-many $T_x \neq T_y$ |  | Machine translation |

# VDSM



$n$ — Individual

$t$ — Video Frame

$\mathbf{d}^n$ — Action Performed

$\mathbf{s}^n$ — Identity

$\mathbf{z}_t^n$ — Per-Frame Pose

$\mathbf{h}_t^n$ — Per-image embedding

# Generative Adversarial Network

Invented by Ian Goodfellow

*How much are you ready to pay for it ?*

Christie's New York 2018

# Principle



**GENERATOR**
"The Artist"
A neural network trying to create pictures of cats that look real.

**GENERATOR**

Thousands of real-world images labeled "CAT"

**DISCRIMINATOR**
"The Art Critic"
A neural network examining cat pictures to determine if they're real or fake.

**DISCRIMINATOR**

# Principle

# Principle

Step 1: Train Discriminator and 'Freeze' Generator parameters
Step 2: Train Generator and 'Freeze' Discriminator parameters

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

https://youtu.be/-Zi5SReze6U

# Principle

Step 1: Train Discriminator and 'Freeze' Generator parameters
Step 2: Train Generator and 'Freeze' Discriminator parameters

$$\min_{G} \max_{D} V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

Discriminator gradient for update (gradient ascent):

predict well on real images
=> want probability close to 1

predict well on fake images
=> want probability close to 0

$$\nabla_{\mathbf{w}_D} \frac{1}{n} \sum_{i=1}^{n} \left[ \log D\left(\boldsymbol{x}^{(i)}\right) + \log\left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right) \right]$$

https://youtu.be/-Zi5SReze6U

# Principle

Step 1: Train Discriminator and 'Freeze' Generator parameters
Step 2: Train Generator and 'Freeze' Discriminator parameters

$$\min_G \max_D V(D, G) = \mathbb{E}_{\boldsymbol{x} \sim p_{\text{data}}(\boldsymbol{x})}[\log D(\boldsymbol{x})] + \mathbb{E}_{\boldsymbol{z} \sim p_{\boldsymbol{z}}(\boldsymbol{z})}[\log(1 - D(G(\boldsymbol{z})))]$$

## Generator gradient for update (gradient descent):

predict badly on fake images
=> want probability close to 1

$$\nabla_{\mathbf{w}_G} \frac{1}{n} \sum_{i=1}^{n} \log \left(1 - D\left(G\left(\boldsymbol{z}^{(i)}\right)\right)\right)$$

https://youtu.be/-Zi5SReze6U

# Principle

# GAN use case : generate images



Figure 3. Example results by our proposed StackGAN, GAWWN [20], and GAN-INT-CLS [22] conditioned on text descriptions from CUB test set. GAWWN and GAN-INT-CLS generate 16 images for each text description, respectively. We select the best one for each of them to compare with our StackGAN.

# State Of The Art in GANs
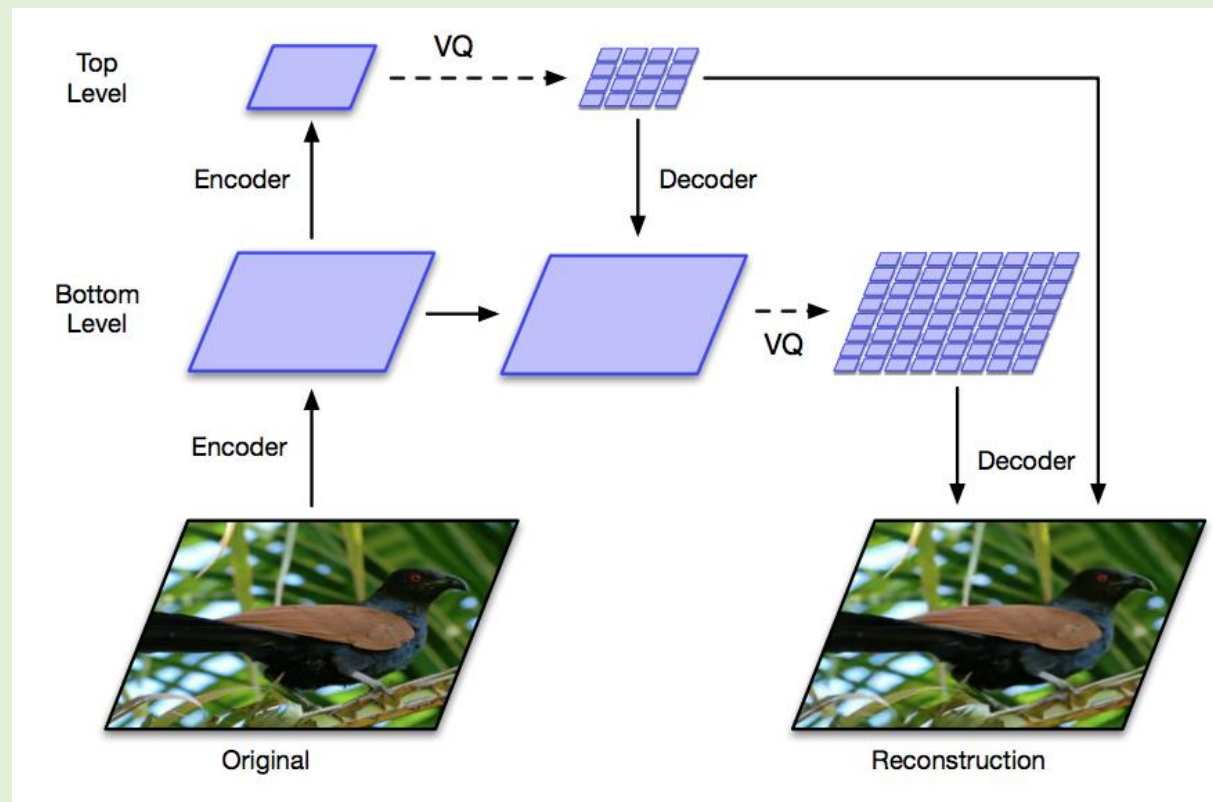


(Karras et al, 2018)

(Brock et al, 2018)

# Adversarial Principles are Widely Applicable
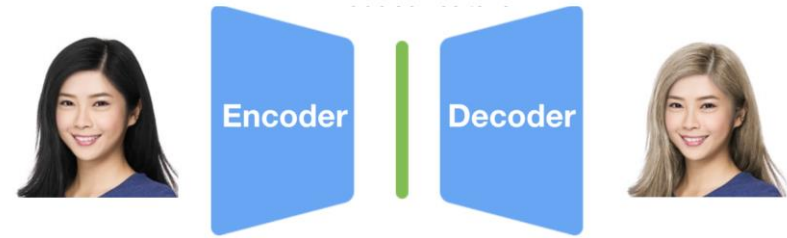
- Fairness / Privacy Preservation
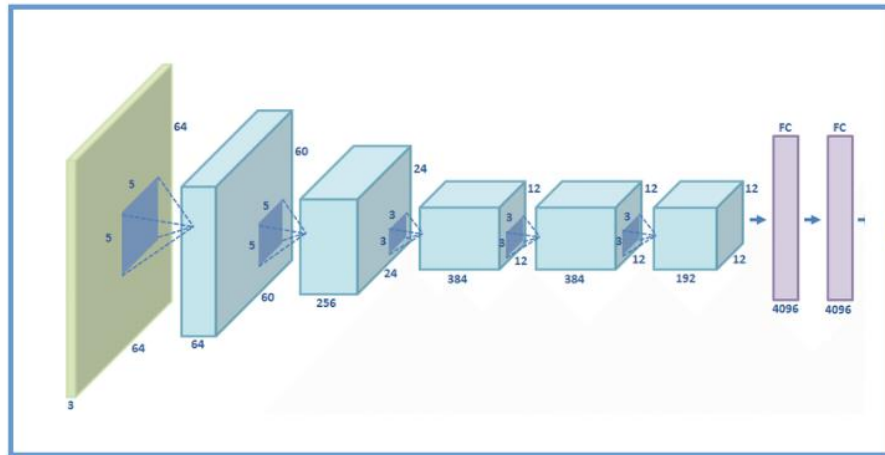- Disentanglement

# Variational AutoEncoder

# Principle

## AutoEncoder

• generate appropriately novel data



Input

Latent space/
Feature

Reconstruction

Encoder

Decoder

# Principle

- Pass images/data in through encoder *'bottleneck'*

- *Parameterize this bottleneck encoding so we can sample from it afterwards*

- Reconstruct the original image/data from the encoding

$p(x)$



Inference Encoder

$\phi$

$\boldsymbol{\mu}_\phi$  $\boldsymbol{\sigma}_\phi$

$q_\phi(z|x)$

+  *  $\epsilon \sim \mathcal{N}(0, I)$

$p(z)$

$z$

$$\mathrm{KL}\left[q_\phi(z|x)||p(z)\right]$$

$\theta$

Generation Decoder

Reconstruction Loss



$p_\theta(x|z)$
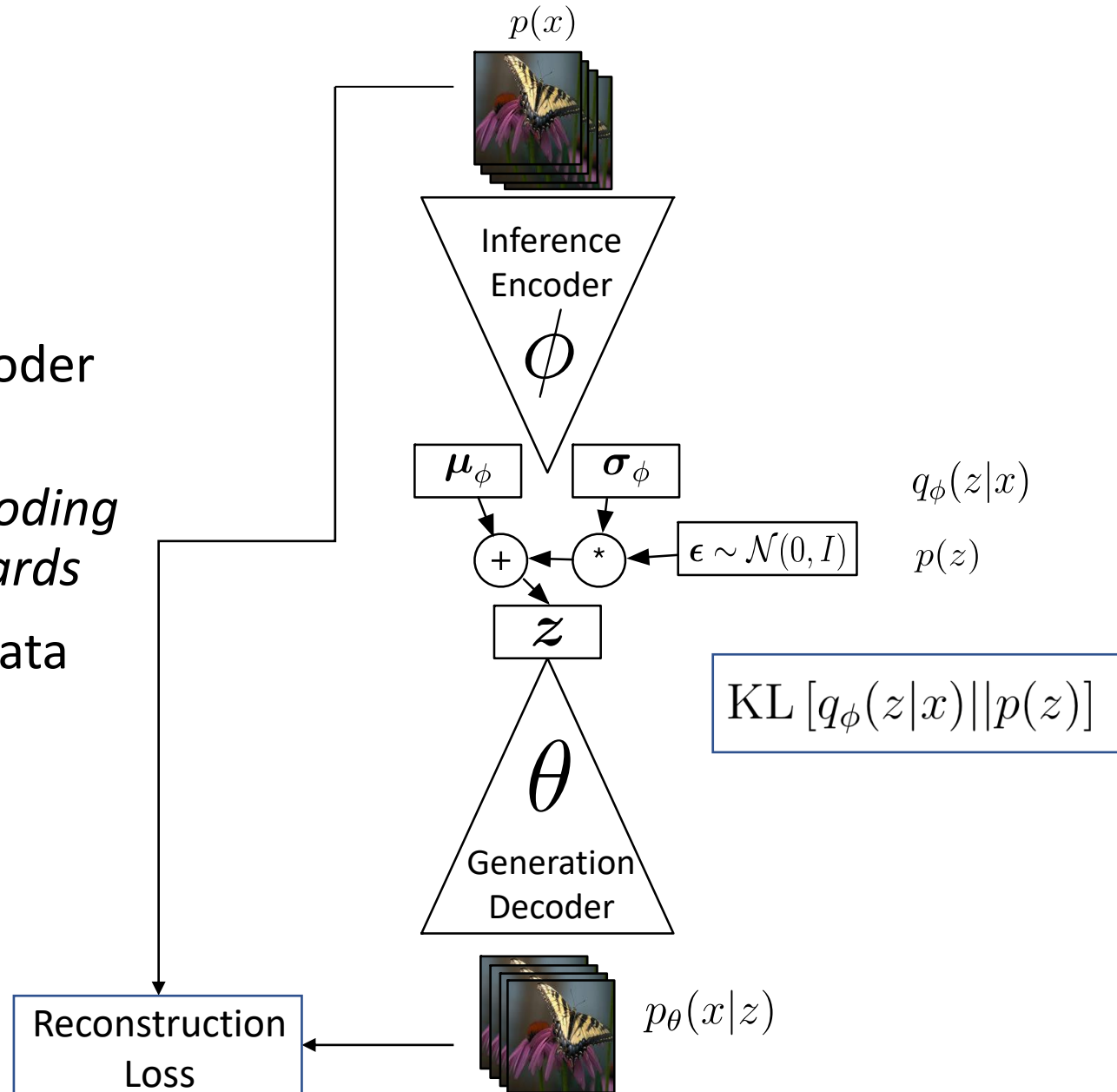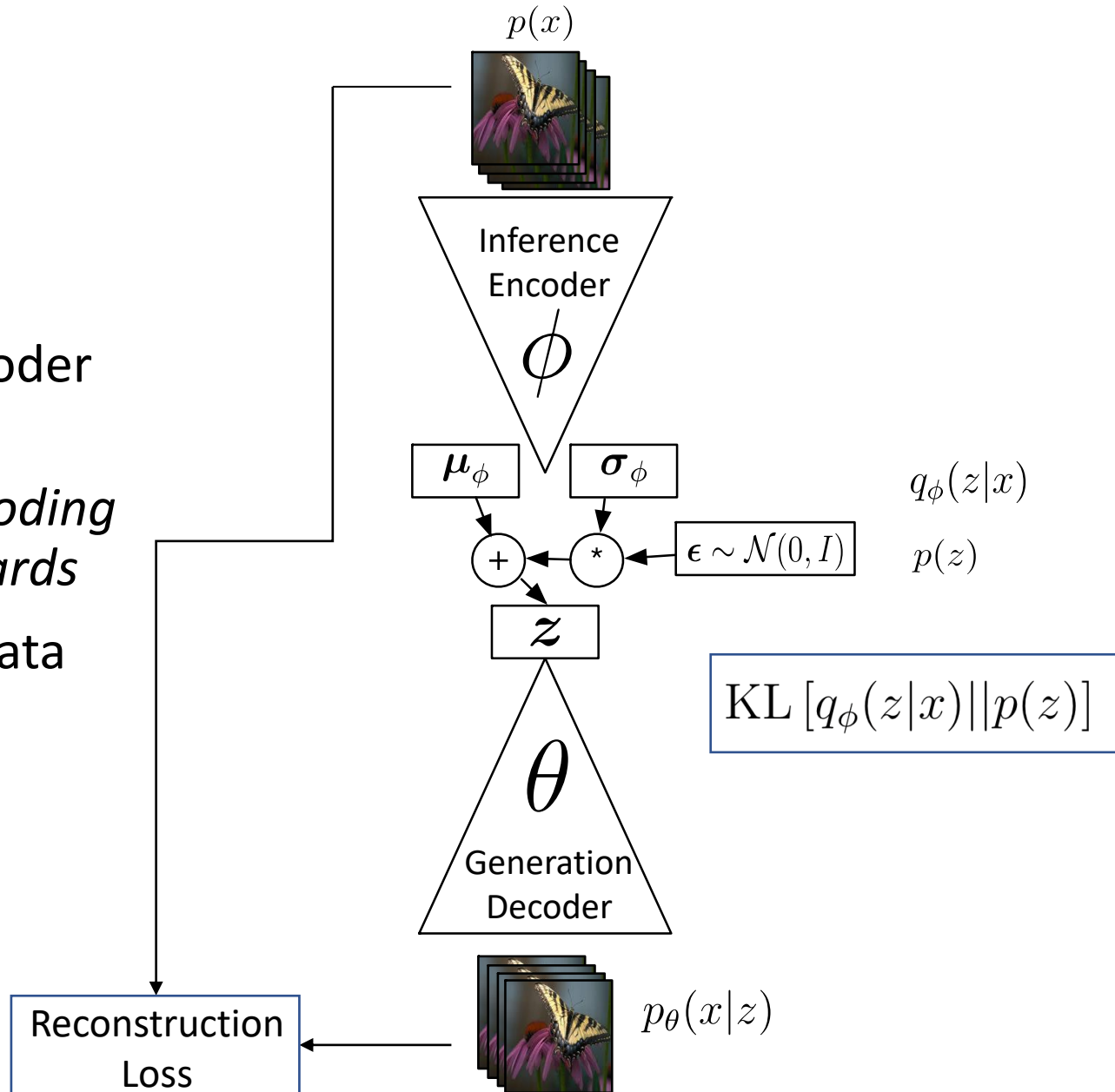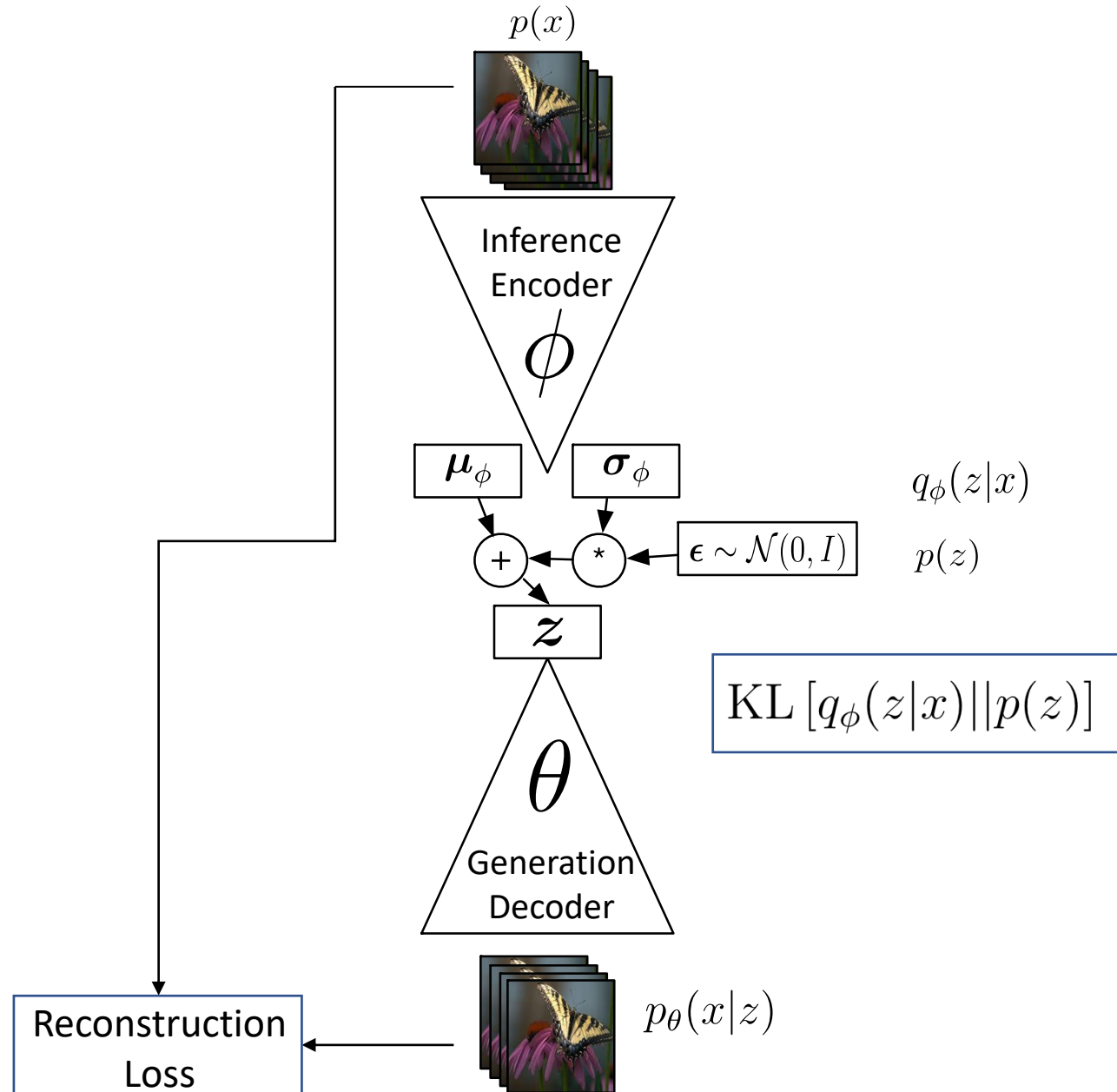
# Principle

- Pass images/data in through encoder *'bottleneck'*

- *Parameterize this bottleneck encoding so we can sample from it afterwards*

- Reconstruct the original image/data from the encoding

$p(x)$



Inference
Encoder

$\phi$

$\boldsymbol{\mu}_\phi$    $\boldsymbol{\sigma}_\phi$

$q_\phi(z|x)$

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

$+$    $*$

$\boldsymbol{z}$

$\mathrm{KL}\left[q_\phi(z|x)||p(z)\right]$

$\theta$

Generation
Decoder


$p_\theta(x|z)$

Reconstruction
Loss

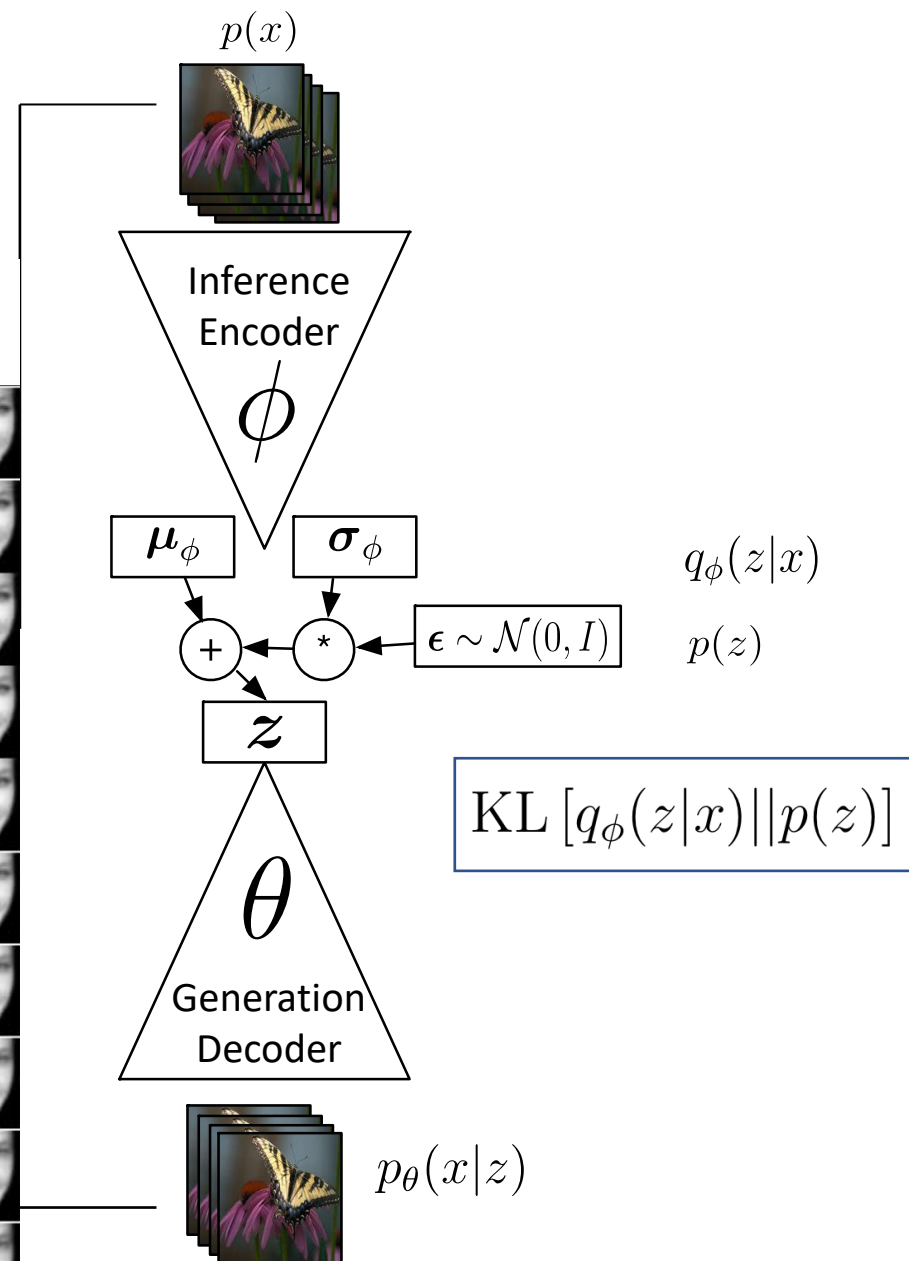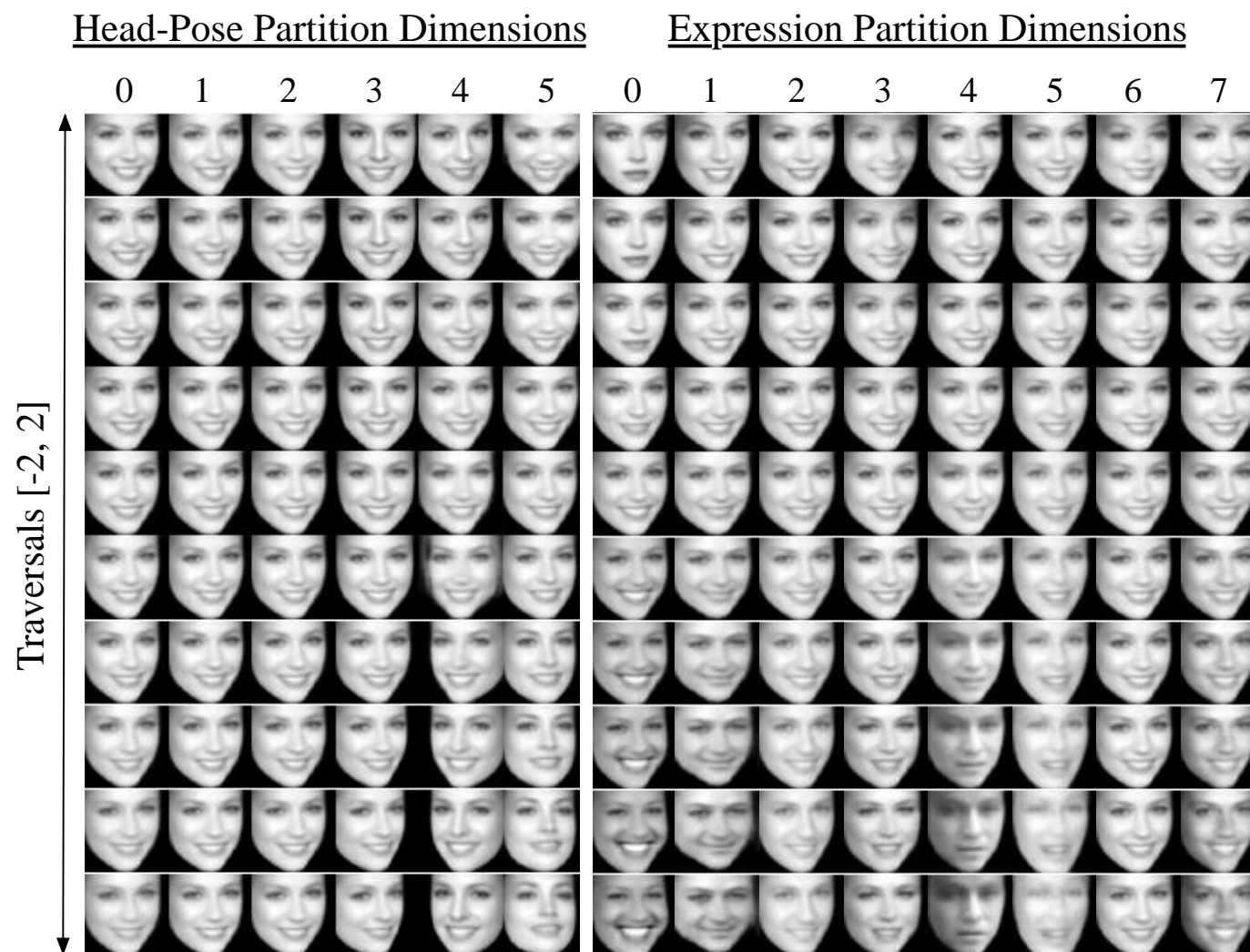# VAE use cases

- Compression

- Data generation

- Latent variable modeling

- Density estimation

# VAE use cases



Head-Pose Partition Dimensions

Expression Partition Dimensions

$p(x)$

Inference Encoder $\phi$

$q_\phi(z|x)$

$\boldsymbol{\mu}_\phi$ $\boldsymbol{\sigma}_\phi$

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

$\boldsymbol{z}$

$\mathrm{KL}\left[q_\phi(z|x)||p(z)\right]$

$\theta$

Generation Decoder

$p_\theta(x|z)$

Traversals [-2, 2]

# Tutorial / Practical