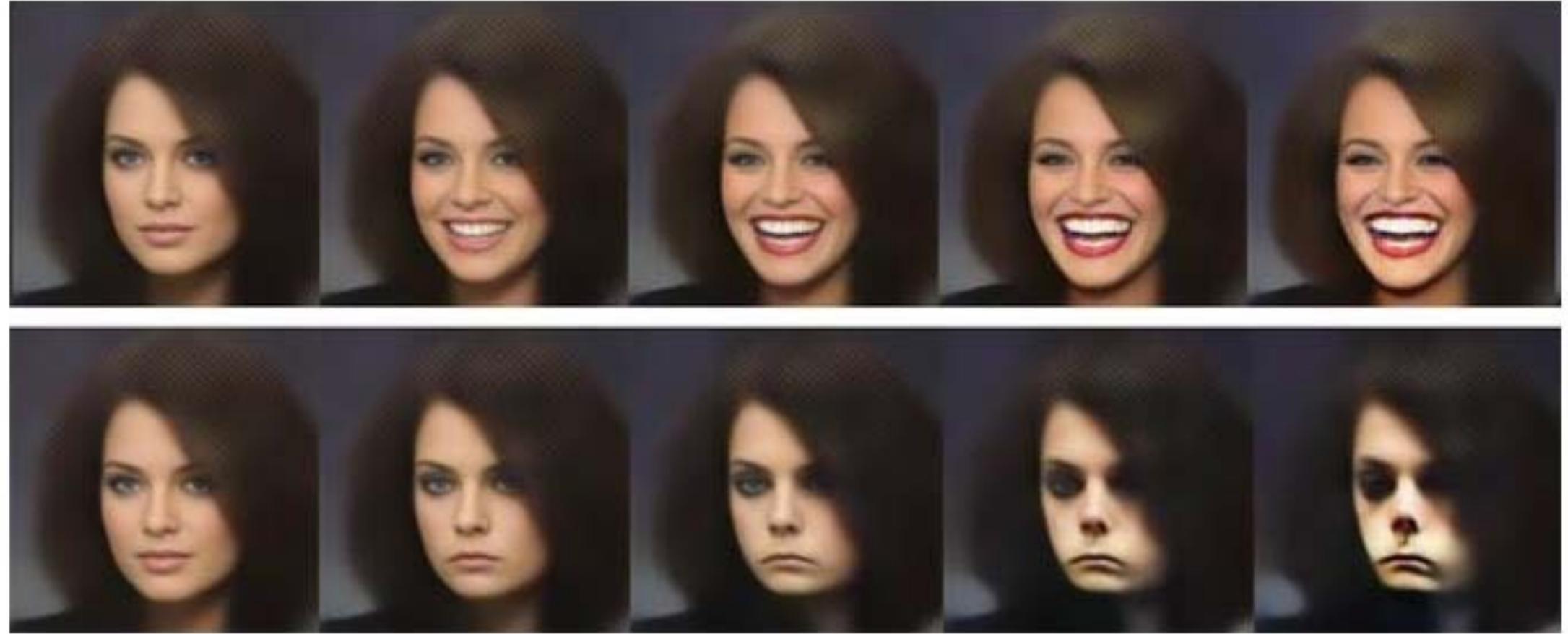


Generative Models



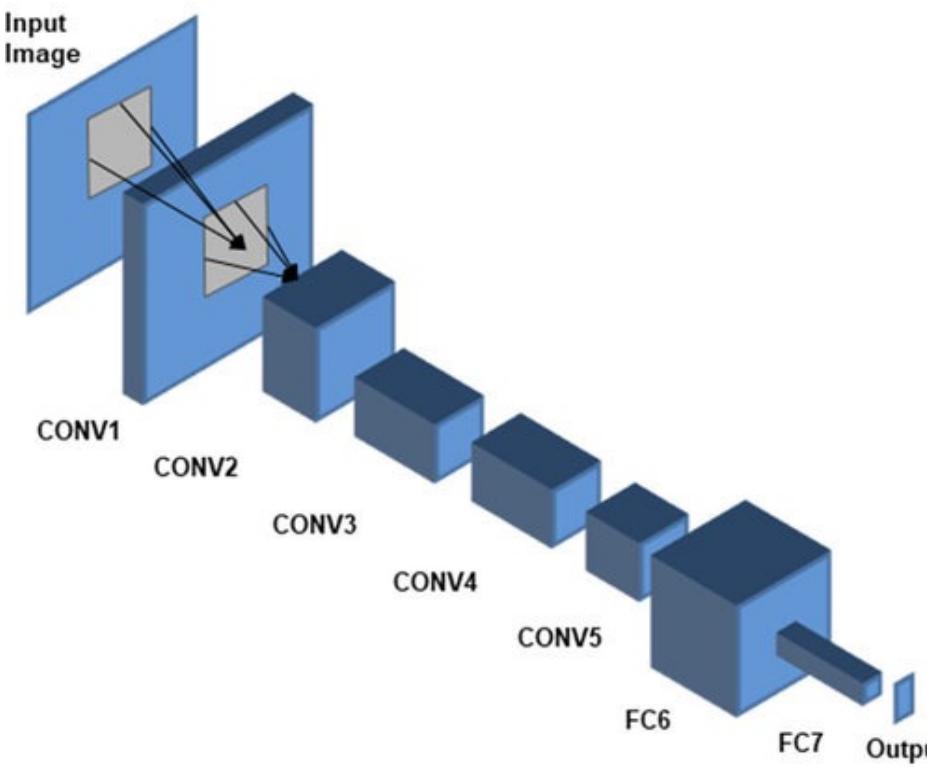
11. November 2023

Today : Autoencoders, Variational Autoencoders and Generative Adversarial Networks



Previously: Image Classification and Semantic Segmentation

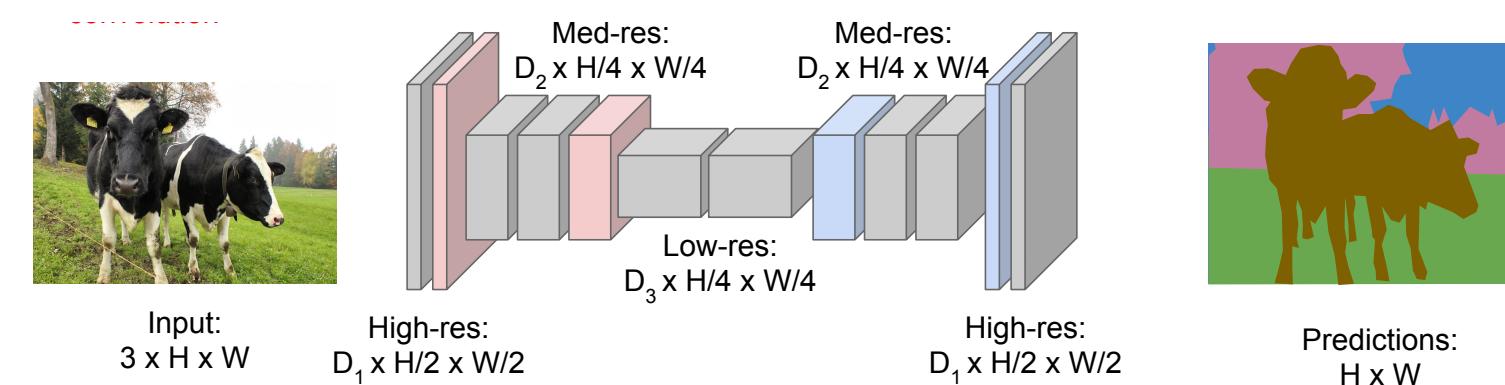
Image classification



Input: pixel values (high-dimensional)
Output: class probabilities (low-dimensional)

Classification = **compression** into low-dimensional summary of image

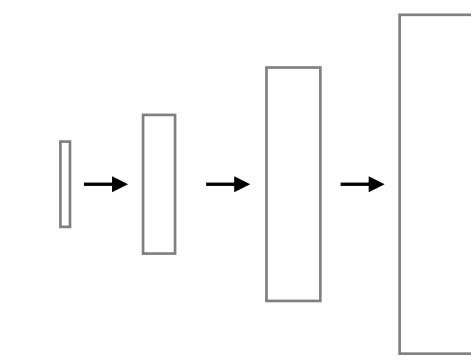
Semantic segmentation



Input: pixel values (high-dimensional)
Output: per-pixel class probabilities (high-dimensional)

Semantic segmentation = **translation** into high-dimensional, but interpretable space

Image generation



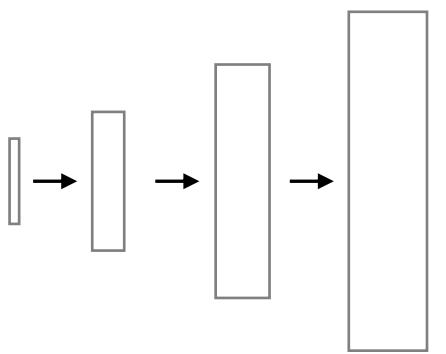
Input: low-dimensional image summary
Output: pixel values (high-dimensional)

Image generation = **expansion** of low-dimensional image summary into high-dimensional pixel values

Goal: Learn low-dimensional image representations and expansion into images from available images

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

Learned
image generation



But how learn the structure of the dataset,
so that new images can be generated?

You already know how to go from low-dimensional image representation to pixel values (cf. U-Net):

- Nearest-neighbor interpolation
- Linear interpolation
- Max un-pooling
- Transposed convolutions
- Etc.

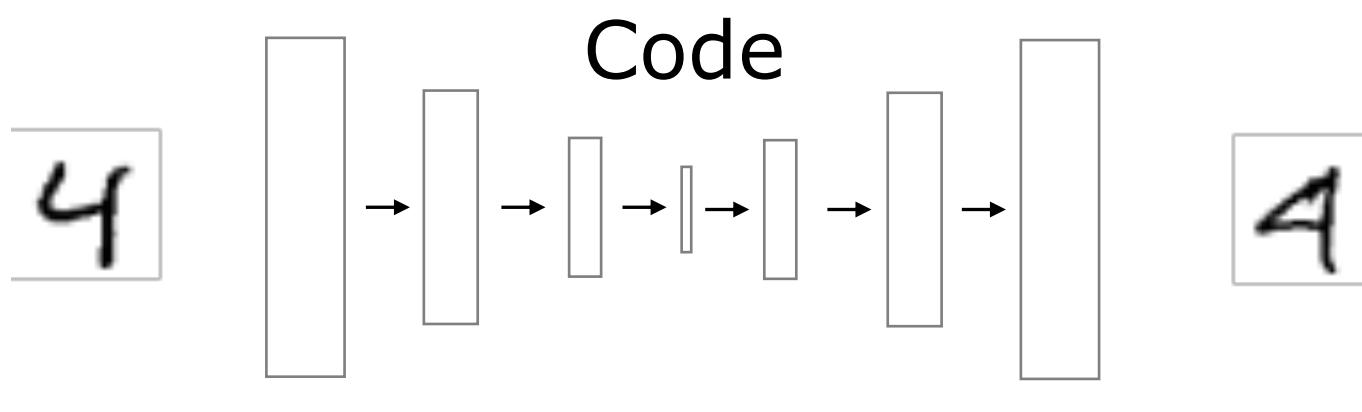
Autoencoder

Dataset

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

Input image

Output image



Autoencoder

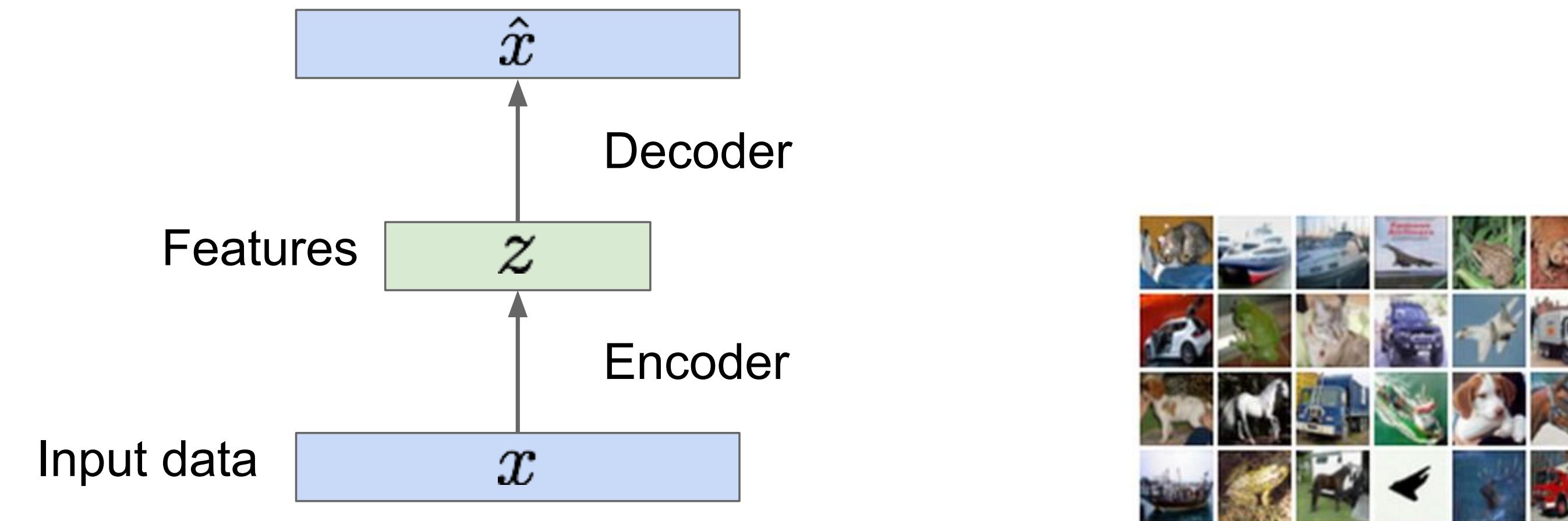
- Encoder E: a neural network
- Decoder D: a neural network
- Train such that $D(E(\text{image})) = \text{image}$ for each image in training set

Important: no labels required!

When done: throw away the encoder

Some Background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

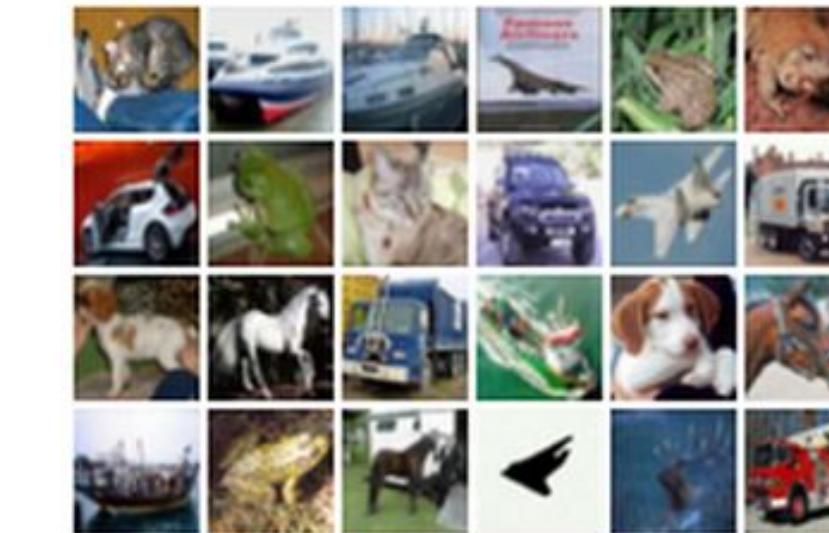
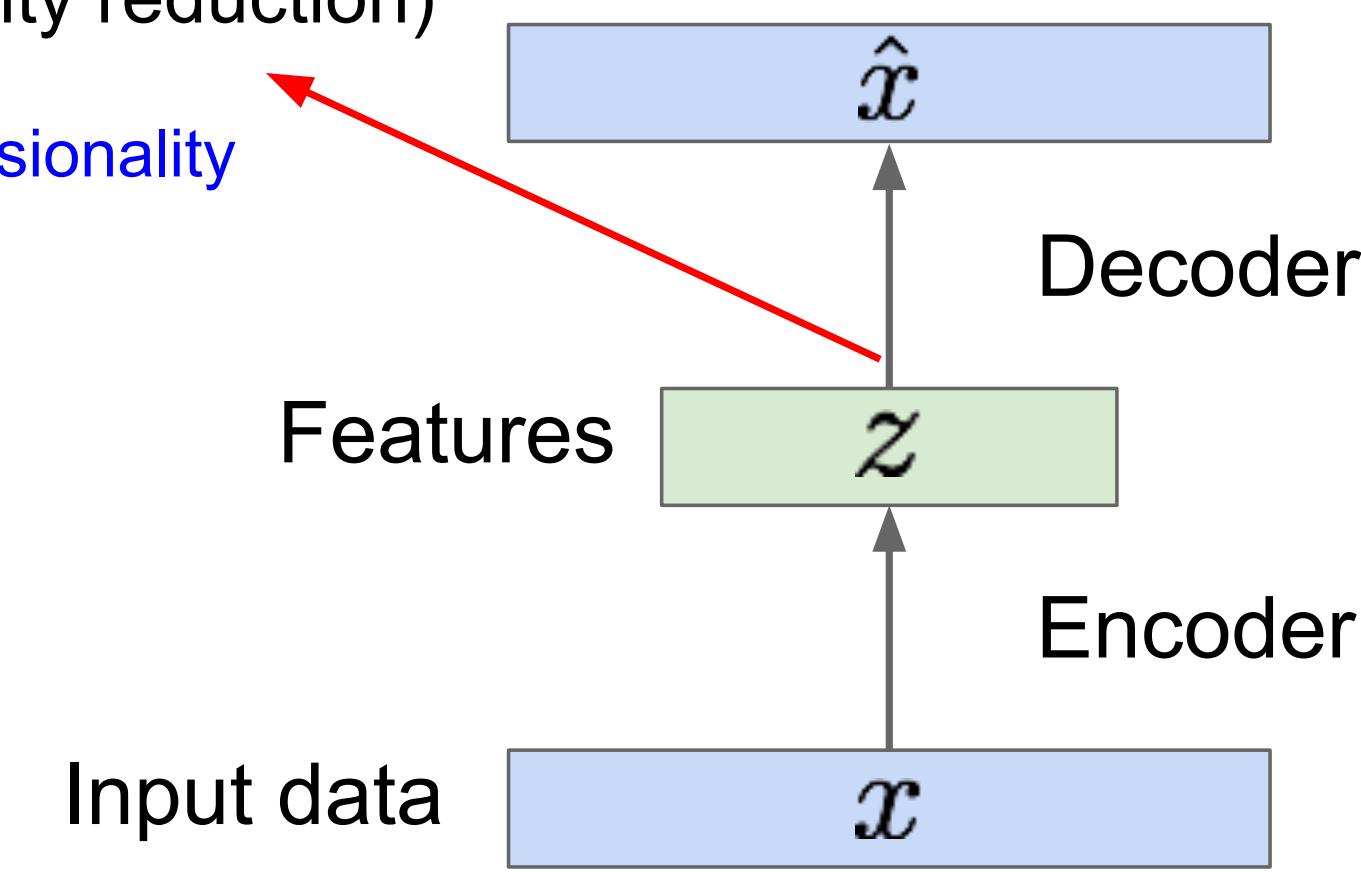


Some Background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?



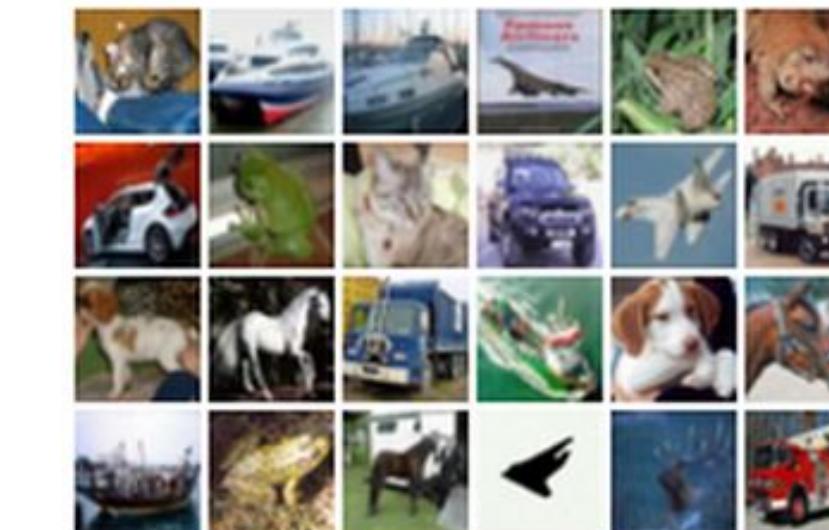
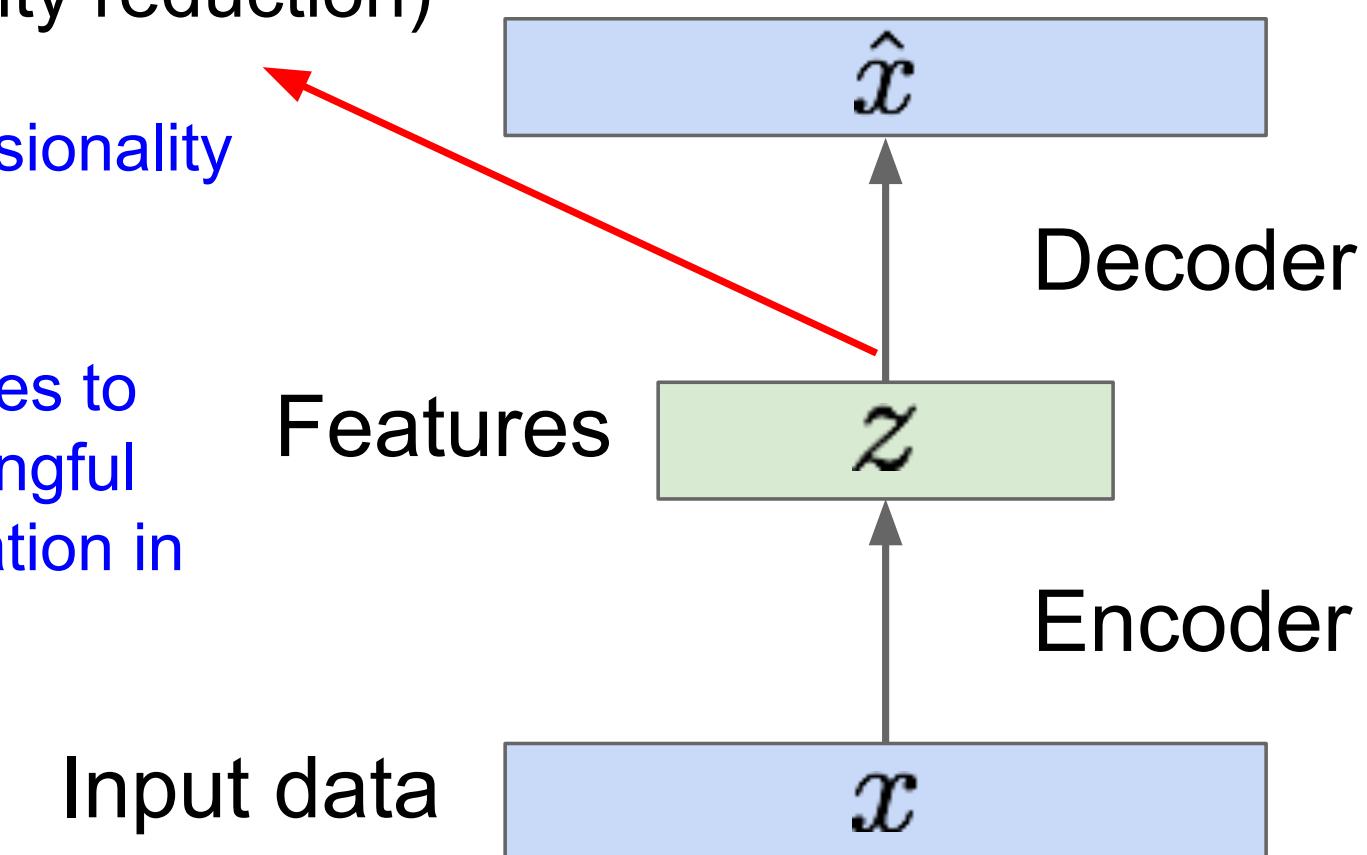
Some Background first: Autoencoders

Unsupervised approach for learning a lower-dimensional feature representation from unlabeled training data

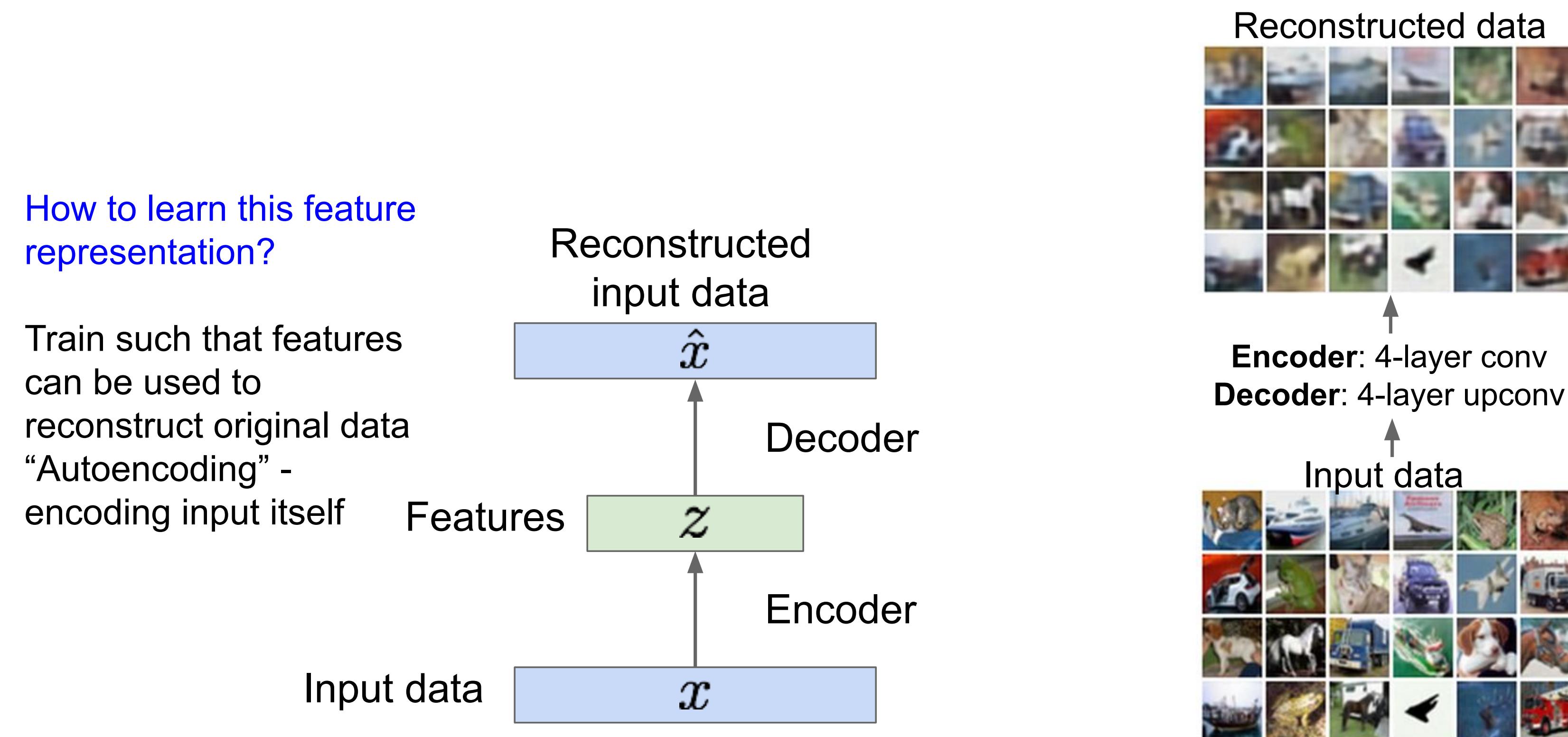
z usually smaller than x
(dimensionality reduction)

Q: Why dimensionality reduction?

A: Want features to capture meaningful factors of variation in data

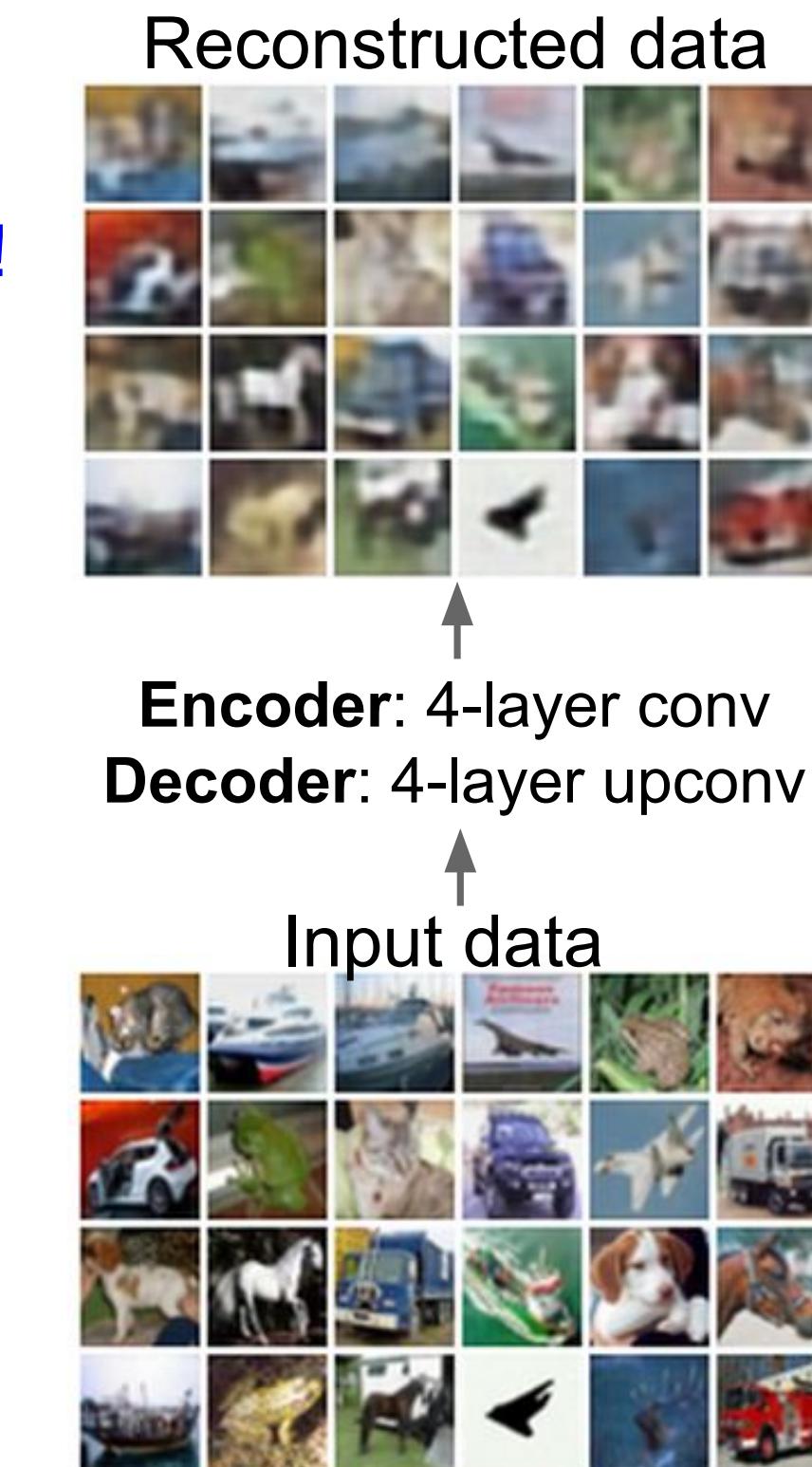
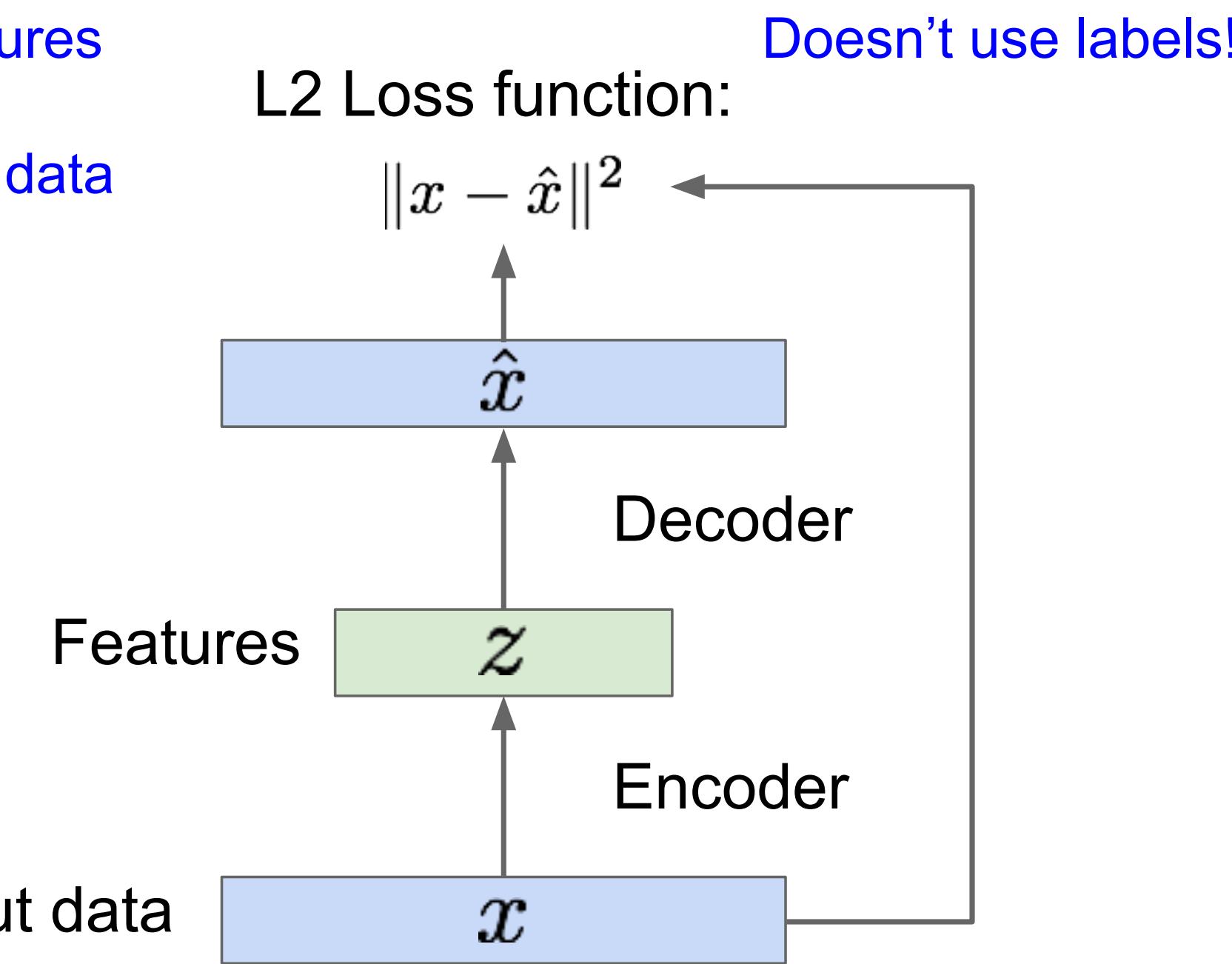


Some Background first: Autoencoders

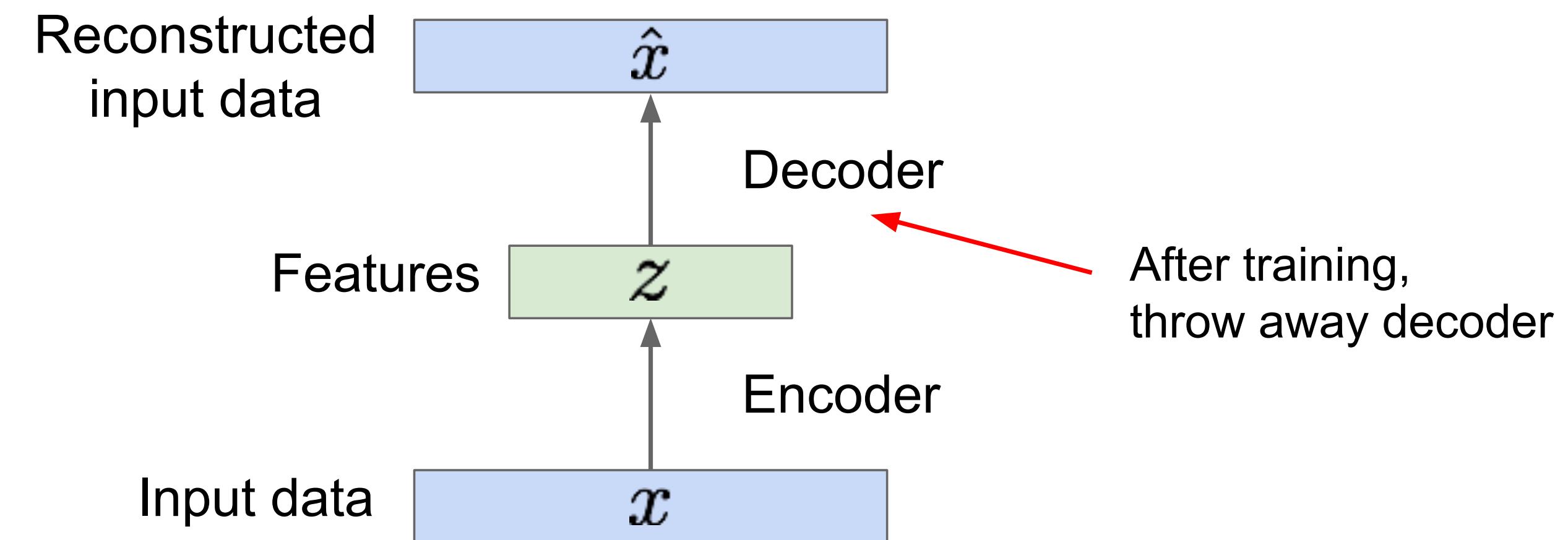


Some Background first: Autoencoders

Train such that features
can be used to
reconstruct original data



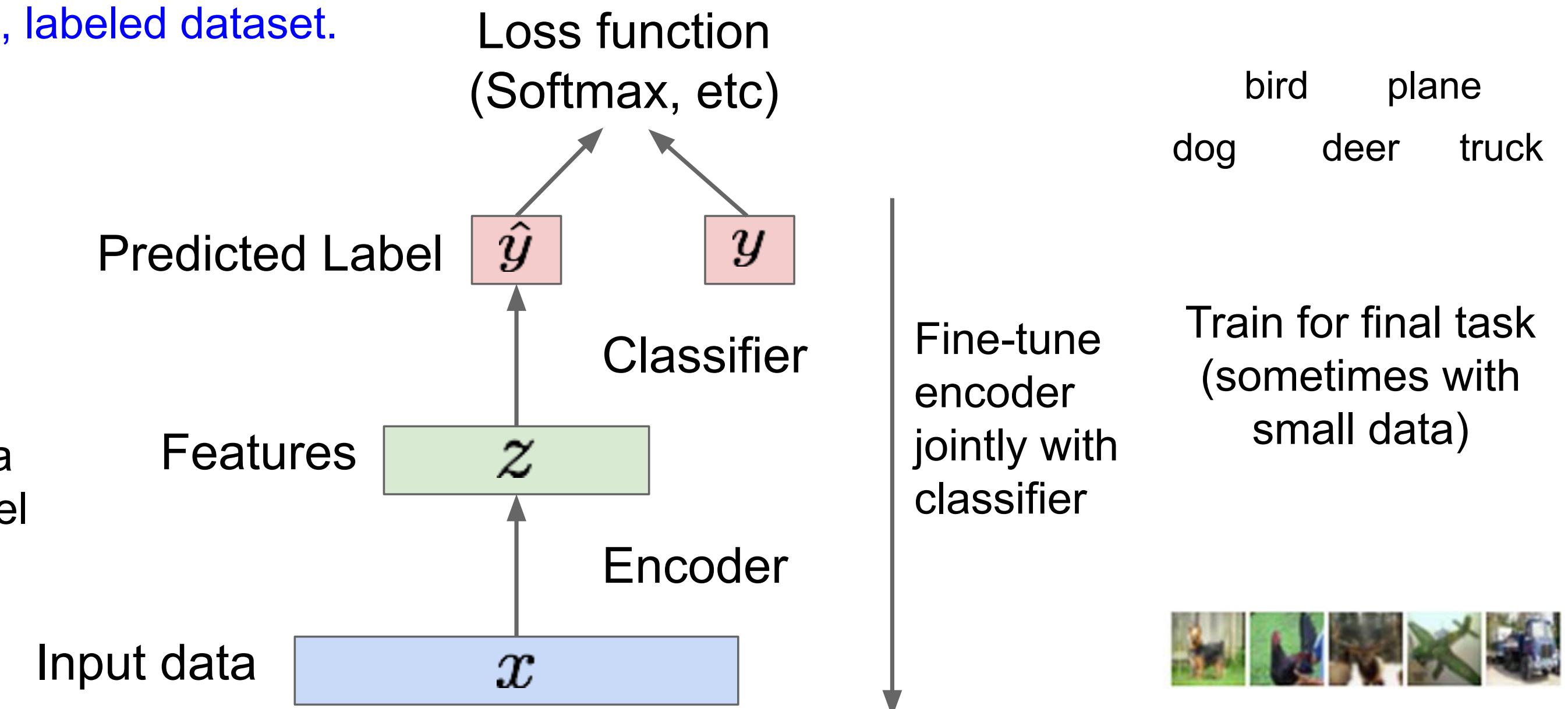
Some Background first: Autoencoders



Some Background first: Autoencoders

- Learn to encode / compress and decode / decompress from training images (without labels)
- Semi-supervised learning:
Code (learned without labels) can be used as input for small supervised machine learning model
- Autoencoder is trained to **reconstruct training images**
 - No way to generate new data from learned model

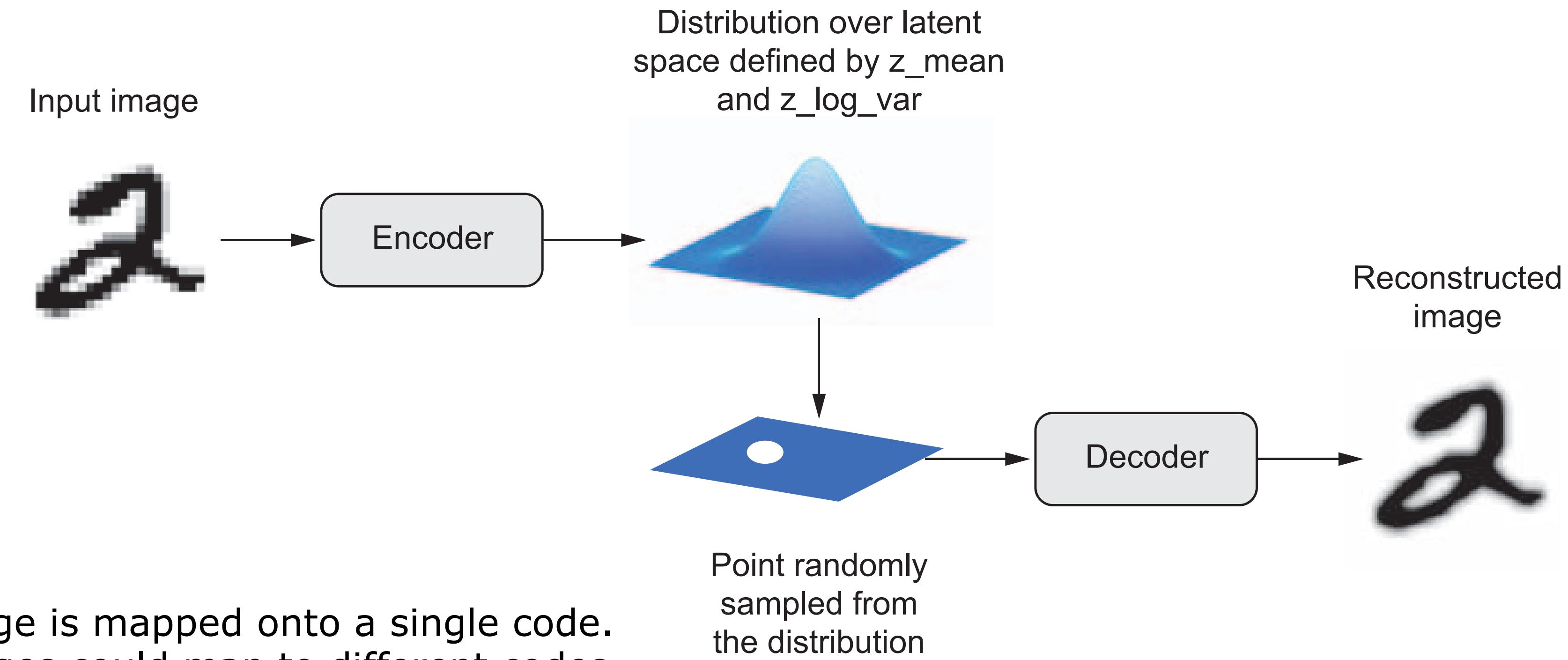
Transfer from large, unlabeled dataset to small, labeled dataset.



Encoder can be used to initialize a **supervised** model

Exercises

Variational Autoencoders



AE: Input image is mapped onto a single code.
 → Similar images could map to different codes
 → Similar codes could map to different images

VAE: Input image is mapped onto a smooth **distribution over codes**.
 → Thus, a input image is mapped onto a distribution over output images
 → Similar codes should decode to similar images
 → Imposes structure on latent space (e.g., independent features)

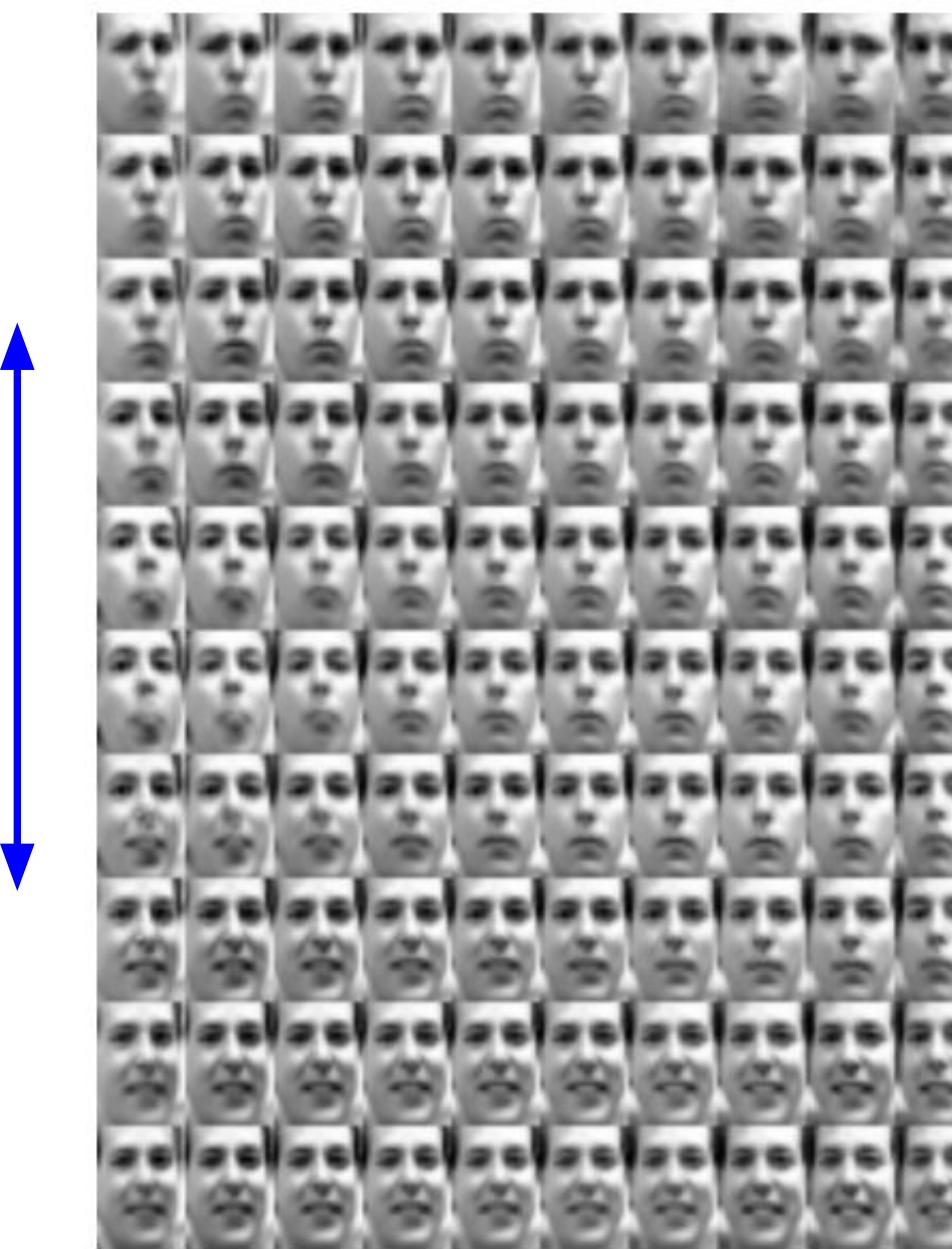
Variational Autoencoder : Generating Data

Diagonal prior on \mathbf{z}
=> independent
latent variables

Different
dimensions of \mathbf{z}
encode
interpretable factors
of variation

Degree of smile

Vary \mathbf{z}_1



Vary \mathbf{z}_2

Head pose

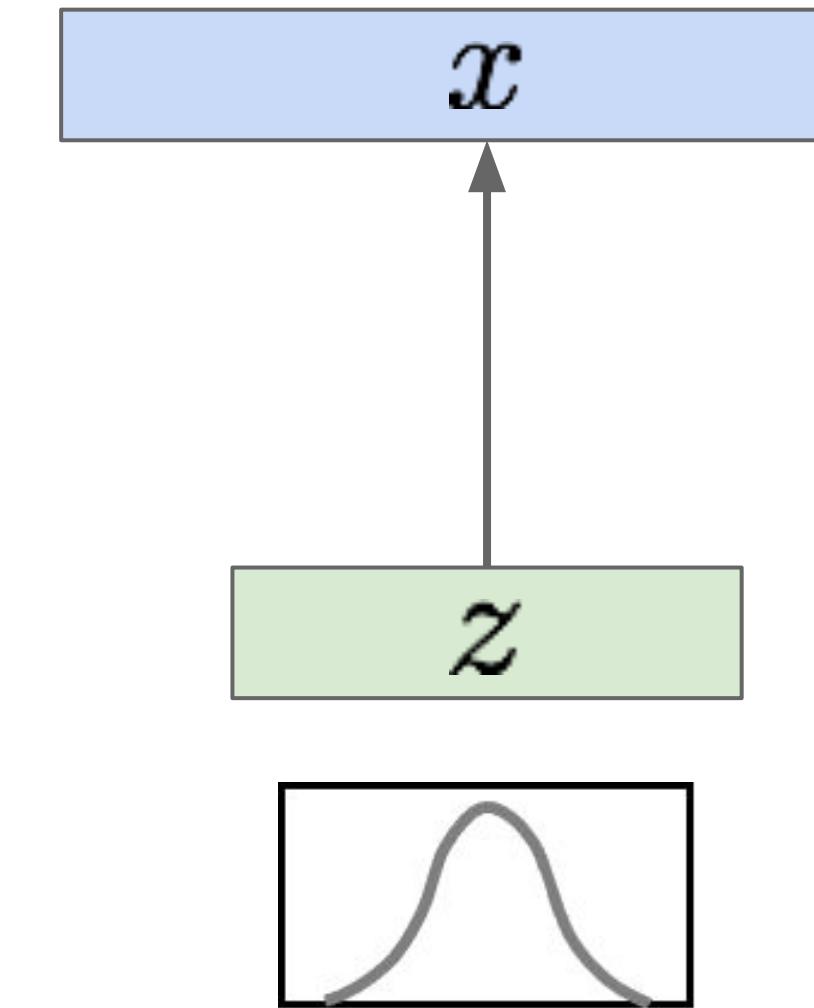
Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoder : Generating Data

We want to estimate the true parameters θ^* of this generative model given training data x .

Sample from
true conditional
 $p_{\theta^*}(x \mid z^{(i)})$

Sample from
true prior
 $z^{(i)} \sim p_{\theta^*}(z)$



How should we represent this model?

Choose prior $p(z)$ to be simple, e.g.
Gaussian. Reasonable for latent attributes,
e.g. pose, how much smile.

Kingma and Welling, "Auto-Encoding Variational Bayes", ICLR 2014

Variational Autoencoder

Decoder:
reconstruct
the input data

$$\begin{aligned}
 \log p_\theta(x^{(i)}) &= \mathbf{E}_{z \sim q_\phi(z|x^{(i)})} [\log p_\theta(x^{(i)})] \quad (p_\theta(x^{(i)}) \text{ Does not depend on } z) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \right] \quad (\text{Bayes' Rule}) \\
 &= \mathbf{E}_z \left[\log \frac{p_\theta(x^{(i)} | z)p_\theta(z)}{p_\theta(z | x^{(i)})} \frac{q_\phi(z | x^{(i)})}{q_\phi(z | x^{(i)})} \right] \quad (\text{Multiply by constant}) \\
 &= \mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z)} \right] + \mathbf{E}_z \left[\log \frac{q_\phi(z | x^{(i)})}{p_\theta(z | x^{(i)})} \right] \quad (\text{Logarithms}) \\
 &= \underbrace{\mathbf{E}_z [\log p_\theta(x^{(i)} | z)] - D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z))}_{\mathcal{L}(x^{(i)}, \theta, \phi)} + \underbrace{D_{KL}(q_\phi(z | x^{(i)}) || p_\theta(z | x^{(i)}))}_{\geq 0}
 \end{aligned}$$

Encoder:
make approximate
posterior distribution
close to prior

Tractable lower bound which we can take gradient of and optimize! ($p_\theta(x|z)$ differentiable, KL term differentiable)

Evidence lower bound (ELBO)

Variational:

- Encoder $q(z | x)$ is an approximation of the true distribution $p(z | x)$
- The approximation error is computed by the **Kullback-Leibler divergence** (D_{KL})

Variational Autoencoder : Generating Data

Probabilistic spin to traditional autoencoders => allows generating data

Defines an intractable density => derive and optimize a (variational) lower bound

Pros:

- Principled approach to generative models
- Interpretable latent space
- Allows inference of $q(z|x)$, can be useful feature representation
for other tasks

Cons:

Samples blurrier and lower quality compared to state-of-the-art (GANs)

Exercises

Generative Adversarial Networks (Goodfellow et al. 2014)



Painter: she would like to learn to learn to paint as Van Gogh



Critic: she would like to learn be able to tell real Van Gogh paintings from fakes apart.

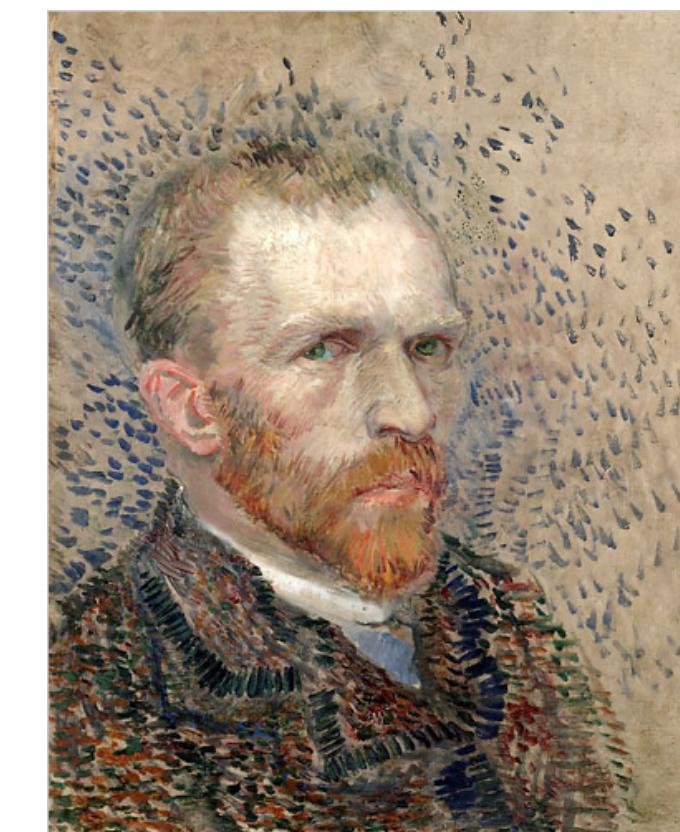
Step 1



Painter: she would like to learn to learn to paint as Van Gogh

The painter becomes slowly better

Produce a fake portrait that is not perfect



FAKE

Then the critic tells the painter what was clearly wrong.



Critic: she would like to learn be able to tell real Van Gogh paintings from fakes apart.

The critic try to determine if the painting is fake or not.

Step 2 – The critic must become better too



TRUE



FAKE

The critic trains with
real Van Gogh
paintings



The critic becomes
slowly better

The critic trains with
clearly fake Van
Gogh paintings

Reference – Loss Function (Cross Entropy, CE)

The diagram illustrates the Cross Entropy Loss function. A central equation is displayed: $y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$. Four blue arrows point from surrounding text to parts of this equation: one arrow points to the term $y \log(\hat{y})$ from the text 'Class Label: 0,1'; another arrow points to the term $(1 - y) \log(1 - \hat{y})$ from the text 'Predicted Class probability of painting of being real (class 1)'; a third arrow points to the first \hat{y} from the text 'Predicted Class probability of painting of being real (class 1)'; and a fourth arrow points to the second \hat{y} from the same text.

Class Label: 0,1
0 → Painting is fake
1 → Painting is real

$y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})$

Predicted Class probability of painting of being real (class 1)

Ideas of the painter
(today I will paint a portrait,
a landscape, etc.)

Step 2



X_{real}

Step B - Update

θ_D

Discriminator

$$L_D = \text{CE}(\hat{Y}_{\text{fake}}, 0) + \text{CE}(\hat{Y}_{\text{real}}, 1)$$



\hat{Y}_{real}

$$\text{CE}(\hat{Y}_{\text{real}}, 1)$$

Same weights

Discriminator

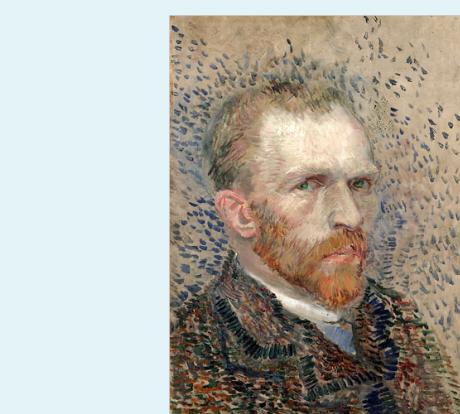
\hat{Y}_{fake}

$$\text{CE}(\hat{Y}_{\text{fake}}, 0)$$

Noise

ξ

Generator

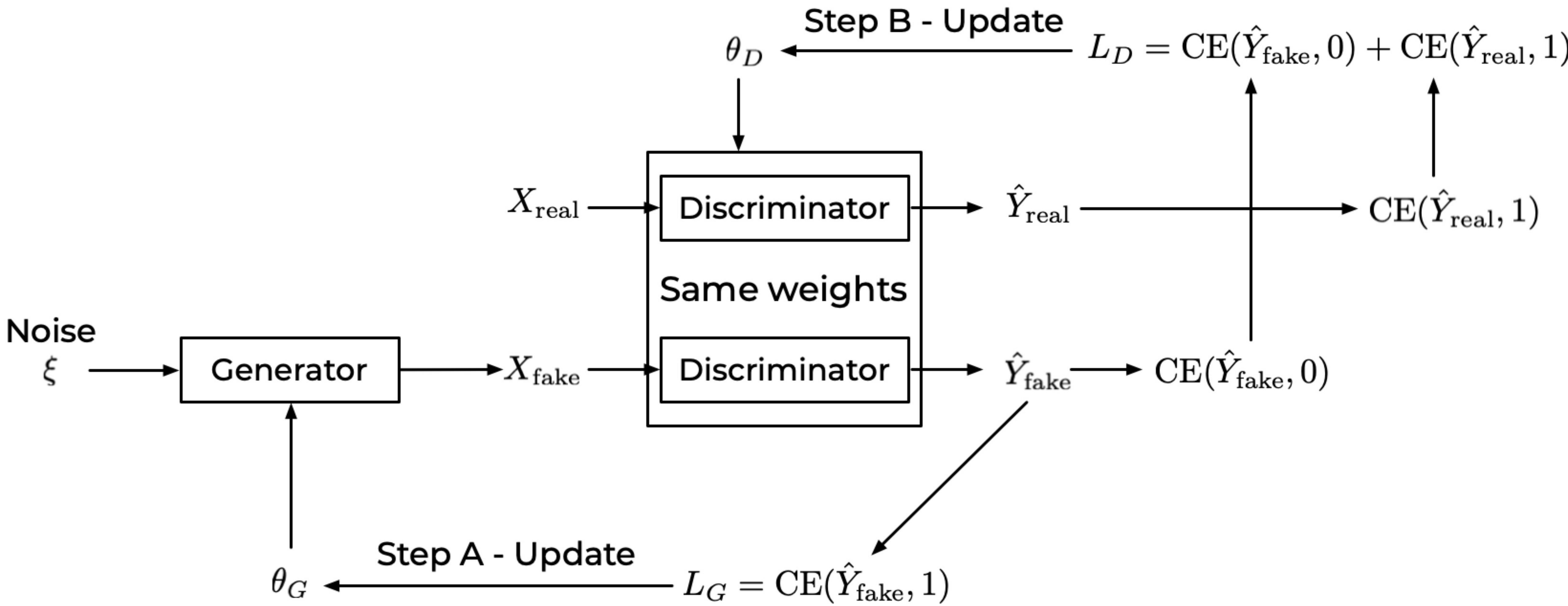


Step A - Update

θ_G

$$L_G = \text{CE}(\hat{Y}_{\text{fake}}, 1)$$

Step 1



How to use Keras for GANs

The most important aspect of using Keras for GANs is the necessity of building a custom training loop (and not using **compile()/fit()** approach.

```
def train_step(images):
    noise = tf.random.normal([BATCH_SIZE, noise_dim])

    with tf.GradientTape() as gen_tape, tf.GradientTape() as disc_tape:
        # Calculation of X_fake
        generated_images = generator(noise, training=True) ←  $X_{fake}$ 

        # Calculation of \hat{Y}_{real}
        real_output = discriminator(images, training=True) ←  $X_{real}$ 

        # Calculation of \hat{Y}_{fake}
        fake_output = discriminator(generated_images, training=True)

        # Calculation of L_G
        gen_loss = generator_loss(fake_output) ←  $L_G = CE(\hat{Y}_{fake}, 1)$ 

        # Calculation of L_D
        disc_loss = discriminator_loss(real_output, fake_output) ←  $L_D = CE(\hat{Y}_{fake}, 0) + CE(\hat{Y}_{real}, 1)$ 

        #
        # Gradients Calculation
        #
        # Calculation of the gradients of L_G for backpropagation
        gradients_of_generator = gen_tape.gradient(gen_loss, generator.trainable_variables)

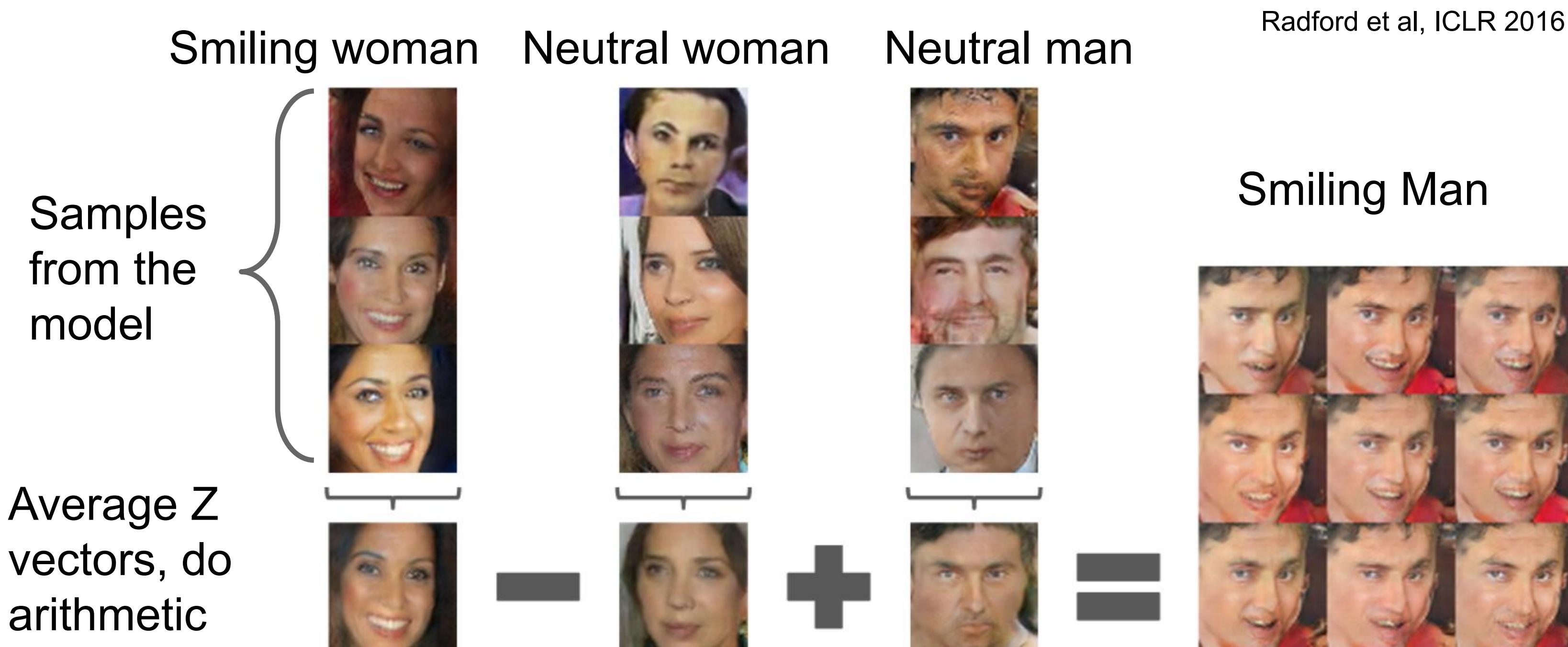
        # Calculation of the gradients of L_D for backpropagation
        gradients_of_discriminator = disc_tape.gradient(disc_loss, discriminator.trainable_variables)

        #
        # Training Steps A and B
        #
        # Step A
        generator_optimizer.apply_gradients(zip(gradients_of_generator, generator.trainable_variables))

        # Step B
        discriminator_optimizer.apply_gradients(zip(gradients_of_discriminator, discriminator.trainable_variables))
```

Exercises

Generative Adversarial Networks



Generative Adversarial Networks

Glasses man



No glasses man

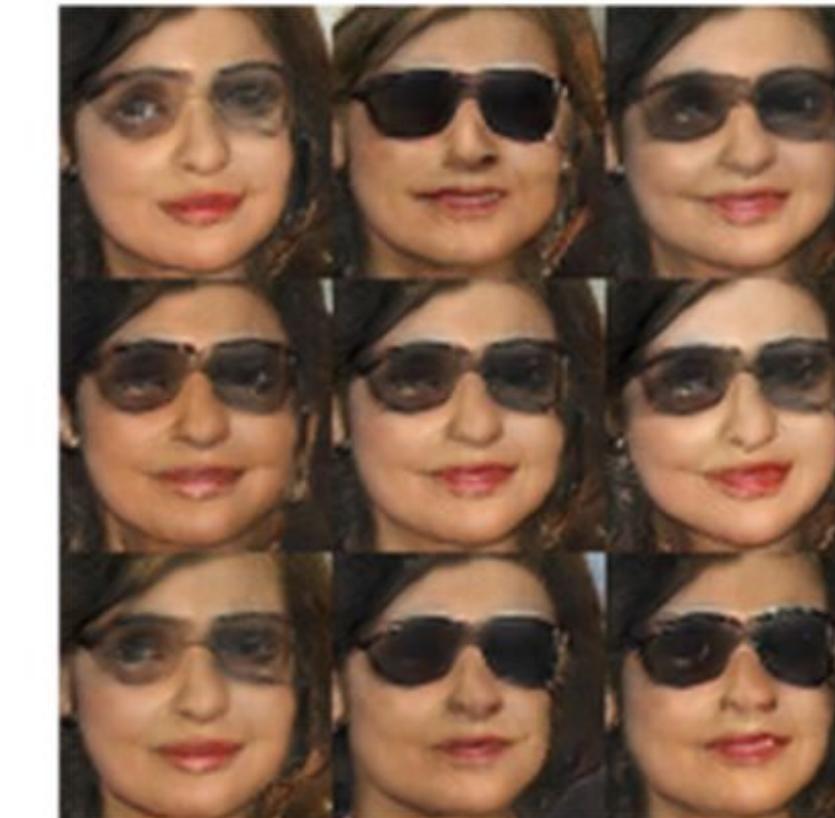


No glasses woman



Radford et al,
ICLR 2016

Woman with glasses



Generative Adversarial Networks

Glasses man



No glasses man

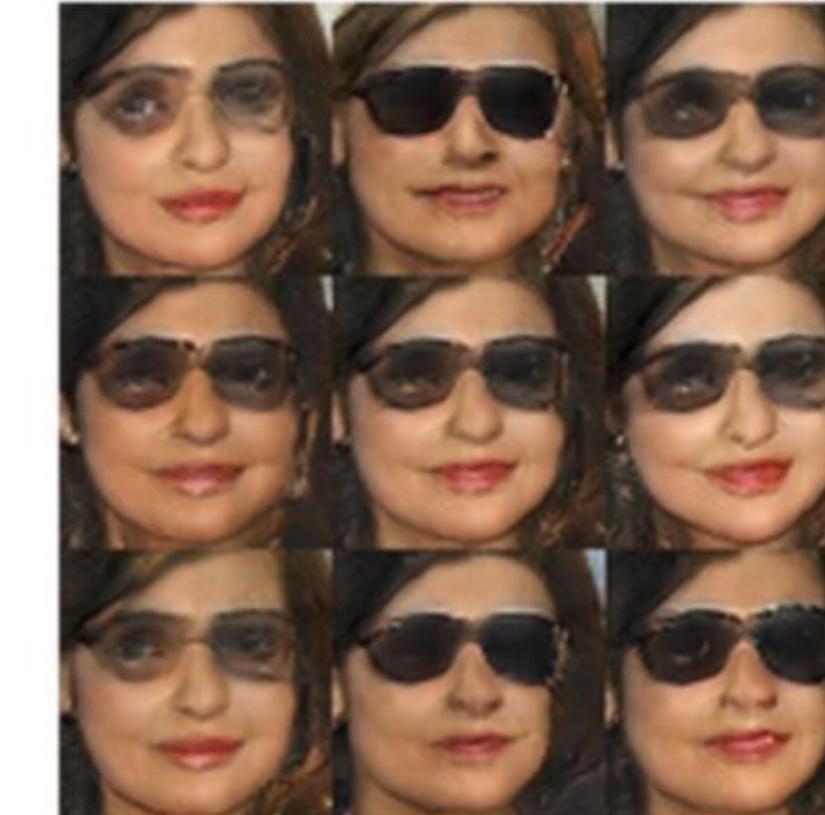


No glasses woman

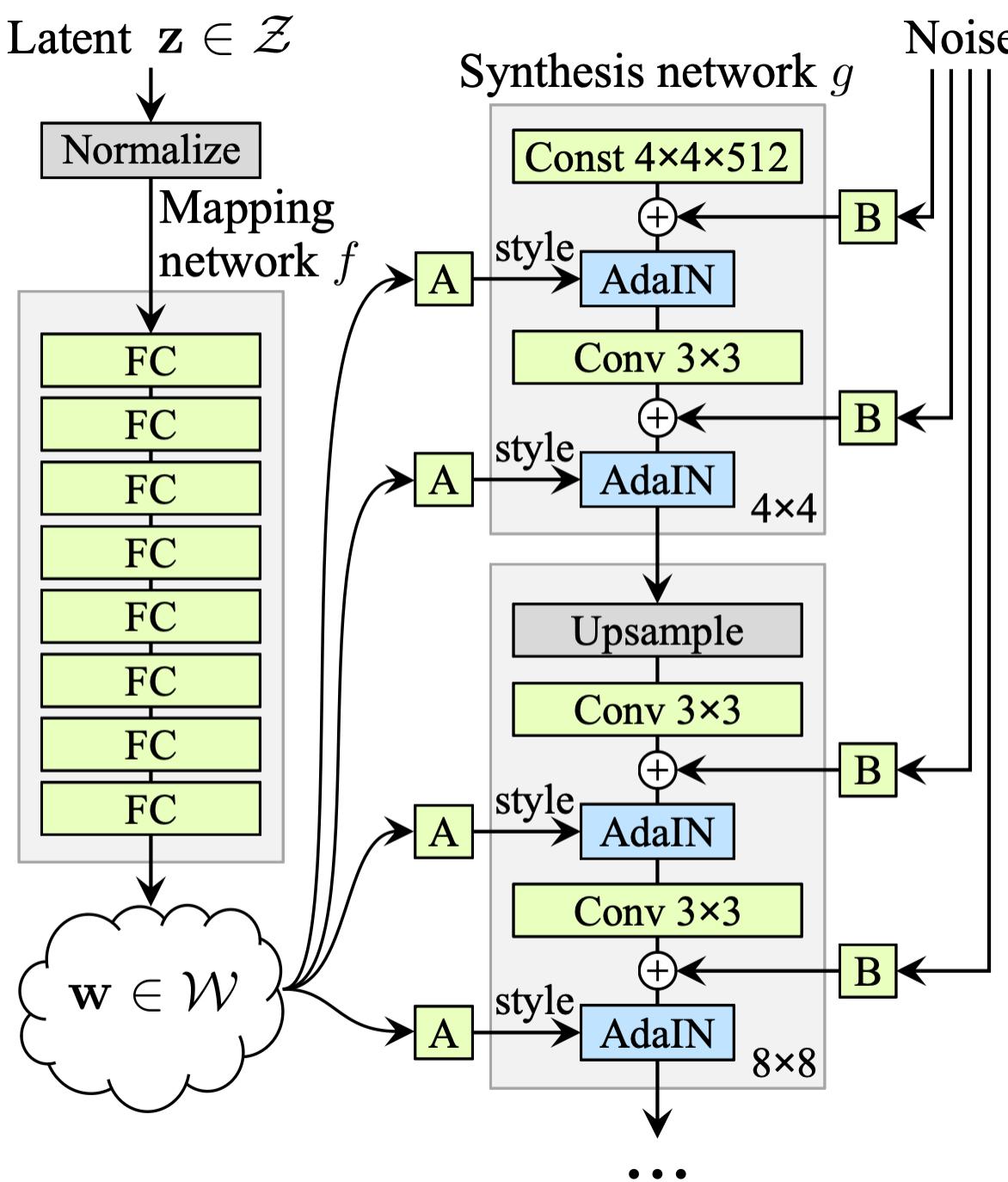


Radford et al,
ICLR 2016

Woman with glasses



StyleGAN (Karras et al. 2019)



The output image is **progressively built up** from 4x4 pixels, to 8x8 pixels, ..., to 1024x1024 pixels (Progressive GAN; Karras et al. 2018)

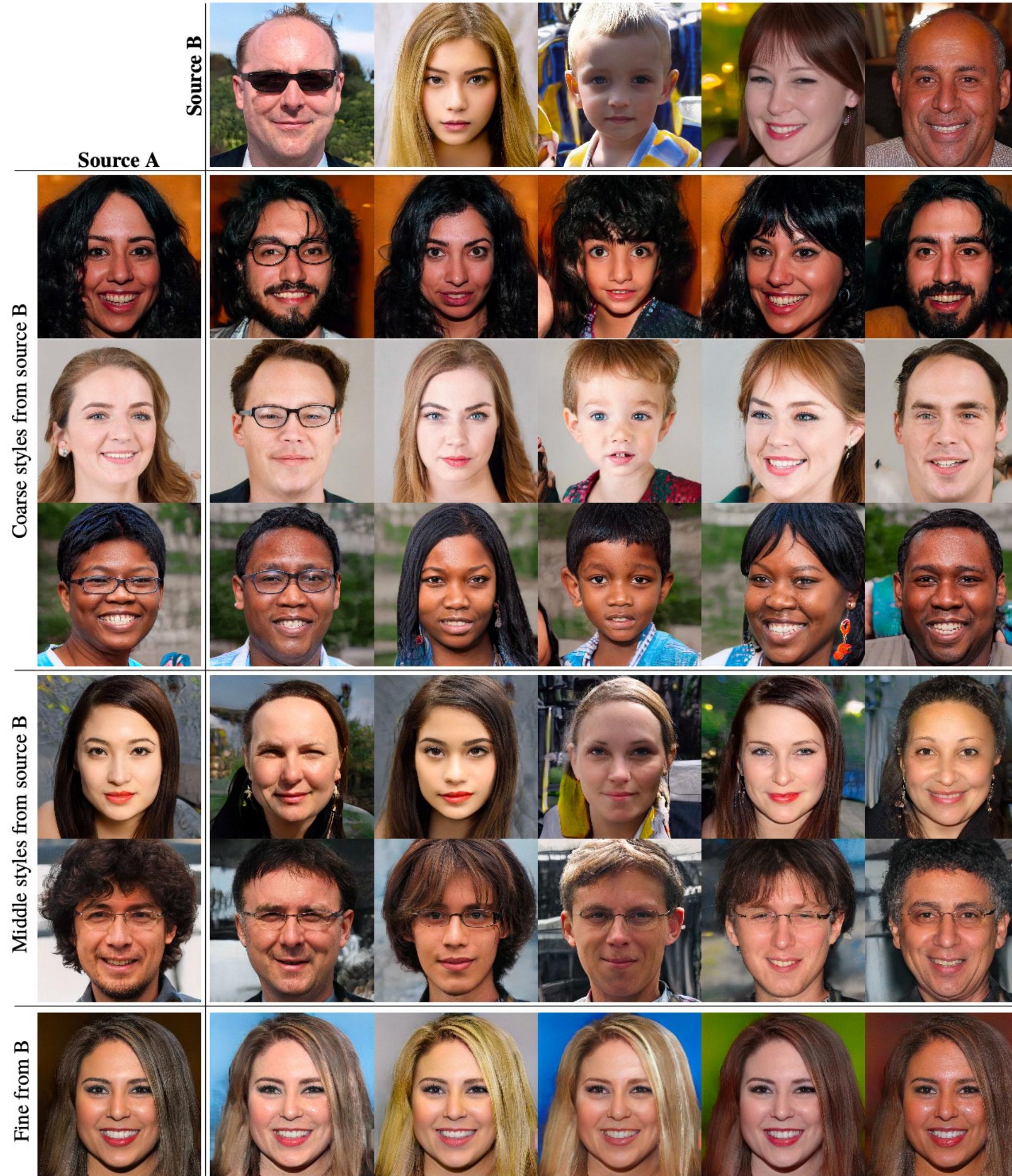
StyleGAN starts with a **fixed 4x4 pixel image** and **renormalizes the feature maps** at various points through image generation.

The **target mean values and standard deviations of the renormalization are computed from the latent representation** z by a convolutional neural network.

- The “early” normalizations define the high-level image content (“coarse style”)
- The “late” normalizations define the fine-level details of the image (“fine style”)

You can use coarse styles from one latent representation and fine style from another latent representation!

StyleGAN (Karras et al. 2019)



The output image is **progressively built up** from 4x4 pixels, to 8x8 pixels, ..., to 1024x1024 pixels (Progressive GAN; Karras et al. 2018)

StyleGAN starts with a **fixed 4x4 pixel image** and **renormalizes the feature maps** at various points through image generation.

The **target mean values and standard deviations of the renormalization are computed from the latent representation** z by a convolutional neural network.

- The “early” normalizations define the high-level image content (“coarse style”)
- The “late” normalizations define the fine-level details of the image (“fine style”)

You can use coarse styles from one latent representation and fine style from another latent representation!

Generative Adversarial Networks (GANs)

- GANs consists of **two competing components**
 - Generator tries to generate representative images
 - Discriminator tries to identify generated images
 - Components are optimized by **separate backpropagation steps**
- GANs typically produce **more photo-realistic images than variational autoencoders** (VAEs)
- However, the **training of GANs is very challenging**
 - If the discriminator becomes too good:
vanishing gradients for generator
 - If the generator becomes too good:
mode collapse for discriminator
 - **Limited variety** is one of the main limitations of GANs