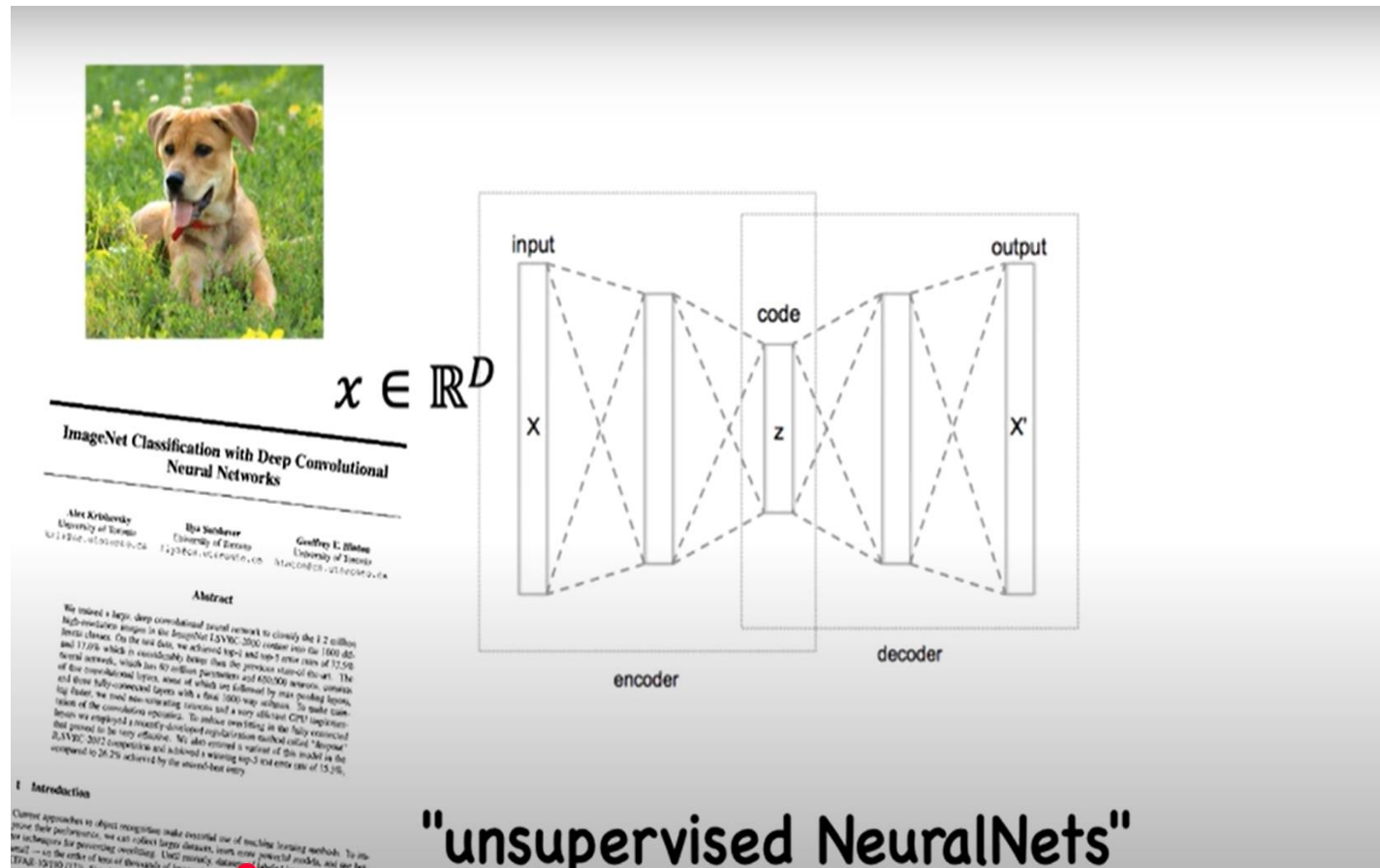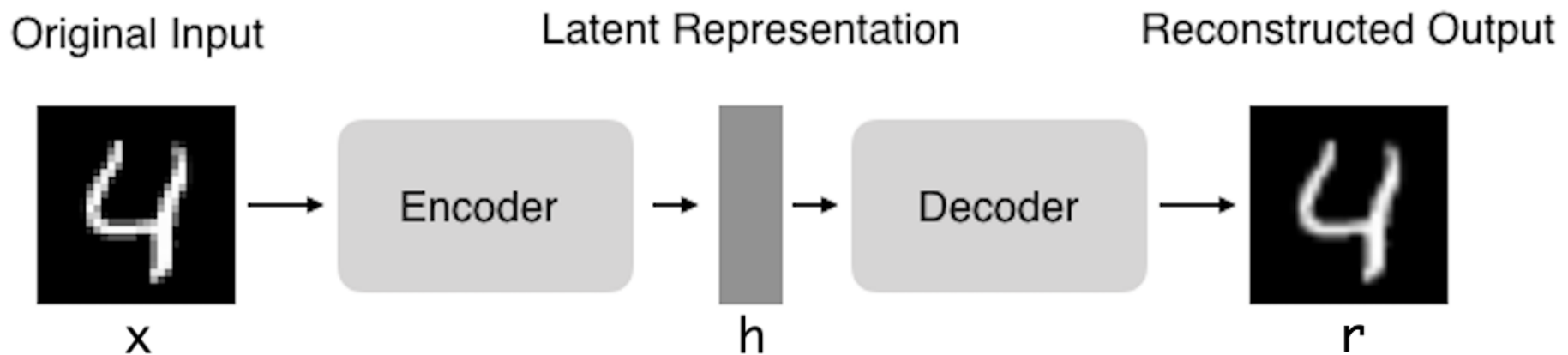# Autoencoders

# Autoencoders

- Autoencoders are designed to reproduce their input, especially for images.
  - Key point is to reproduce the input from a learned encoding.

# Autoencoders: structure

- Encoder:  compress input into a latent-space of usually smaller dimension.  h = f(x)

- Decoder: reconstruct input from the latent space.   r = g(f(x)) with r as close to x as possible

Original Input          Latent Representation          Reconstructed Output

Encoder → h → Decoder

x                              h                              r

# Autoencoders



$x \in \mathbb{R}^D$

$r \in \mathbb{R}^D$

Input layer

Hidden layer

Output layer

"bottleneck"

$h = f(x)$

$r = g(h)$

$\min Loss(x, r)$

# Autoencoders

**Encoding**

$$Z = W1 * X + b1$$

where X is the input data, W1 is the encoding weight matrix, b1 is the encoding bias vector, and Z is the code.

**Decoding:**

$$X' = W2 * Z + b2$$

where W2 is the decoding weight matrix, b2 is the decoding bias vector, and X' is the reconstructed data.

The goal of the linear autoencoder is to minimize the reconstruction error between X and X'.

One common way to do this is to use the mean squared error (MSE) loss function

$$loss = 1/n * \sum(X - X')^2$$

# Under complete autoencoder

- The under complete autoencoder is forced to learn a bottleneck or compressed representation that cannot perfectly reconstruct the original input.

- The idea behind under complete autoencoders is to learn a more efficient and concise representation of the input data by capturing the most important and relevant information.

- The size of the code is smaller than the size of the input,

- Restricting the size of the code, the network is encouraged to learn a compressed representation that captures the most important features of the data and discards irrelevant information.

- undercomplete autoencoder can be used to reduce the dimensionality of the data, which can be useful for visualization, feature extraction, or as input to other machine learning models.

- undercomplete autoencoders can be difficult to train, especially if the code is too small or the data is high-dimensional.

- Regularization techniques such as L1 or L2 regularization, dropout, or sparsity constraints can be used to prevent overfitting and encourage the network to learn more informative representations.



"bottleneck"

$x_1$ $x_2$ $x_3$ $x_4$ $x_5$ $x_6$

$a_1$ $a_2$ $a_3$

$\hat{x}_1$ $\hat{x}_2$ $\hat{x}_3$ $\hat{x}_4$ $\hat{x}_5$ $\hat{x}_6$

$$n_h < d$$

"undercomplete"

Scroll for details

# Sparse Autoencoder

**Sparse Autoencoder** – Uses sparsity constraints to force only a few neurons to be active, improving feature learning.

# Overcomplete encoder

- An overcomplete autoencoder is a type of autoencoder that is trained to learn a compressed representation of the input data, where the size of the code is larger than the size of the input. In other words, the overcomplete autoencoder is capable of representing the input data in a more expressive way than an undercomplete autoencoder.

- overcomplete autoencoders is to learn a more rich and complex representation of the input data by allowing the code to have more dimensions than the input.

- Increasing the number of neurons in the hidden layer or layers, the network can learn to capture more complex patterns and features in the data.

- Overcomplete autoencoders can be useful for tasks such as denoising, where the goal is to remove noise or corruption from the input data.

- By learning a more expressive representation of the input, the overcomplete autoencoder can better differentiate between noise and signal, leading to improved denoising performance.

# Properties of Autoencoders

- **Data-specific**: Autoencoders are only able to compress data similar to what they have been trained on.

- **Lossy:** The decompressed outputs will be degraded compared to the original inputs.

- **Learned automatically from examples:** It is easy to train specialized instances of the algorithm that will perform well on a specific type of input.

# Denoising autoencoders

**Denoising Autoencoder** – Trained with noisy inputs to learn robust feature representations.

- Basic autoencoder trains to minimize the loss between x and the reconstruction g(f(x)).

- Denoising autoencoders train to minimize the loss between x and g(f(x+w)), where w is random noise.

- Same possible architectures, different training data.

- [Kaggle has a dataset on damaged documents.](#)

# Denoising autoencoders

- Denoising autoencoders are a type of autoencoder that is trained to remove noise from corrupted input data.

- The goal of denoising autoencoders is to learn a representation of the input data that captures the underlying structure and removes the noise.

- The basic architecture of a denoising autoencoder is similar to a regular autoencoder, with an encoder that maps the input data to a compressed representation (code), and a decoder that maps the compressed representation back to the original input.

- Denoising autoencoders have been used in a variety of applications, such as image denoising, speech denoising, and text denoising.

# Applications of Autoencoders

- **Dimensionality Reduction** (alternative to PCA)
- **Anomaly Detection** (e.g., fraud detection, medical diagnosis)
- **Denoising Images**
- **Feature Learning** for downstream ML tasks
- **Generative Modeling** (e.g., VAEs for image synthesis)