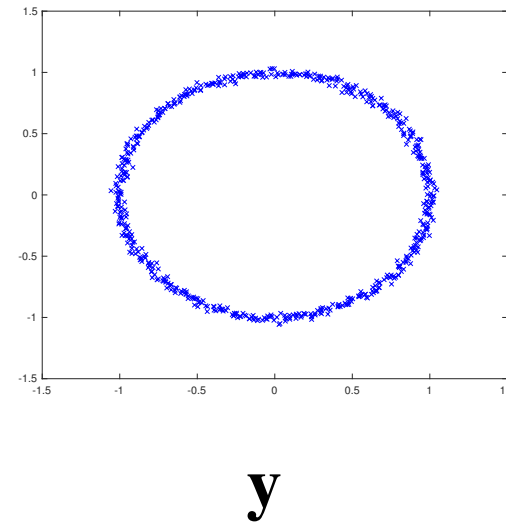
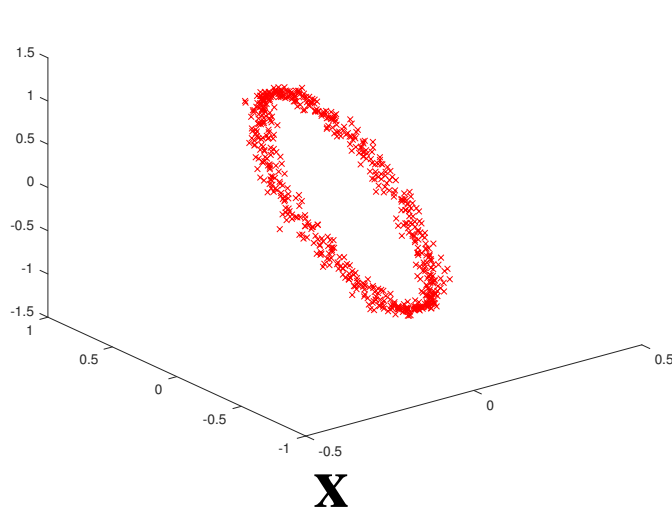


Non Linear Dimensionality Reduction

Pascal Fua
IC-CVLab

Reminder: Dimensionality Reduction

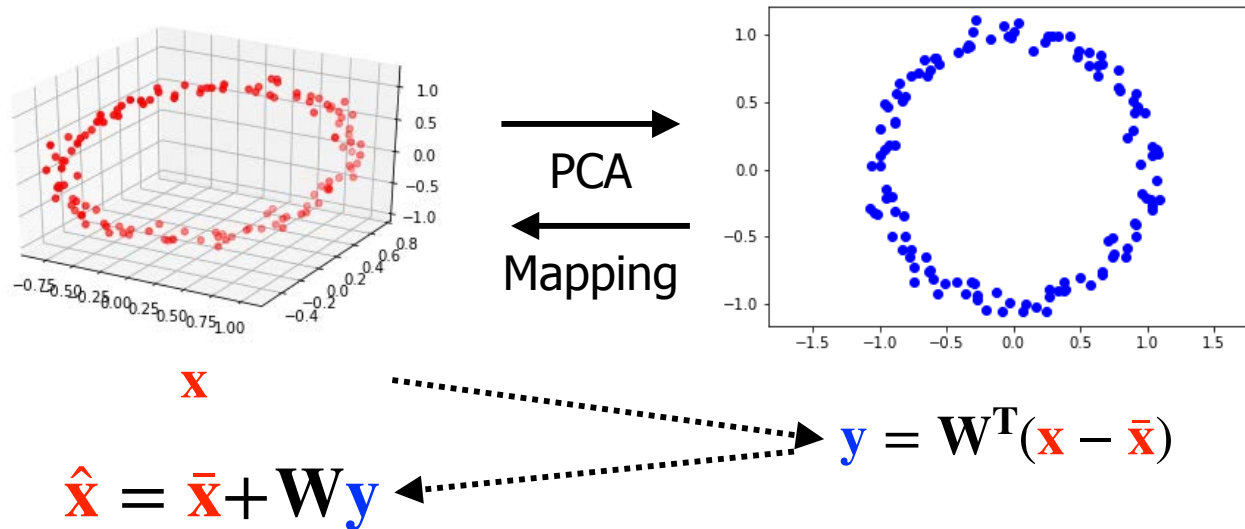


Our goal is to find a mapping $\mathbf{y}_i = f(\mathbf{x}_i)$

- $\mathbf{x}_i \in \mathbb{R}^D$: High-dimensional data sample
- $\mathbf{y}_i \in \mathbb{R}^d$: Low-dimensional representation

How about a linear one $\mathbf{y}_i = \mathbf{W}^T \mathbf{x}_i$?

Reminder: Optimal Linear Mapping



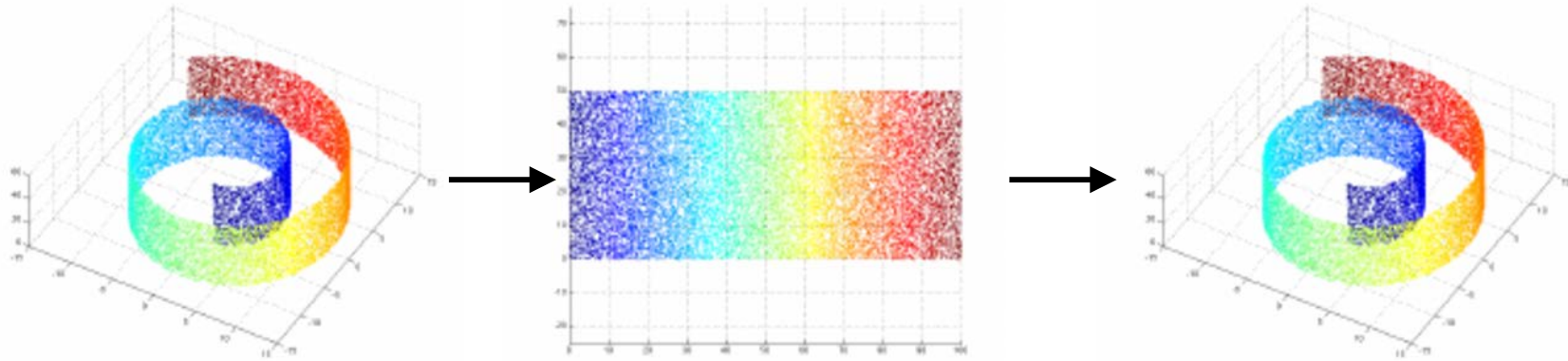
- This mapping incurs some loss of information.
- However, the corresponding rectangular matrix \mathbf{W} is the orthogonal matrix that minimizes the reconstruction error

$$e = \|\hat{\mathbf{x}} - \mathbf{x}\|^2$$

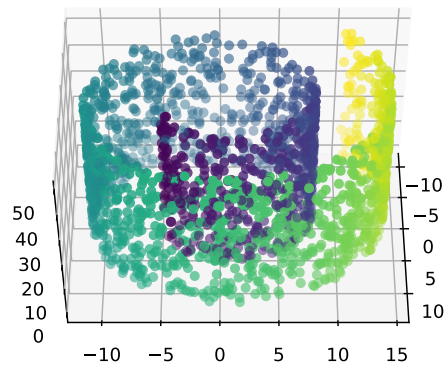
where

$$\hat{\mathbf{x}} = \bar{\mathbf{x}} + \mathbf{W}\mathbf{y} = \bar{\mathbf{x}} + \mathbf{W}\mathbf{W}^T(\mathbf{x} - \bar{\mathbf{x}})$$

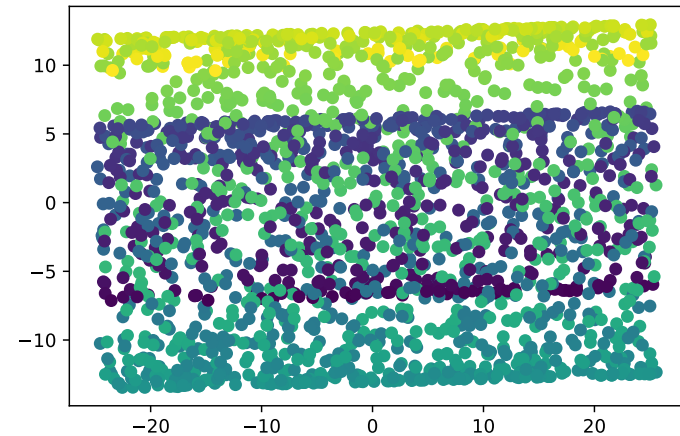
Limitation of the Linear Model



What it should do.



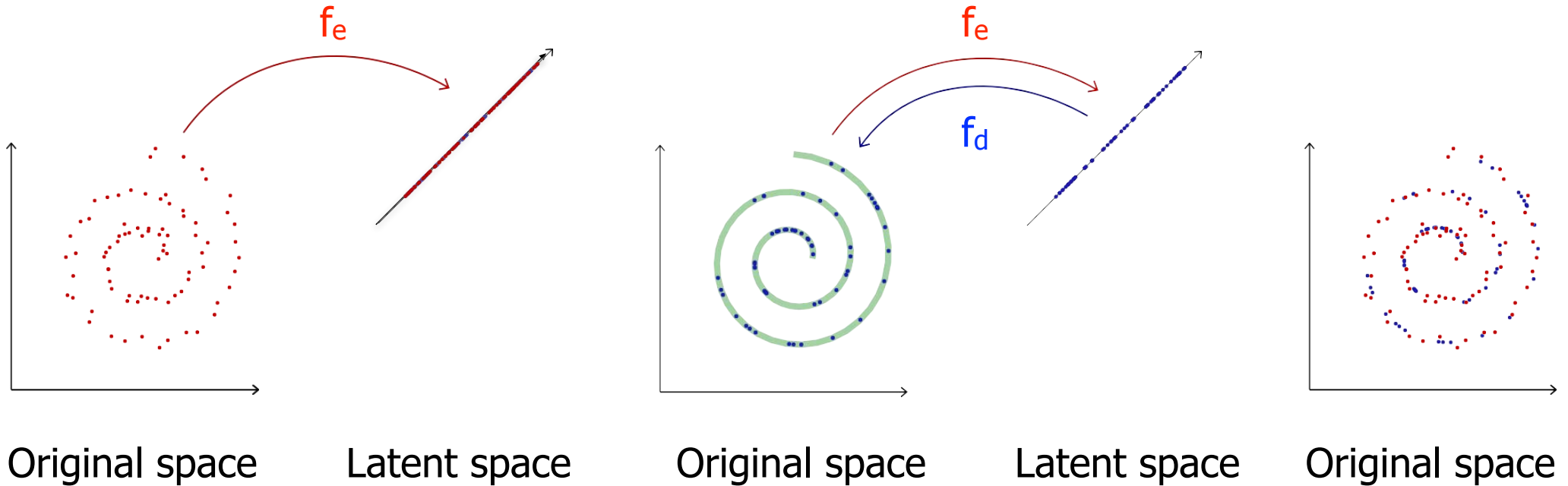
PCA
→



What PCA does.

—> Can we make the mappings non-linear so that cases like this one can be properly handled?

Latent Space



$$\mathbf{z} = f_e(\mathbf{x})$$

$$\hat{\mathbf{x}} = f_d(\mathbf{z})$$

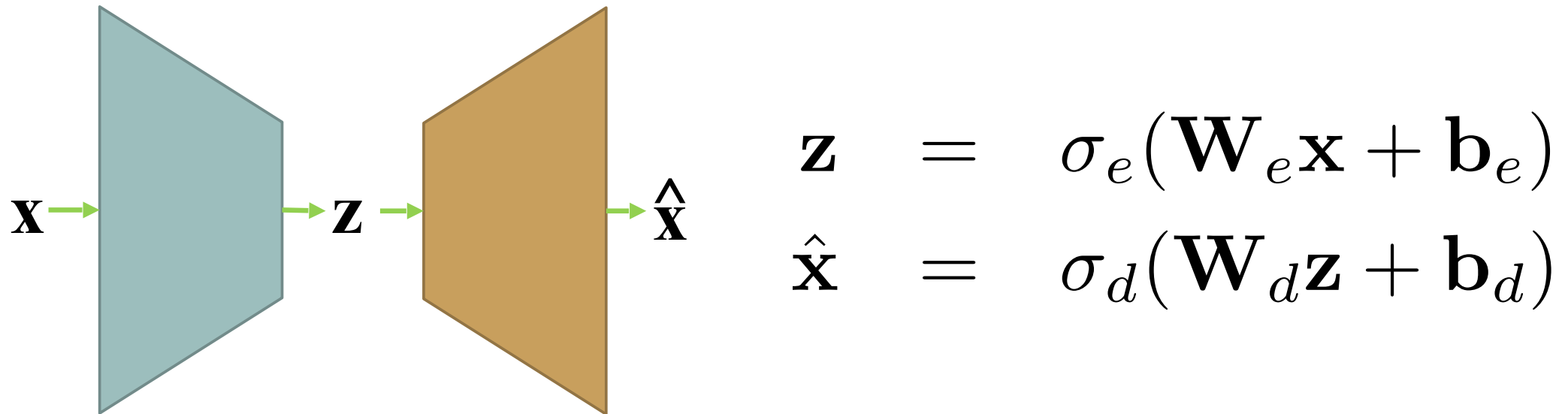
$$\hat{\mathbf{x}} \approx \mathbf{x}$$

- Removes unnecessary degrees of freedom.
- Models correlations between the true ones.
- Makes it possible to denoise the original data.

Linear vs NonLinear

- In PCA and LDA, the functions f_e and f_d are linear.
- We have seen that for classification non-linear functions are often more effective.
- Can we learn non-linear f_e and f_d ?
- There are many ways to do this but we will focus on special purpose deep networks called autoencoders.

Basic Autoencoder

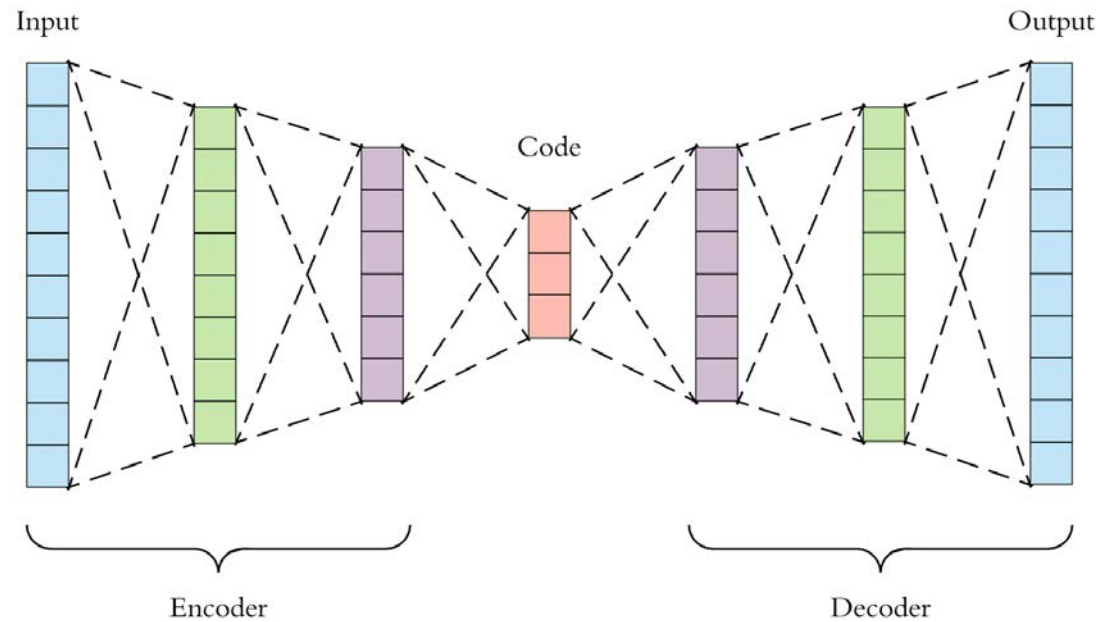


- \mathbf{z} is the **latent vector** representation of \mathbf{x} .
- $\hat{\mathbf{x}}$ is the **reconstruction** of \mathbf{x} . It should be as similar to it as possible.
- \mathbf{W}_e and \mathbf{W}_d can be computed by minimizing

$$\sum_n \|\hat{\mathbf{x}}_n - \mathbf{x}_n\|^2$$

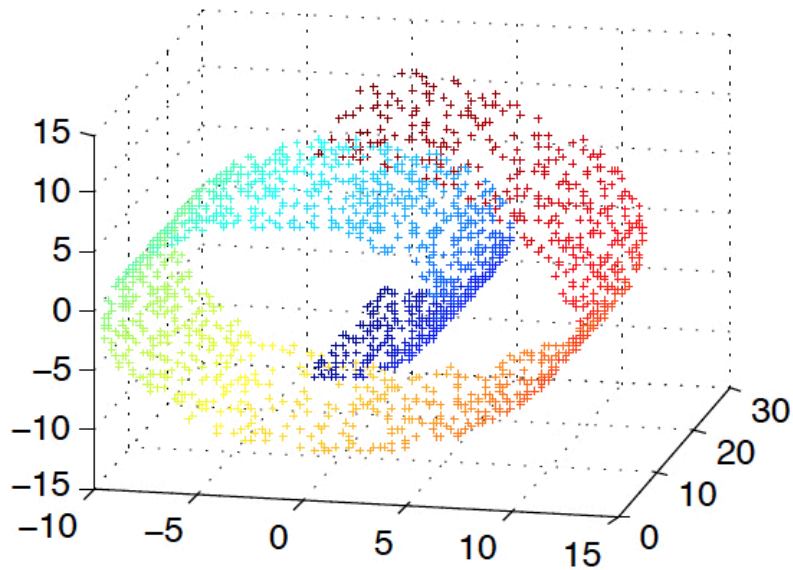
—> Unsupervised learning.

Deep Autoencoder

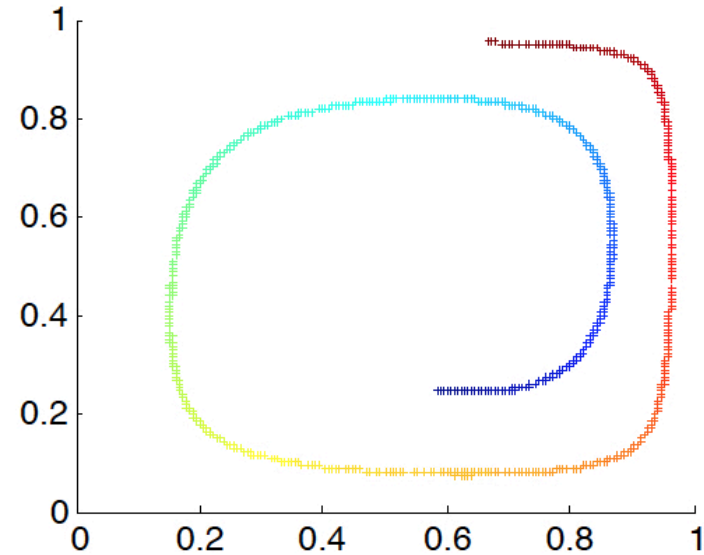


- There is no reason to restrict the encoder and decoder to consist of a single layer
- We can then create deep autoencoders by stacking multiple layers, each with an activation function.

Embedding the Swiss Roll

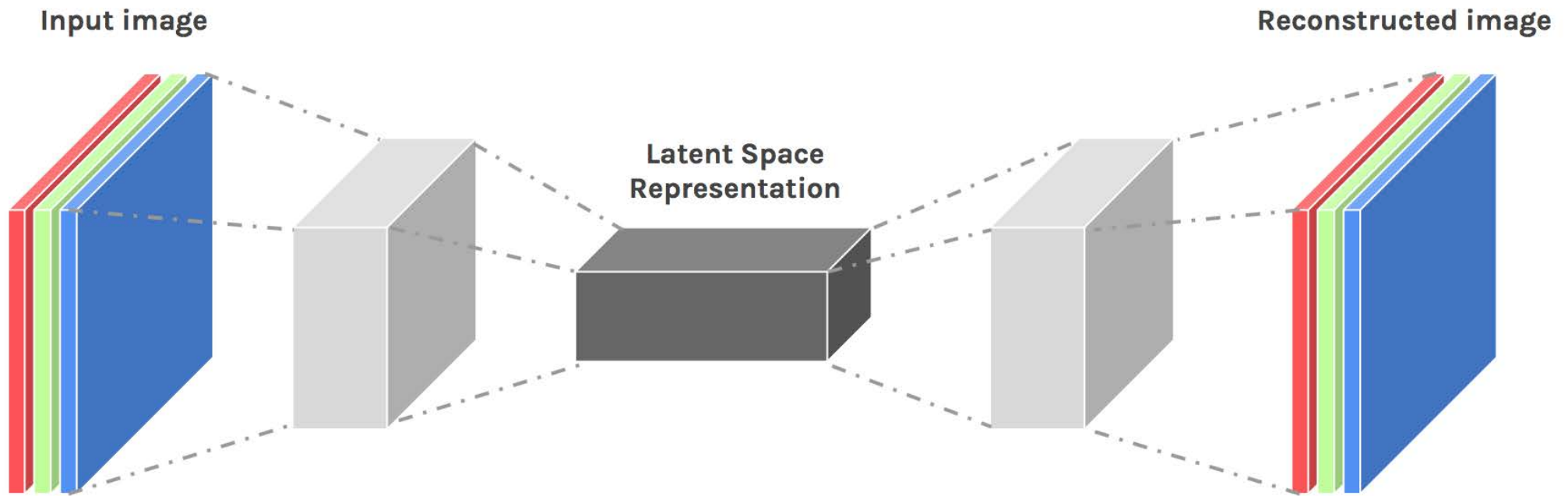


Input



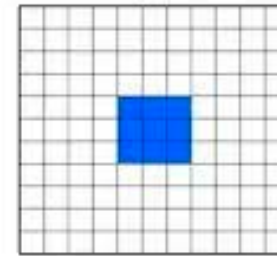
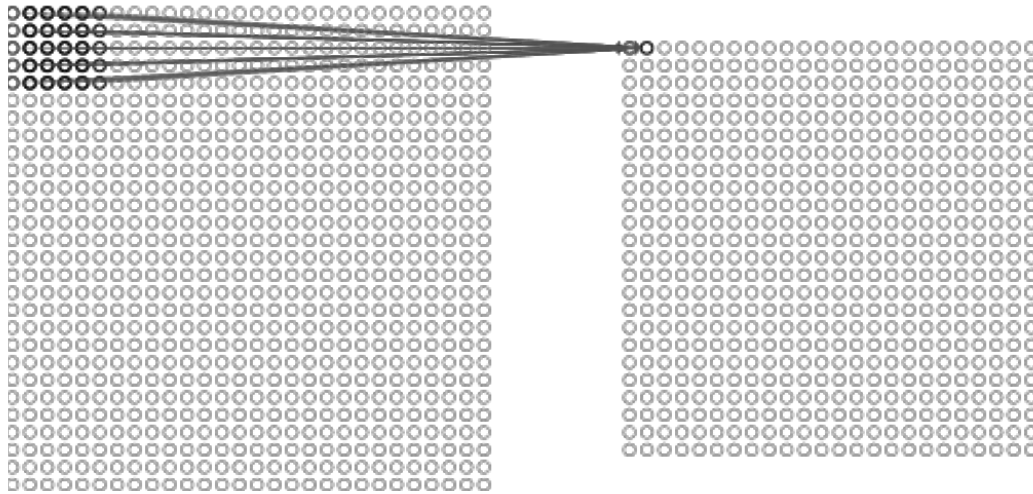
Latent representation

Convolutional Autoencoder

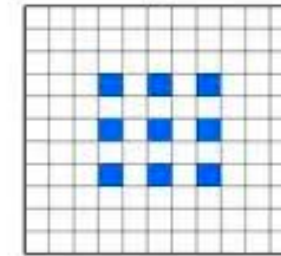


- When dealing with images, the layers can be convolutional ones.

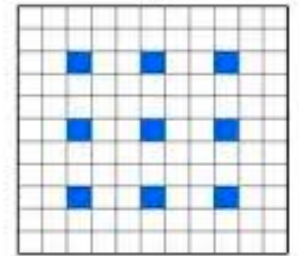
Reminder: 2D Convolutional Layer



stride=1



stride=2

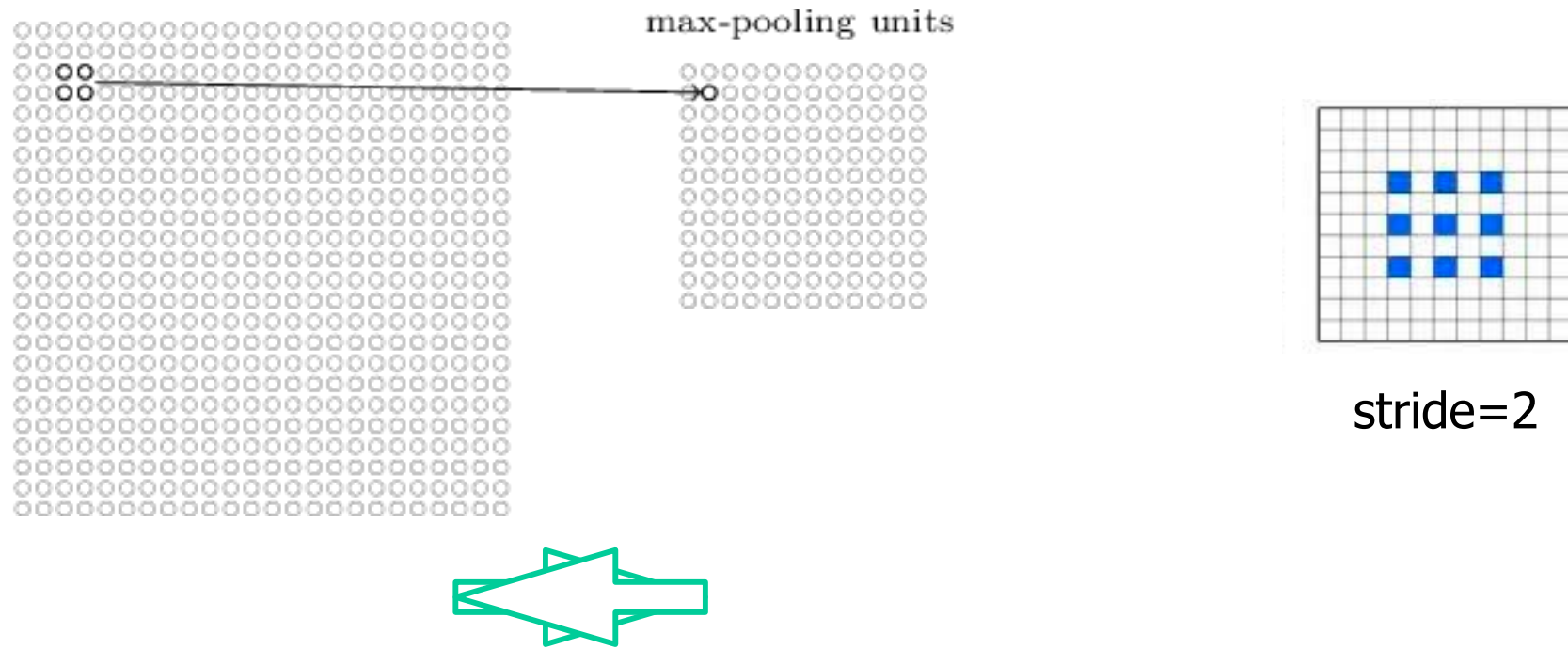


stride=3

$$\sigma \left(b + \sum_{x=0}^{n_x} \sum_{y=0}^{n_y} w_{x,y} a_{i+x,j+y} \right)$$

- Using a stride > 1 reduces the image size.

Reminder: 2D Convolutional Layer



$$\sigma \left(b + \sum_{x=0}^{n_x} \sum_{y=0}^{n_y} w_{x,y} a_{i+x,j+y} \right)$$

- Using a stride > 1 reduces the image size.
- Transposed convolution can be used to increase the image size.

PyTorch Translation

```
class Autoencoder(nn.Module):
```

```
    def __init__(nChannel=10,nHidden=50):
```

```
        self.convE1 = nn.Conv2d(1, nChannel, kernel_size=5,padding=2)
```

```
        self.convE2 = nn.Conv2d(nc, nChannel, kernel_size=5,padding=2)
```

```
        self.convD1 = nn.ConvTranspose2d(nChannel,nChannel,kernel_size=4,stroke=2,padding=1)
```

```
        self.convD2 = nn.ConvTranspose2d(nChannel,1 ,kernel_size=4,stroke=2,padding=1)
```

```
    def encode(self,x):
```

```
        x = self.convE1(x)
```

```
        x = sigma(F.max_pool2d(x,2))
```

```
        x = self.convE2(x)
```

```
        x = sigma(F.max_pool2d(x,2))
```

```
        return x
```

```
    def decode(self, z):
```

```
        z= sigma(self.convD1(z))
```

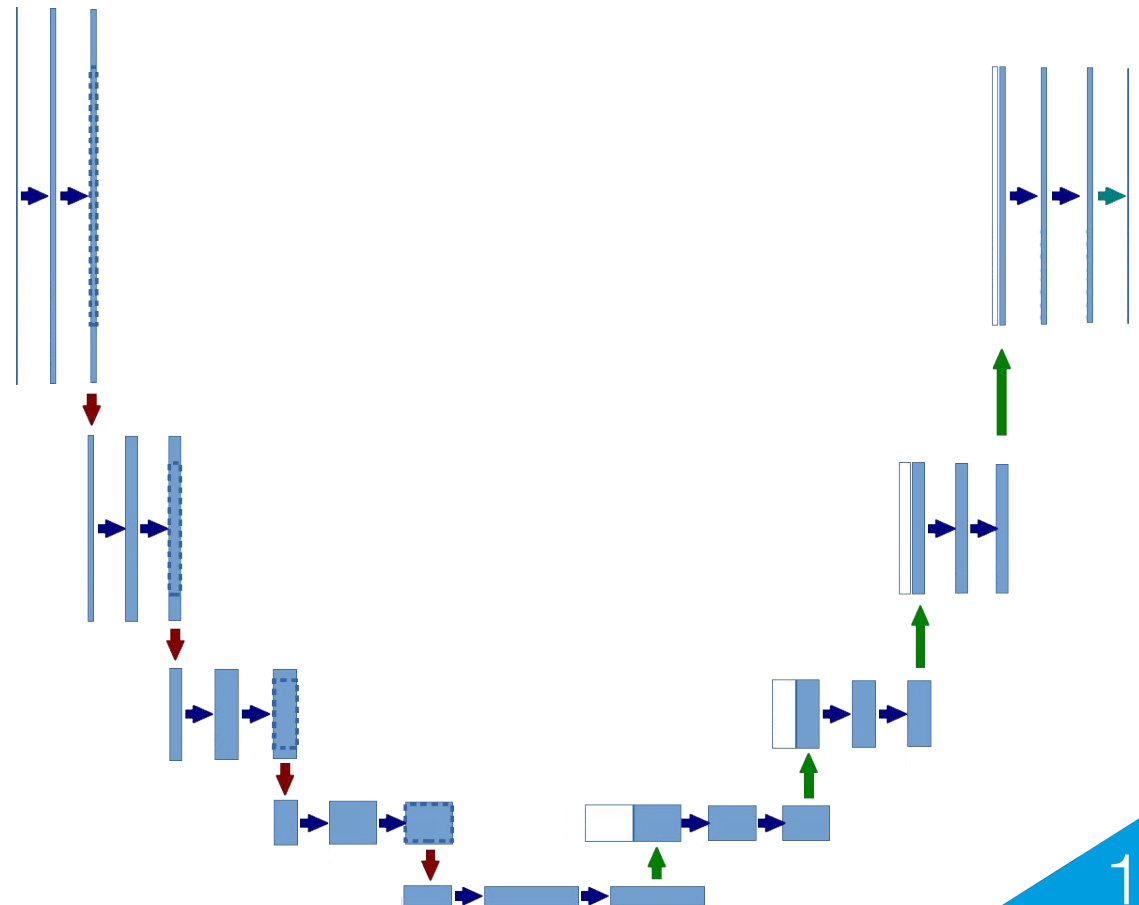
```
        z= self.convD2(z)
```

```
        return z
```

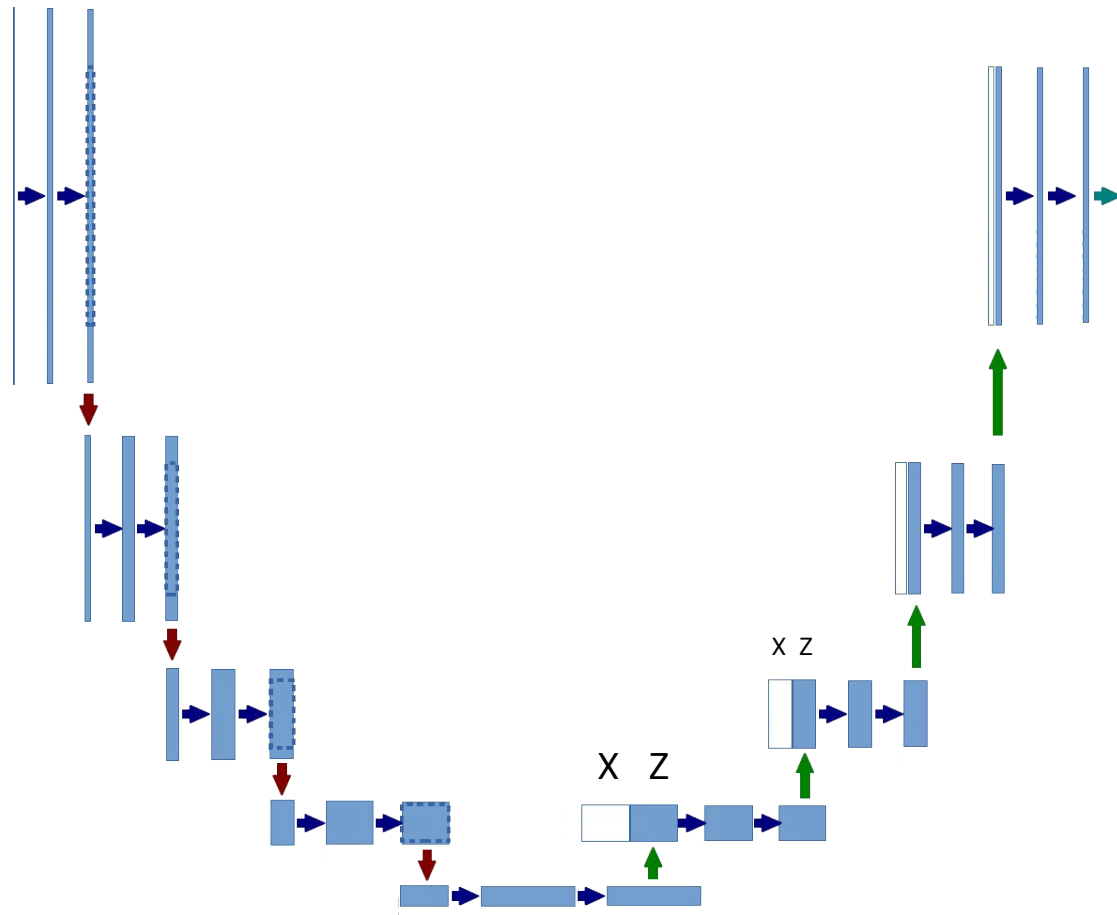
```
    def forward(self,x):
```

```
        z = self.encode(x)
```

```
        return self.decode(z)
```



Reminder: From Autoencoder to UNet



```
def encode(self,x):
```

```
    x = self.convE1(x)
```

```
    x = sigma(F.max_pool2d(x,2))
```

```
    z = self.convE2(x)
```

```
    z = sigma(F.max_pool2d(z,2))
```

```
    return z
```

```
def decode(self, z):
```

```
    z = sigma(self.convD1(z))
```

```
    z = self.convD2(z)
```

```
    return z
```

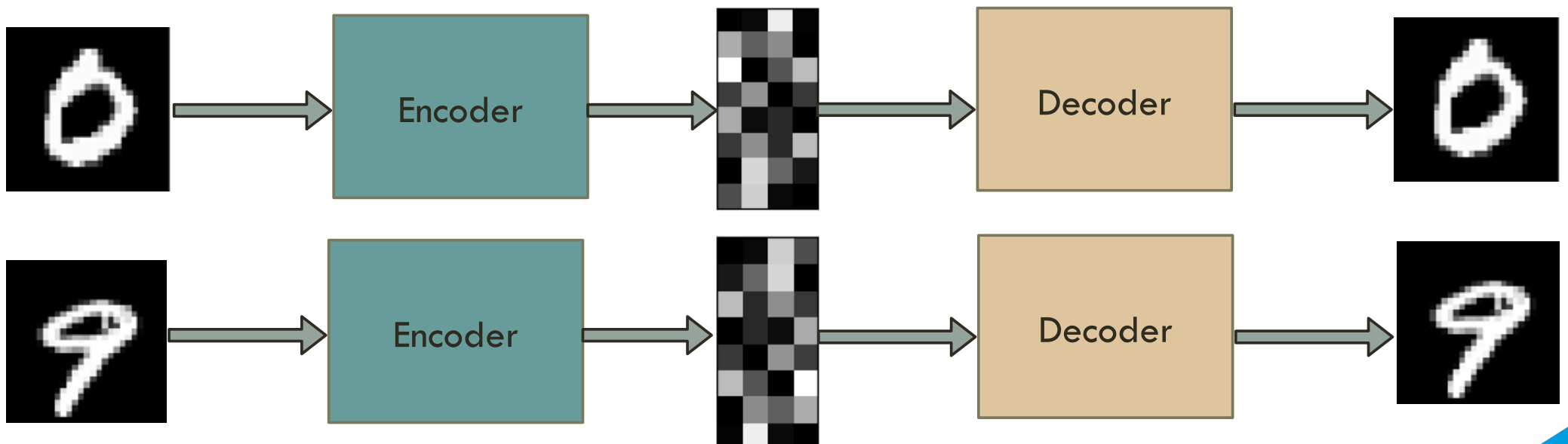
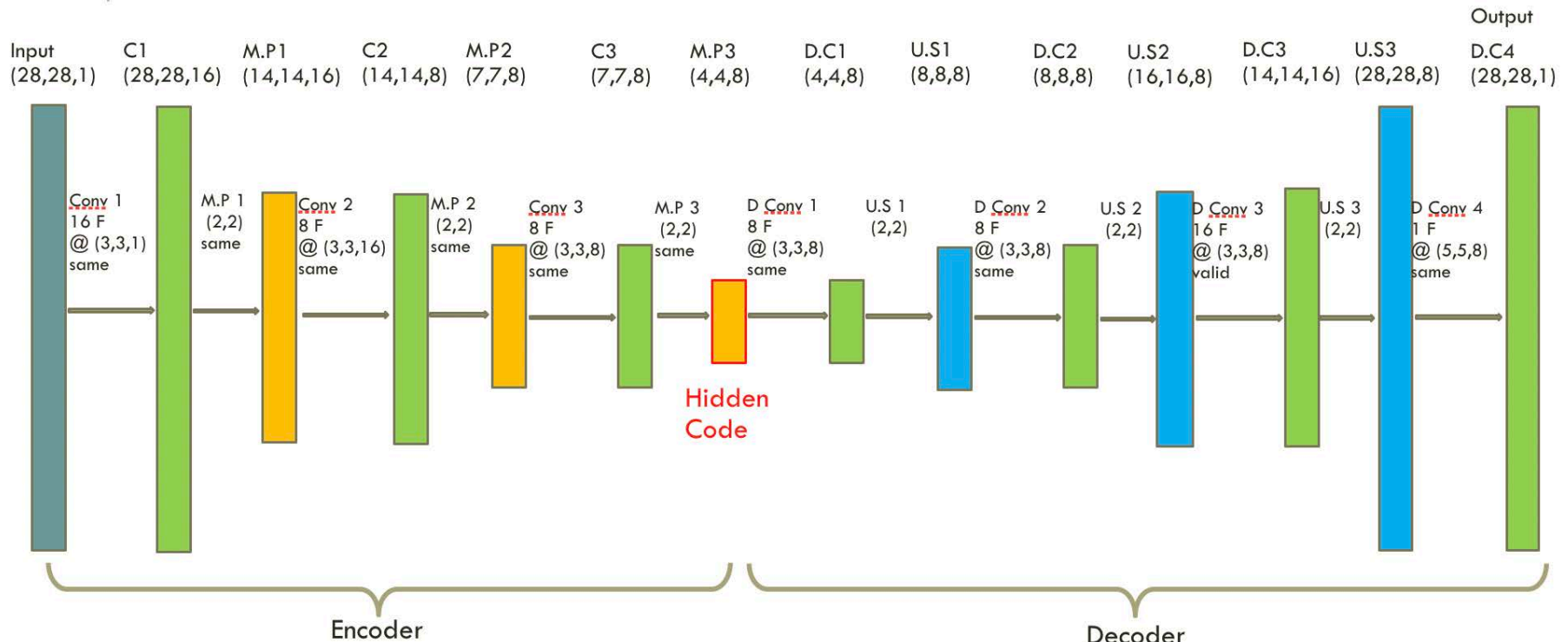
```
def forward(self,x):
```

```
    z = self.encode(x)
```

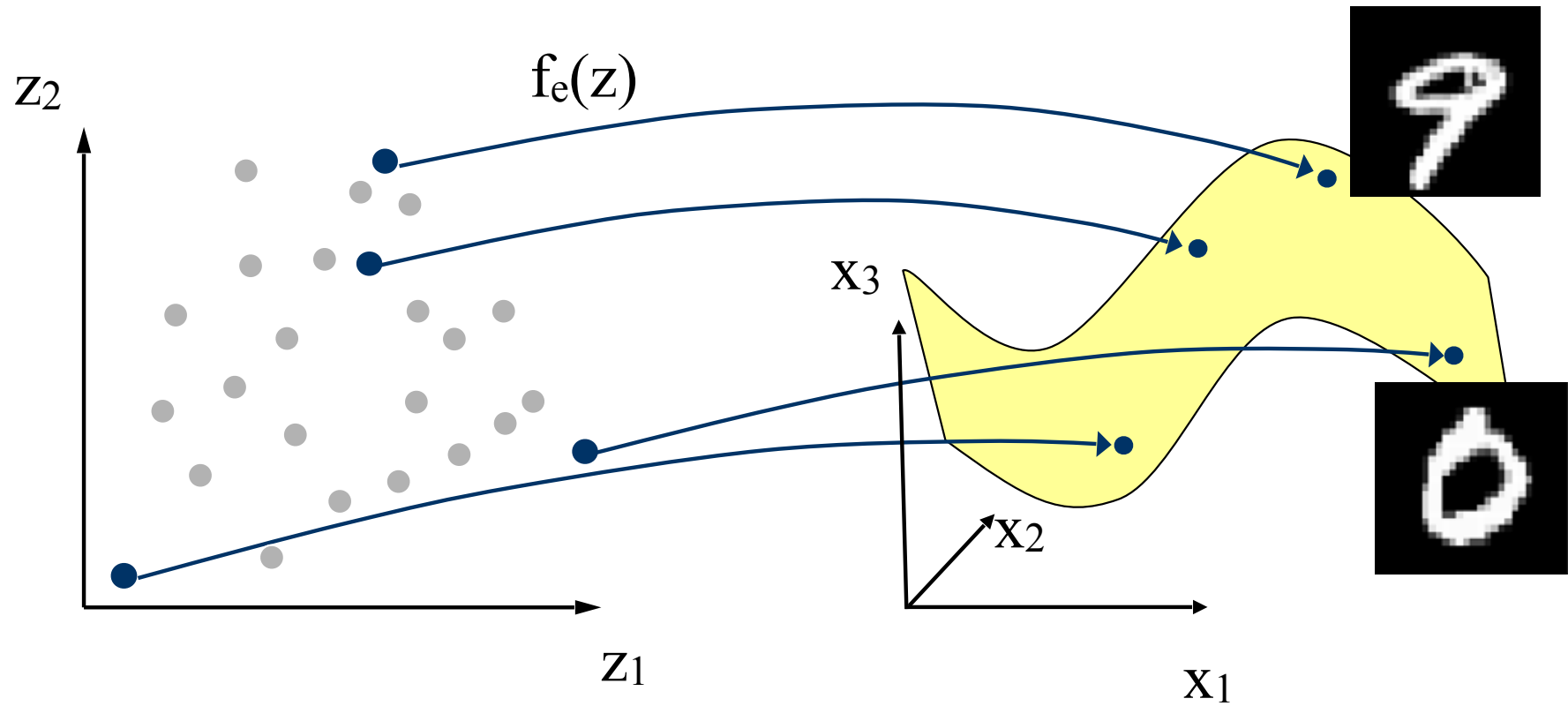
```
    return self.decode(z)
```

Concatenate z on the way up with the equivalent x on the way down, before running self.convD1 and self.convD2

MNIST



Latent Space



Exploring the Latent Space



Wandering through \mathbf{z} -space and observing the effects in \mathbf{x} -space

Encoding / Decoding

X (original samples)

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

$gf(X)$ (CNN, $d=2$)

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

$gf(X)$ (PCA, $d=2$)

9 3 1 0 9 1 9 9 0 9 0 0
9 0 1 3 9 9 3 9 9 0 9 0
9 0 9 9 0 1 3 1 3 0 9 0

Encoding / Decoding

X (original samples)

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

$gf(X)$ (CNN, $d=4$)

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

$gf(X)$ (PCA, $d=4$)

9 2 1 0 9 1 9 9 0 9 0 0
9 0 1 3 9 9 0 9 9 0 4 9
9 0 9 9 0 1 3 1 3 4 9 0

Encoding / Decoding

X (original samples)

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

$gf(X)$ (CNN, $d=8$)

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

$gf(X)$ (PCA, $d=8$)

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

Encoding / Decoding

X (original samples)

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

$g_f(X)$ (CNN, $d = 16$)

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

$g_f(X)$ (PCA, $d = 16$)

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

Encoding / Decoding

X (original samples)

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

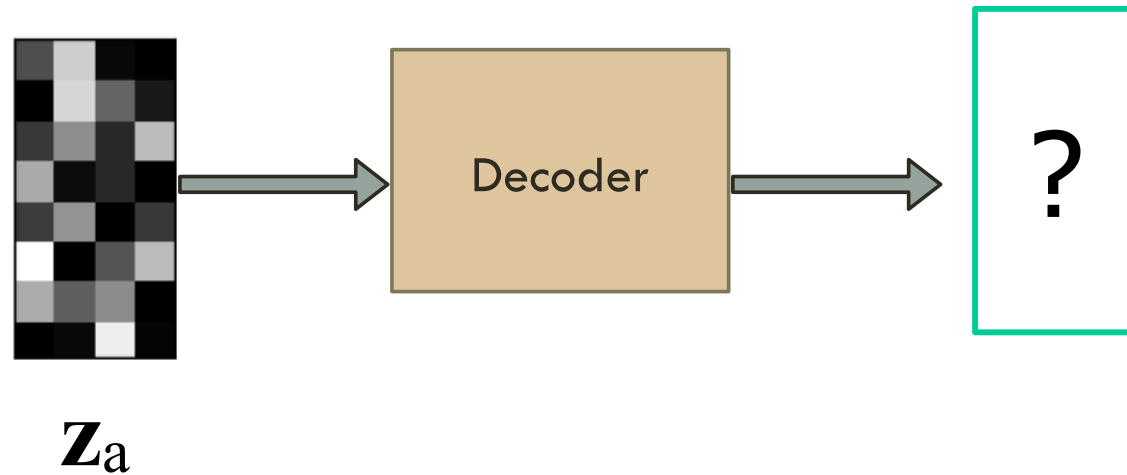
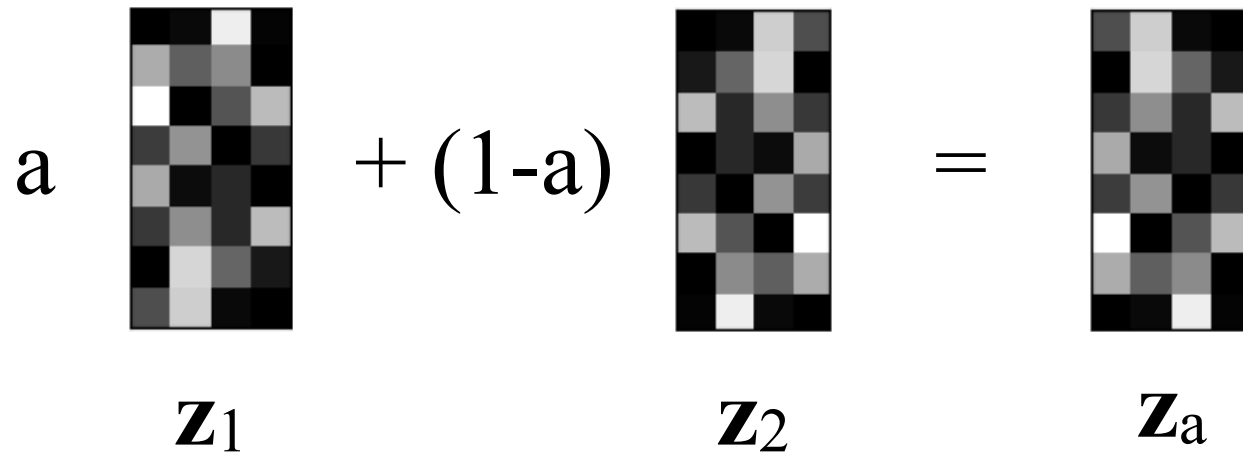
$g_f(X)$ (CNN, $d = 32$)

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

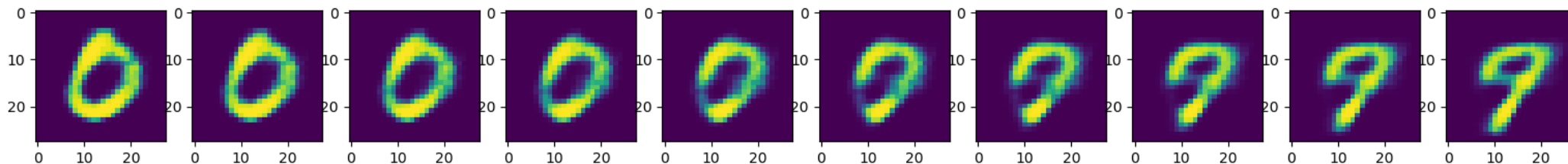
$g_f(X)$ (PCA, $d = 32$)

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

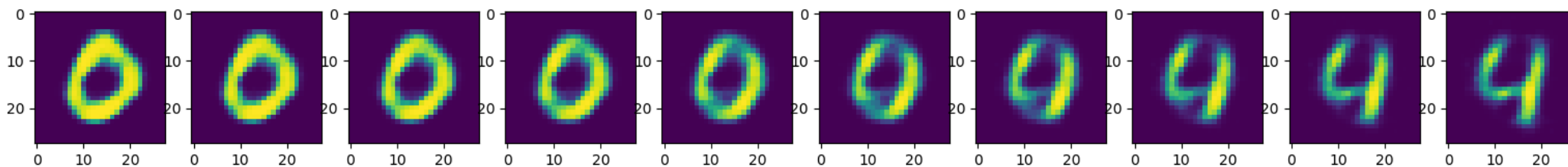
Interpolation



Interpolation



0 \longleftrightarrow 9



0 \longleftrightarrow 4

PCA vs Autoencoder Interpolation

1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
7 7 7 7 7 7 7 7 7 7 7 7

PCA (d=32)

1 1 1 1 1 1 1 1 1 1 1 1
2 2 2 2 2 2 2 2 2 2 2 2
3 3 3 3 3 3 3 3 3 3 3 3
0 0 0 0 0 0 0 0 0 0 0 0
1 1 1 1 1 1 1 1 1 1 1 1
7 7 7 7 7 7 7 7 7 7 7 7

Autoenc (d=8)

0 0 0 0 0 0 0 0 0 0 0 0
7 7 7 7 7 7 7 7 7 7 7 7
4 4 4 4 4 4 4 4 4 4 4 4
7 7 7 7 7 7 7 7 7 7 7 7
2 2 2 2 2 2 2 2 2 2 2 2
4 4 4 4 4 4 4 4 4 4 4 4

Autoenc (d=32)

Application to Image Retrieval

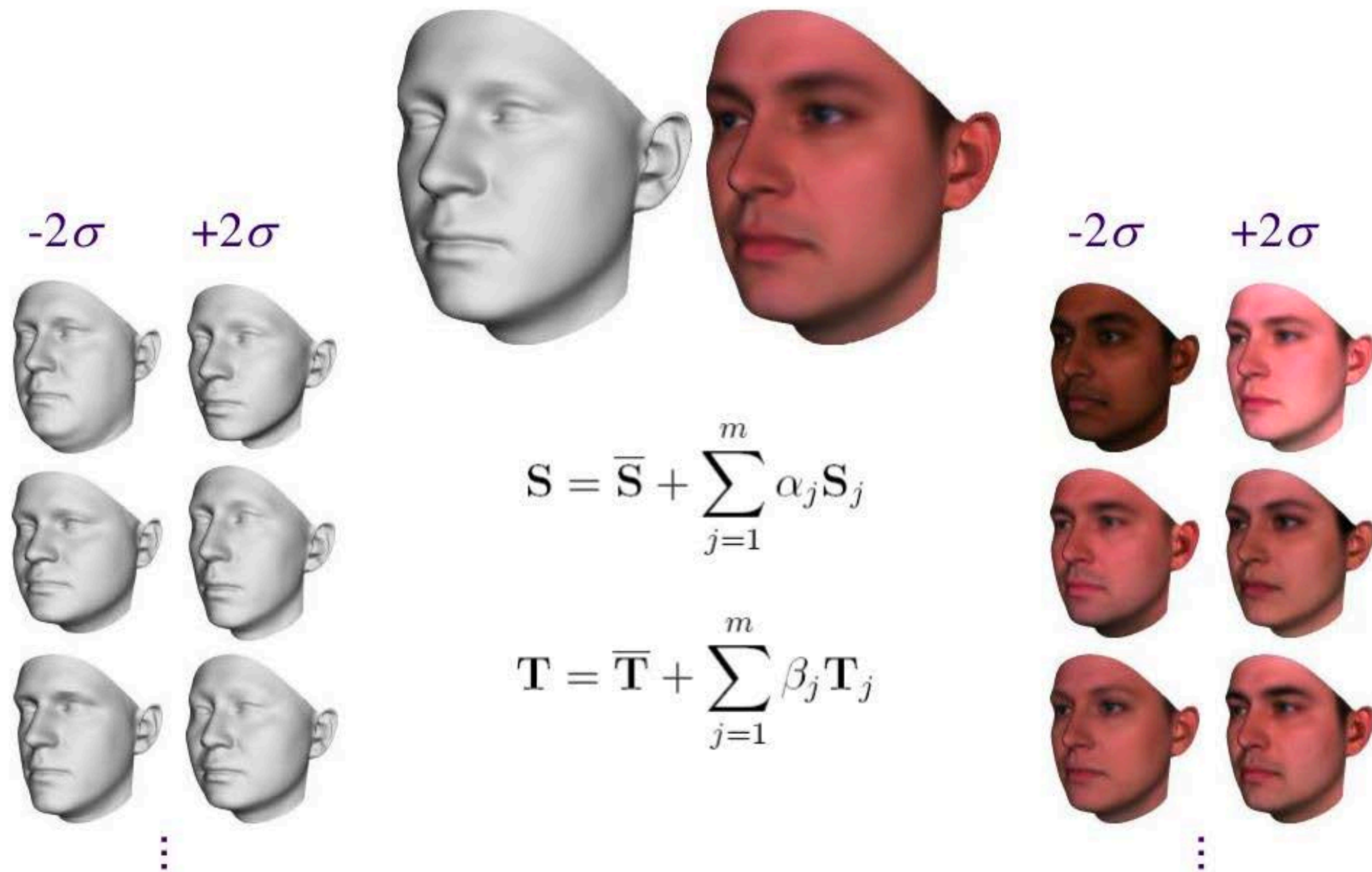
Image Retrieval ($k=5$)



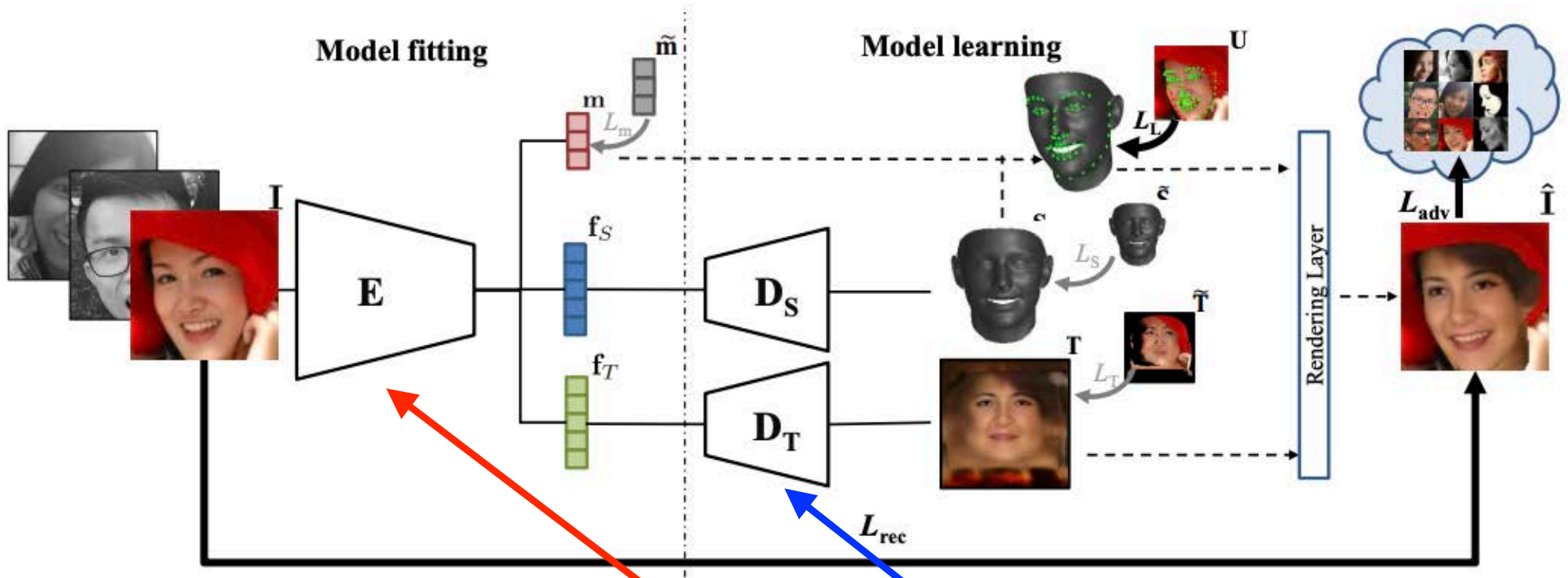
- A code provides a compact representation of the input.
- It can be used for retrieval in a large data collection, e.g., via by k nearest neighbors.

<https://towardsdatascience.com/find-similar-images-using-autoencoders-315f374029ea>

Reminder: PCA Face Modeling



Optional: Non Linear Face Models



- PCA has been replaced by an **encoder** / **decoder** architecture.
- To be discussed next week.

Optional: Linear vs Non Linear

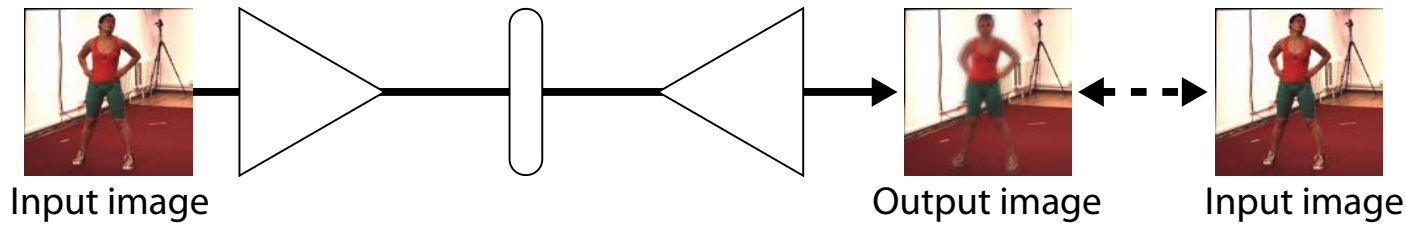


Original

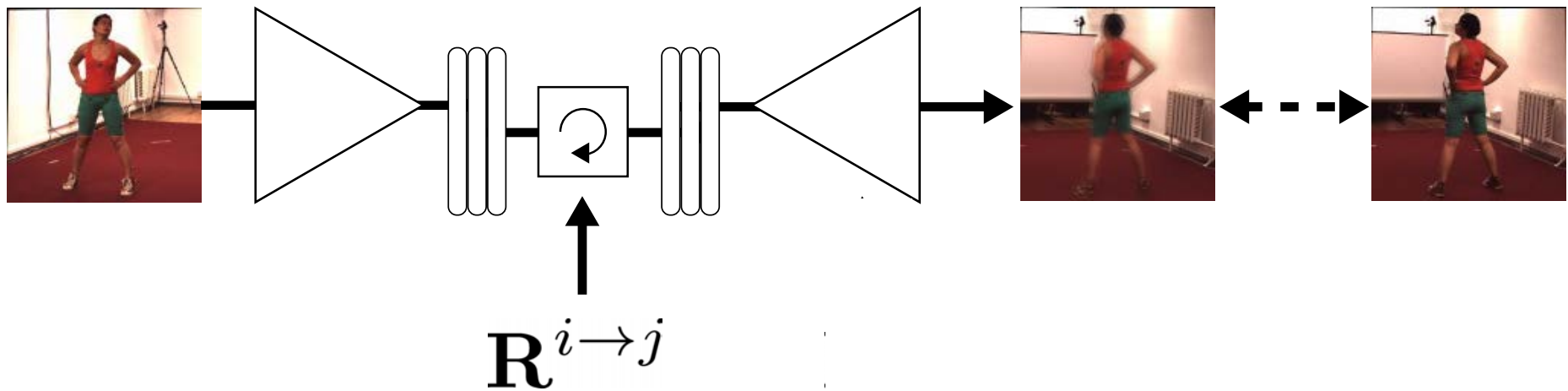
Linear

Non Linear

Optional: Novel View Synthesis

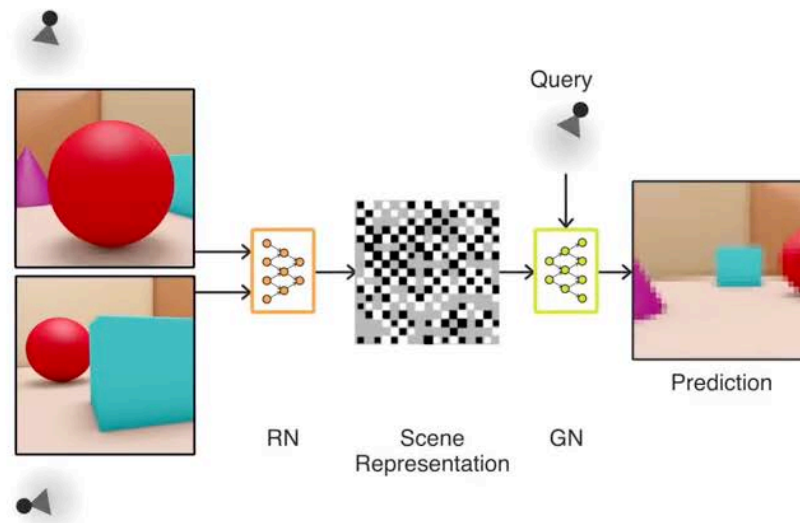


Conventional Autoencoder

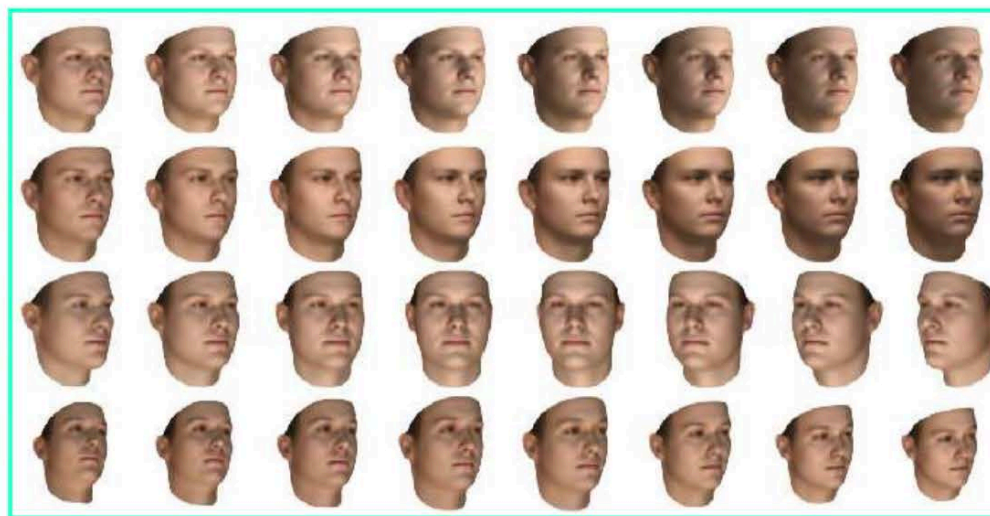
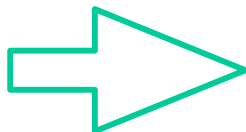
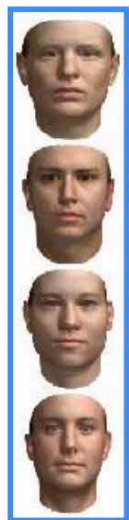


Rotation Aware Autoencoder

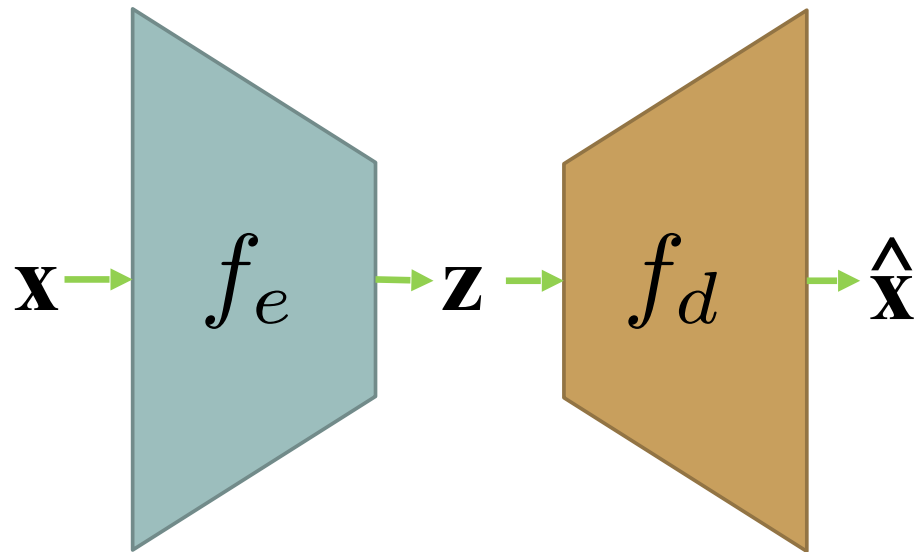
Optional: Novel View Synthesis



DeepMind

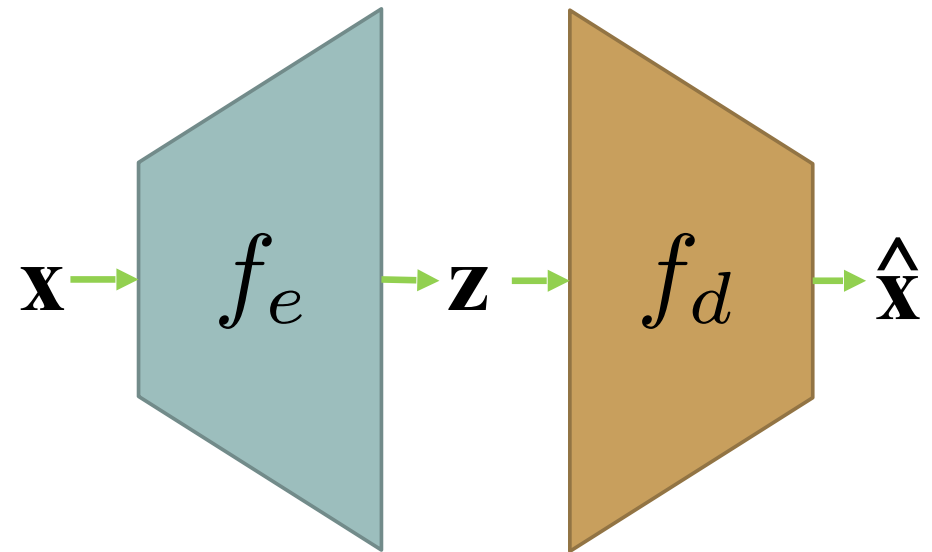


Under- vs Over-Complete



Undercomplete: $\dim(\mathbf{z}) < \dim(\mathbf{x})$

- Compresses the input.
- Captures correlations.

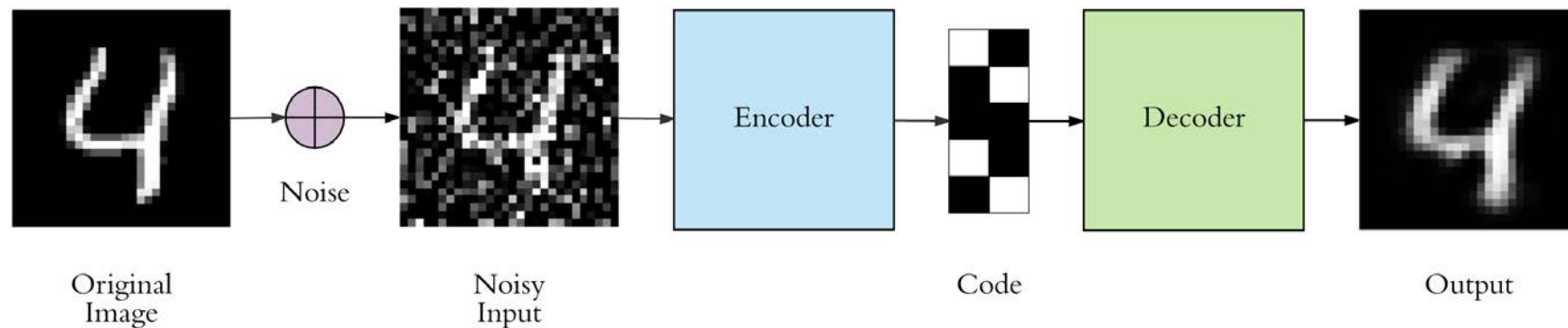


Overcomplete: $\dim(\mathbf{z}) > \dim(\mathbf{x})$

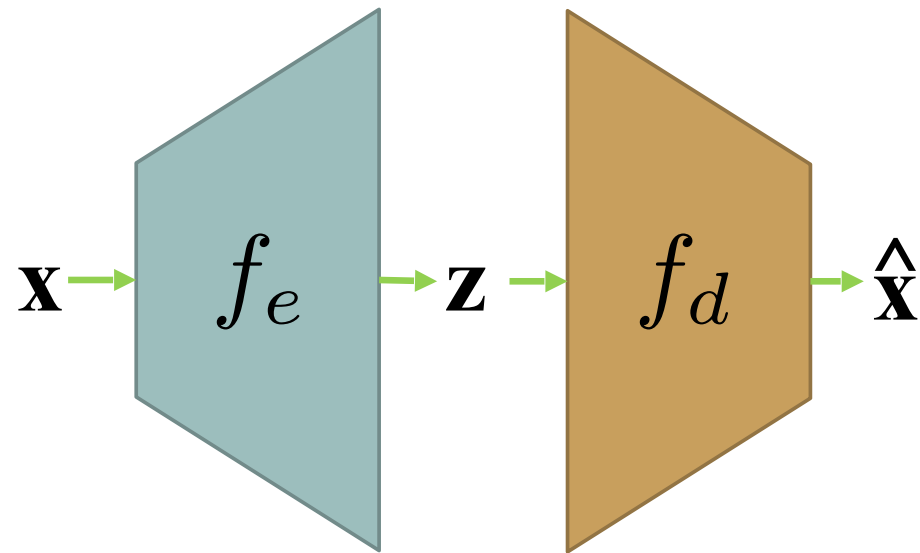
- Higher dimension can help.
- Degenerate solutions possible.
- Need a regularization term.

Denoising Autoencoders

- Having a low-dimensional latent representation encourages the autoencoder to learn an “intelligent” mapping.
- With dimensionality expansion, one could simply learn to copy the input.
- To avoid this, one can add noise to the input and aim to reconstruct a noise-free version of the input:



Contrastive Autoencoders



Overcomplete: $\dim(\mathbf{z}) > \dim(\mathbf{x})$

Another way to avoid trivial solutions is to train with a **regularization** term:

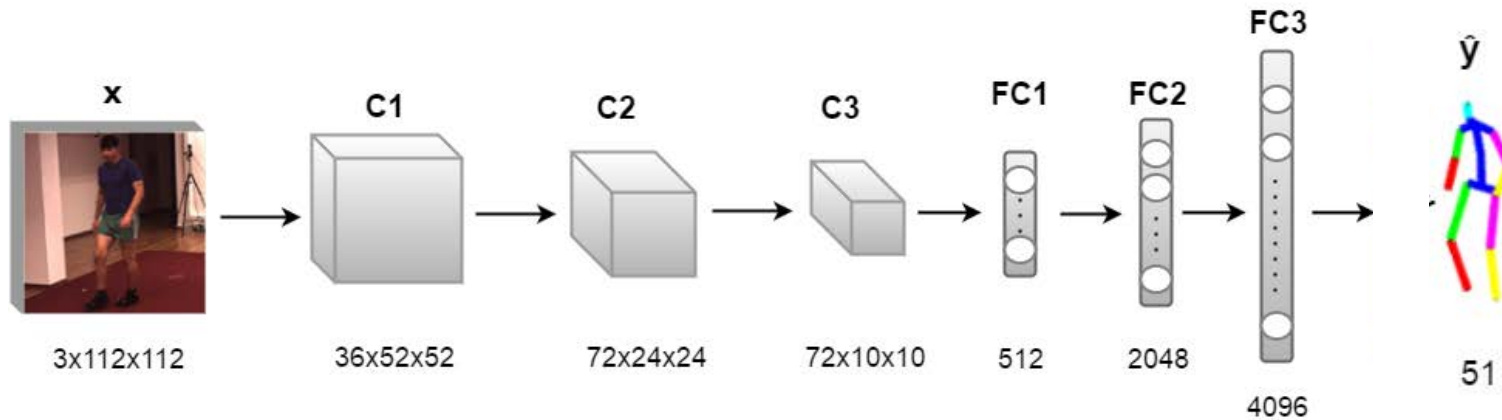
$$R(\mathbf{w}) = \sum_n L(\mathbf{x}_n, \mathbf{w}) + \lambda \Omega(\mathbf{x}_n, \mathbf{w}) ,$$

$$L(\mathbf{x}, \mathbf{w}) = \|f_d(f_e(\mathbf{x}, \mathbf{w})) - \mathbf{x}\|^2 ,$$

$$\Omega(\mathbf{w}) = \sum_{i,j} \left[\frac{\partial f_e(\mathbf{x}, \mathbf{w})_j}{\partial x_i} \right]^2 .$$

- When $\lambda \rightarrow 0$, $f_e \rightarrow$ identity.
- When $\lambda \rightarrow \infty$, $f_e \rightarrow$ constant.

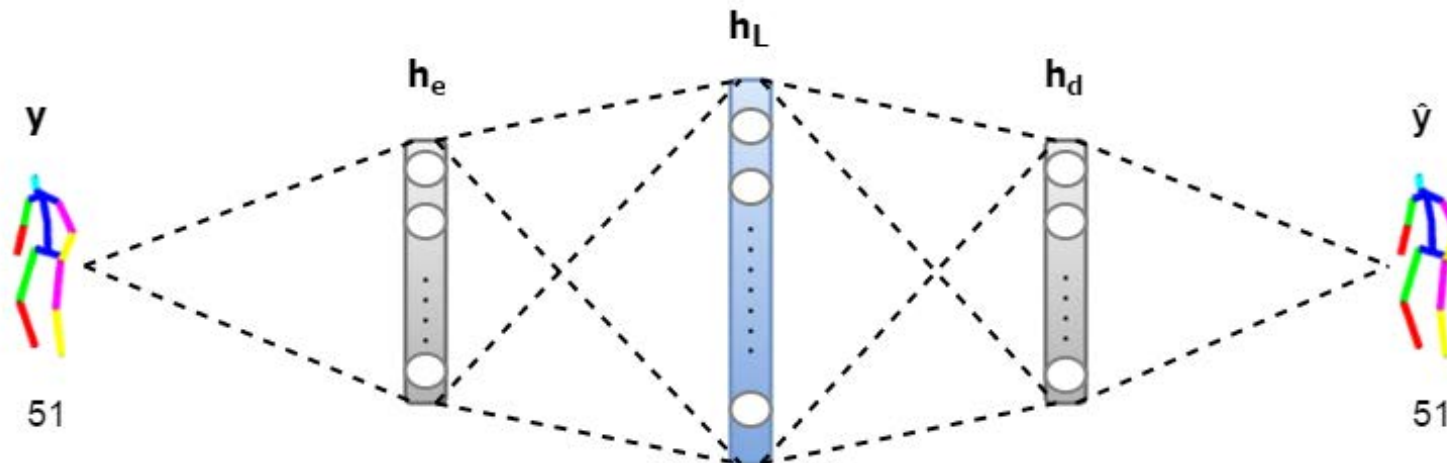
Optional: Body Pose Estimation



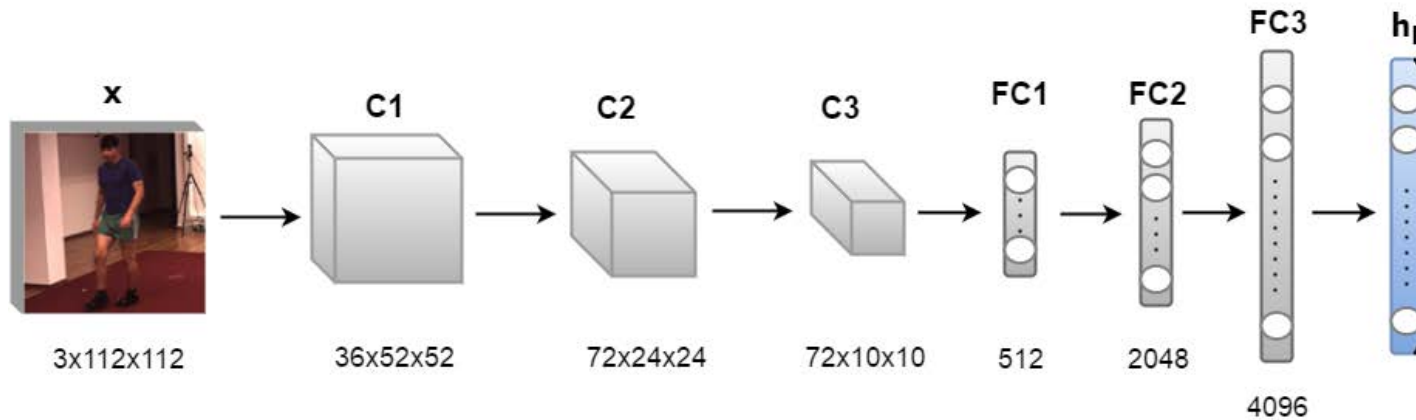
Input: \mathbf{I}

Output: $\{\mathbf{y}_j\}_{1 \leq j \leq J}$

Autoencoder:



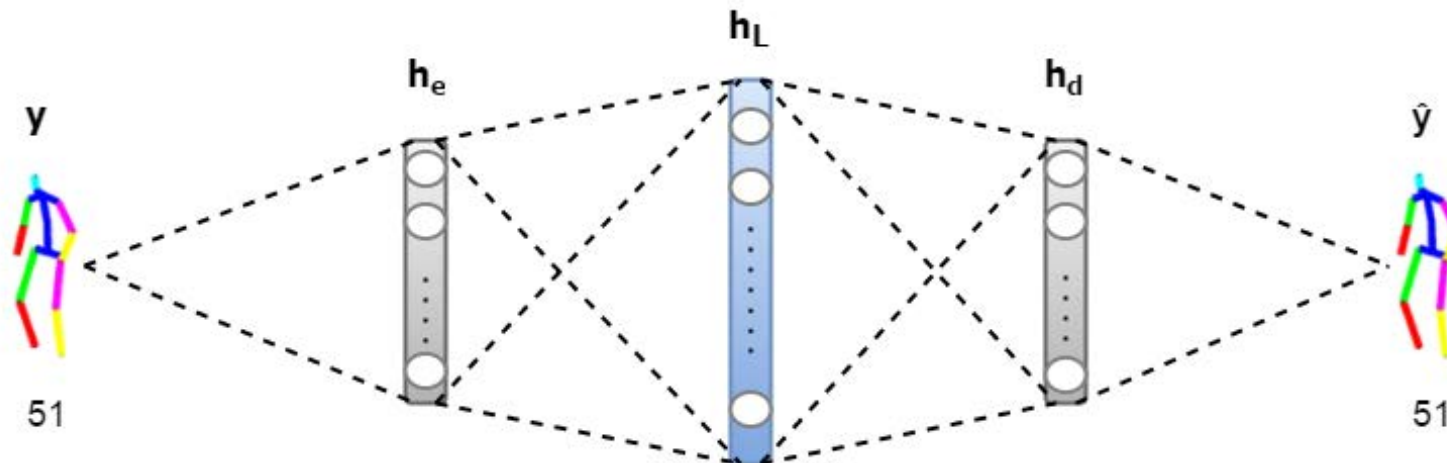
Optional: Structured Prediction



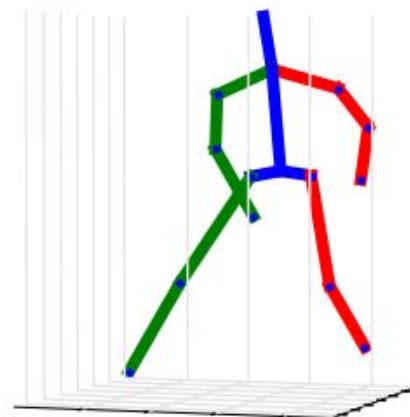
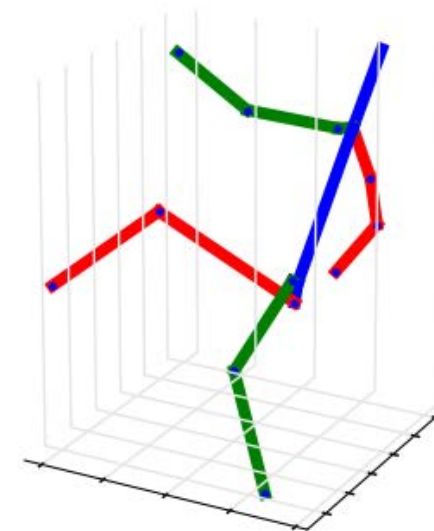
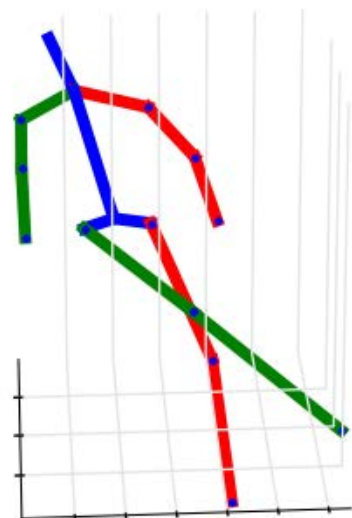
Input: \mathbf{I}

Output: $\{\mathbf{y}_j\}_{1 \leq j \leq J}$

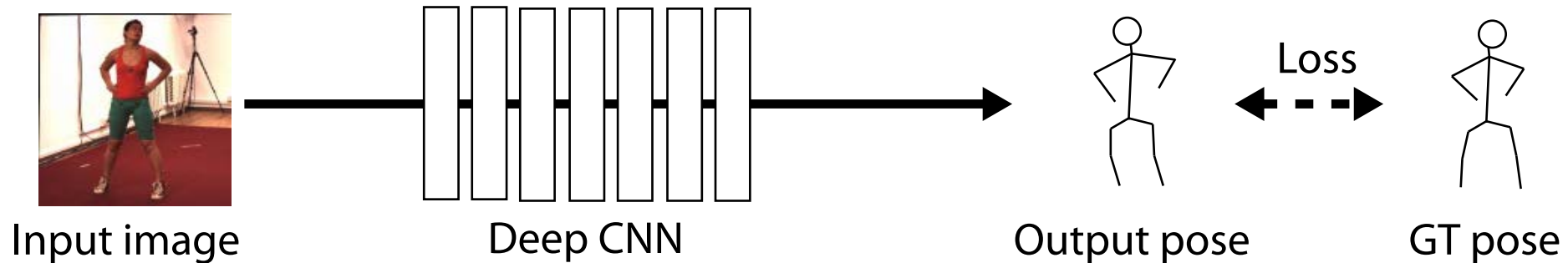
Autoencoder:



Optional: Monocular 3D Pose Estimation



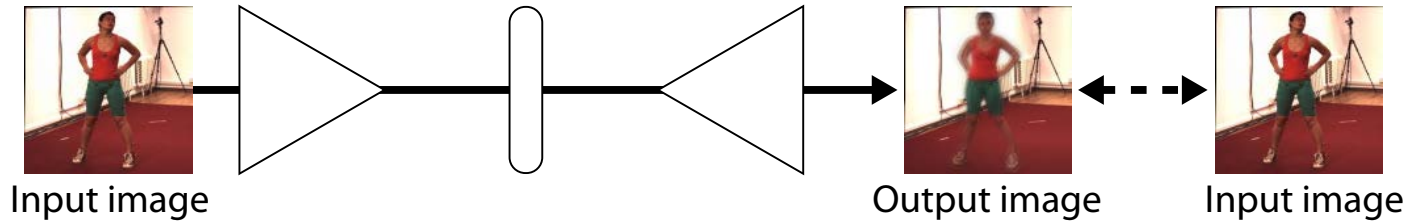
Optional: The Problem with Direct Estimation



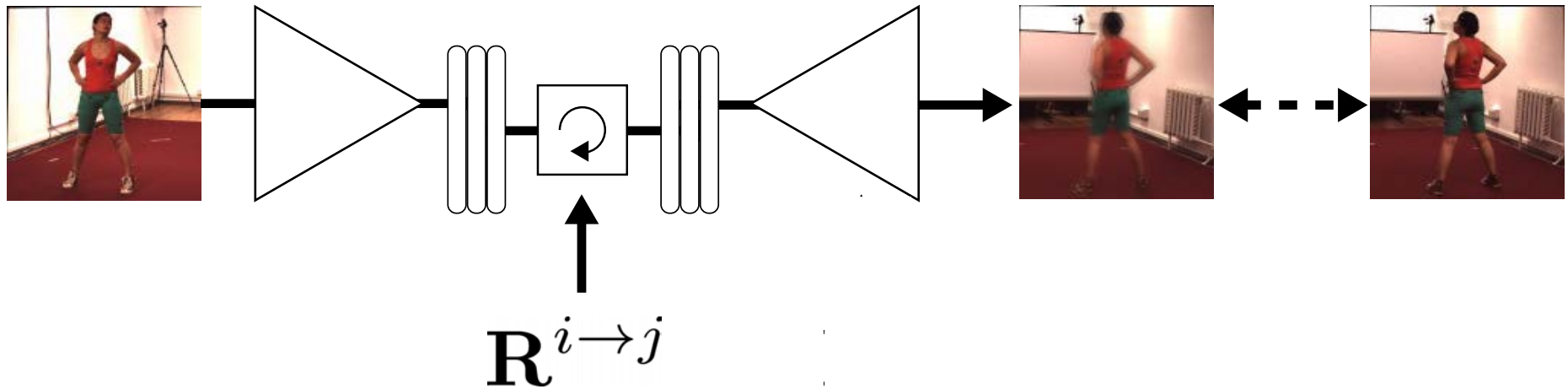
- The human body has many degrees of freedom.
- Going directly from image to 3D pose requires a very deep net.
- Training such a deep net requires a lot of training data.

Can we learn a representation that has fewer degrees of freedom for specific activities?

Optional: Novel View Synthesis Revisited

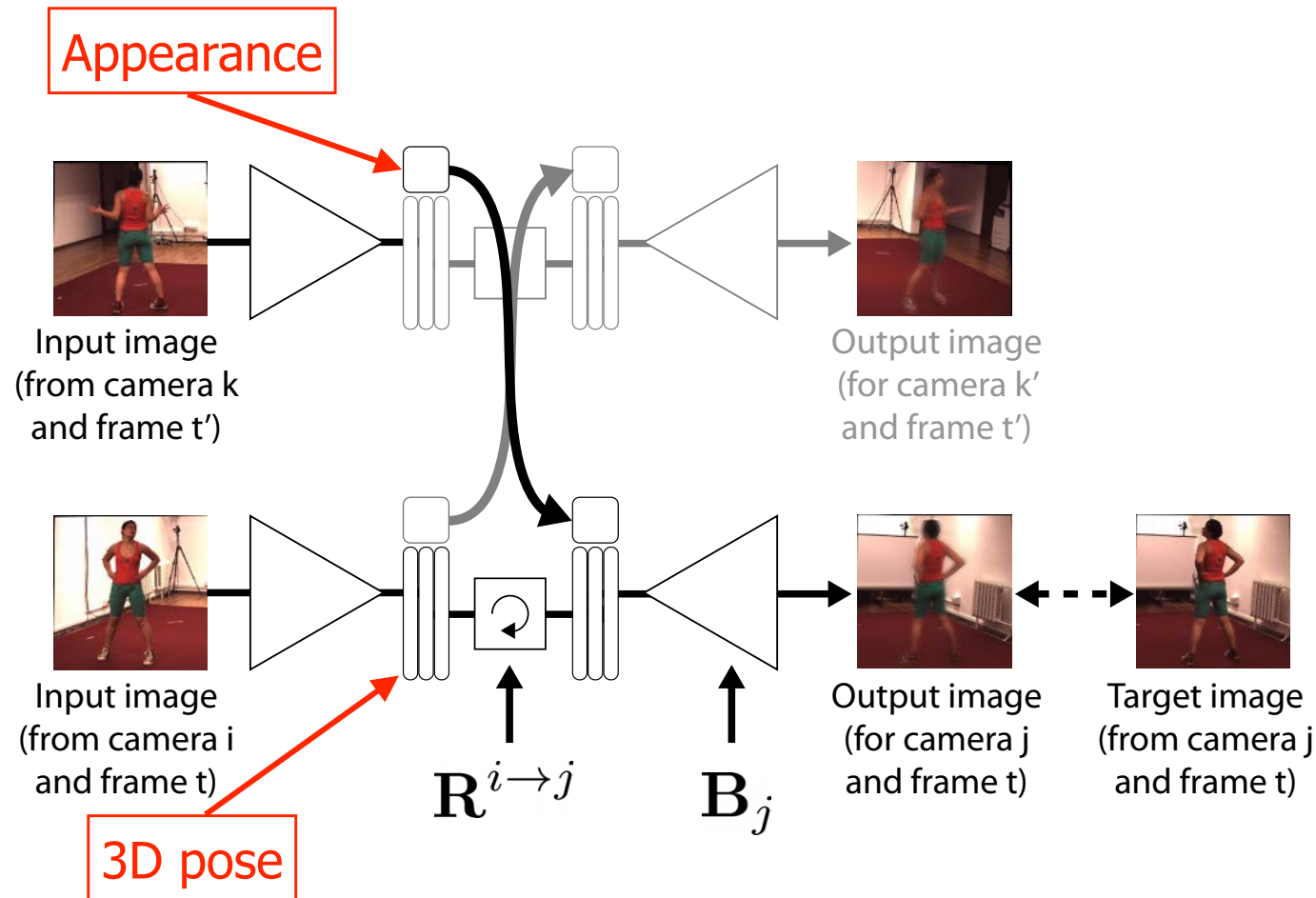


Conventional Autoencoder



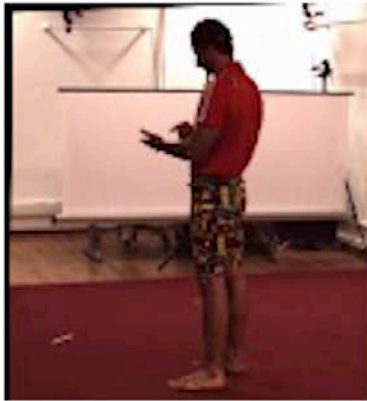
Rotation Aware Autoencoder

Optional: Separating Appearance from Geometry



- The latent representation comprises a $N \times 3$ matrix that encodes the 3D pose and a separate vector that models appearance.
- Before decoding the appearance vectors are swapped to ensure that they are similar in different images.

Optional: Latent Representation



Input



Geometric encoding
(rotating point cloud)



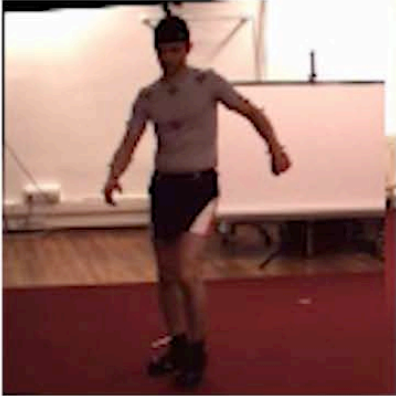
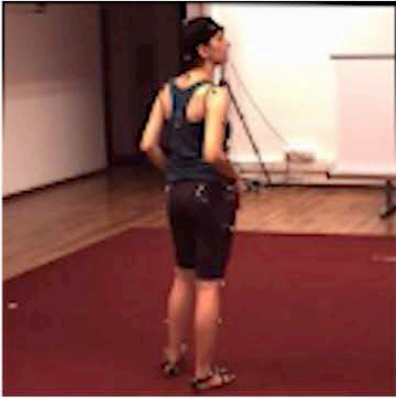
Output

The test subject is reconstructed with the right pose but an approximate appearance.

Optional: Swapping Appearance

Input subject j

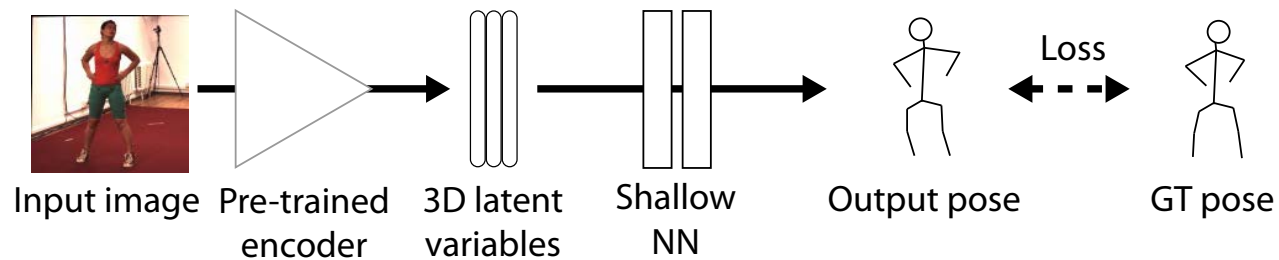
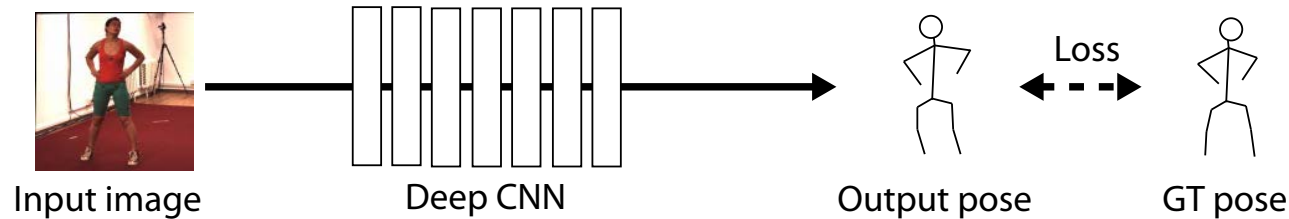
Synthesized pose j
with appearance j



Input subject g

Synthesized pose j
with appearance g

Optional: Direct Estimation vs Using Latent Variables



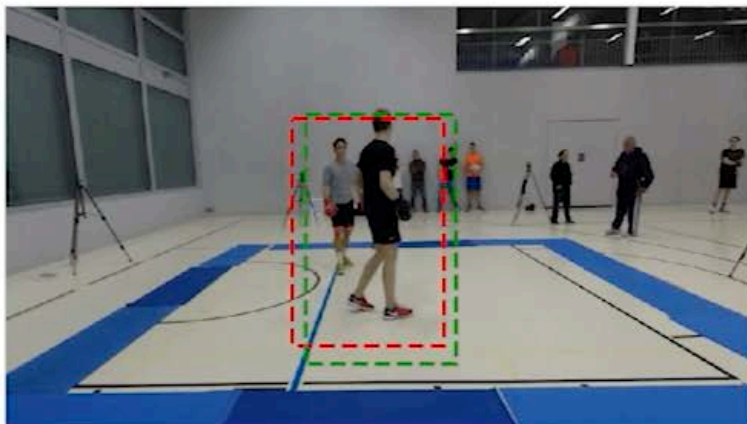
- Training a shallow network requires less training data than a deeper one.
- This is important in specialized areas for which large training databases are hard to build.

Optional: Modeling Boxers

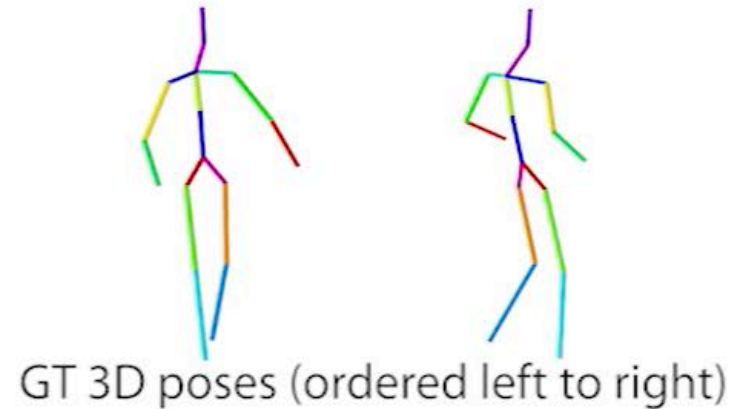
Multi-person 3D pose estimation with NSD



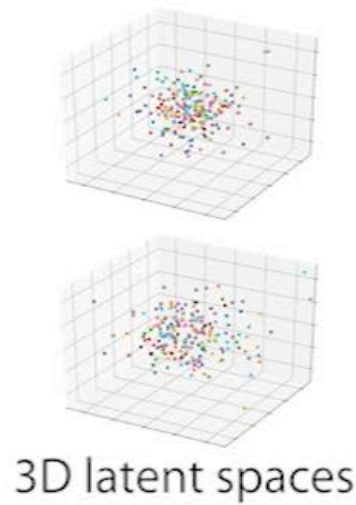
Input



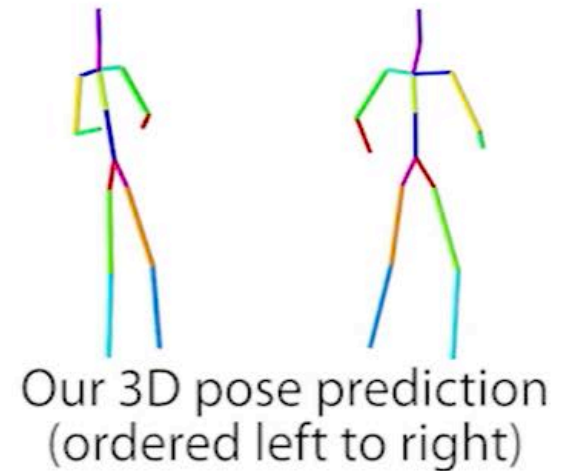
Bounding box



GT 3D poses (ordered left to right)

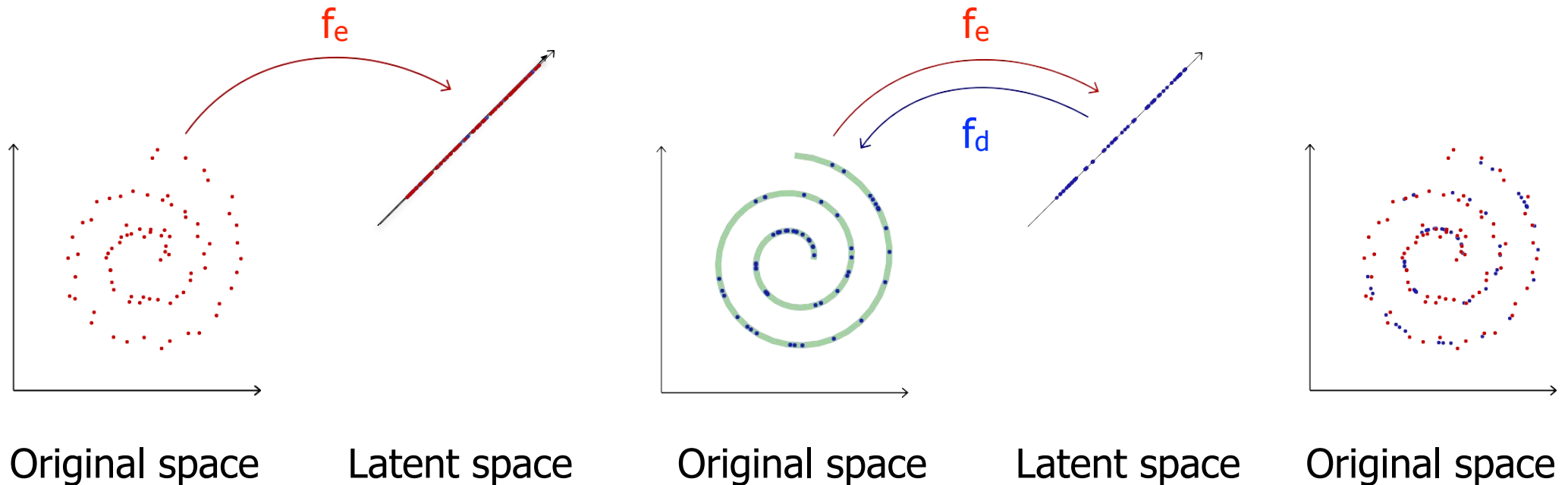


3D latent spaces



Our 3D pose prediction (ordered left to right)

Dimensionality Reduction



$$\mathbf{z} = f_e(\mathbf{x})$$

$$\hat{\mathbf{x}} = f_d(\mathbf{z})$$

$$\hat{\mathbf{x}} \approx \mathbf{x}$$

- Removes unnecessary degrees of freedom.
- Models correlations between the real ones.
- Makes it possible to denoise the original data.
- Can done linearly or non-linearly.