

PAT 498/598 (Fall 2024)

Special Topics: Generative AI for Music and Audio Creation

Lecture 11: Midterm Review

Instructor: Hao-Wen Dong



SCHOOL OF MUSIC, THEATRE & DANCE
PERFORMING ARTS TECHNOLOGY
UNIVERSITY OF MICHIGAN

Review – Background

What is Artificial Intelligence?

AI is the study of how to make computers **do things at which, at the moment, people are better.**

– Elaine Rich and Kevin Knight, 1991

1997



(Source: Britannica)

2016



(Source: The Guardian)

20??



(Source: SC2HL)

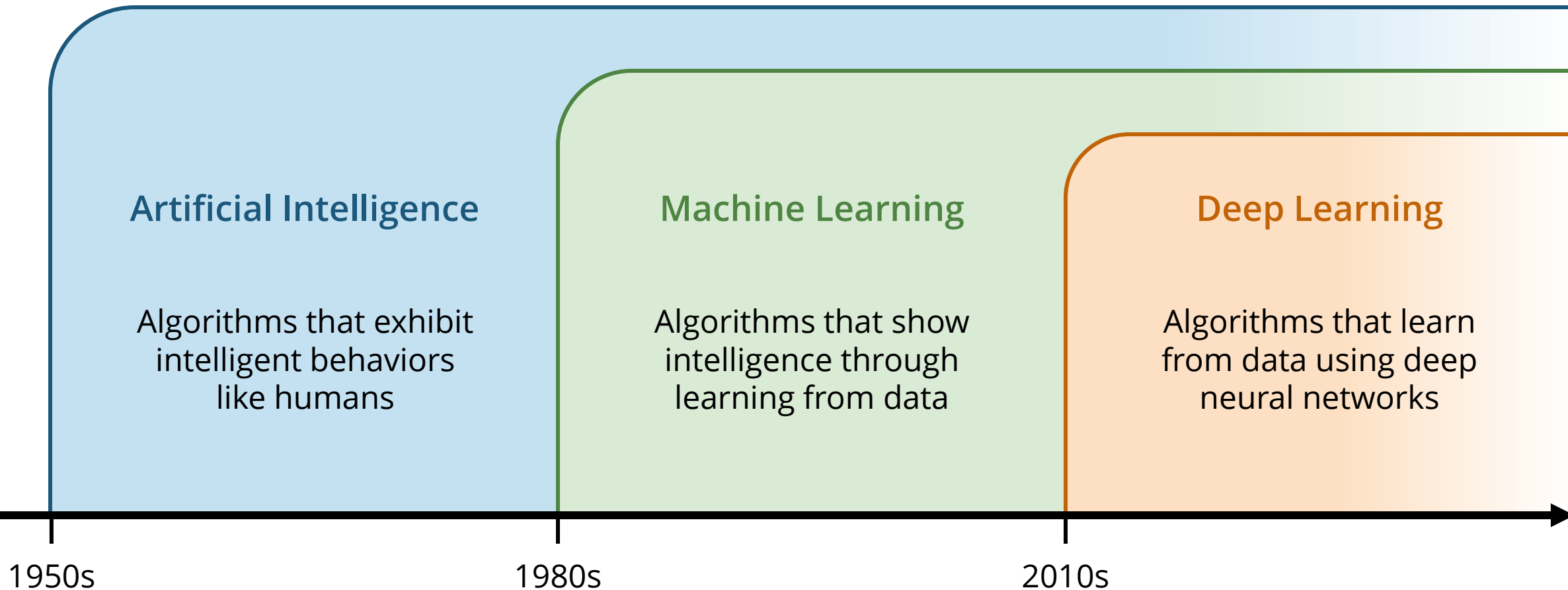
Elaine Rich and Kevin Knight, *Artificial Intelligence*. United Kingdom: McGraw-Hill, 1991.

<https://www.britannica.com/topic/Deep-Blue>

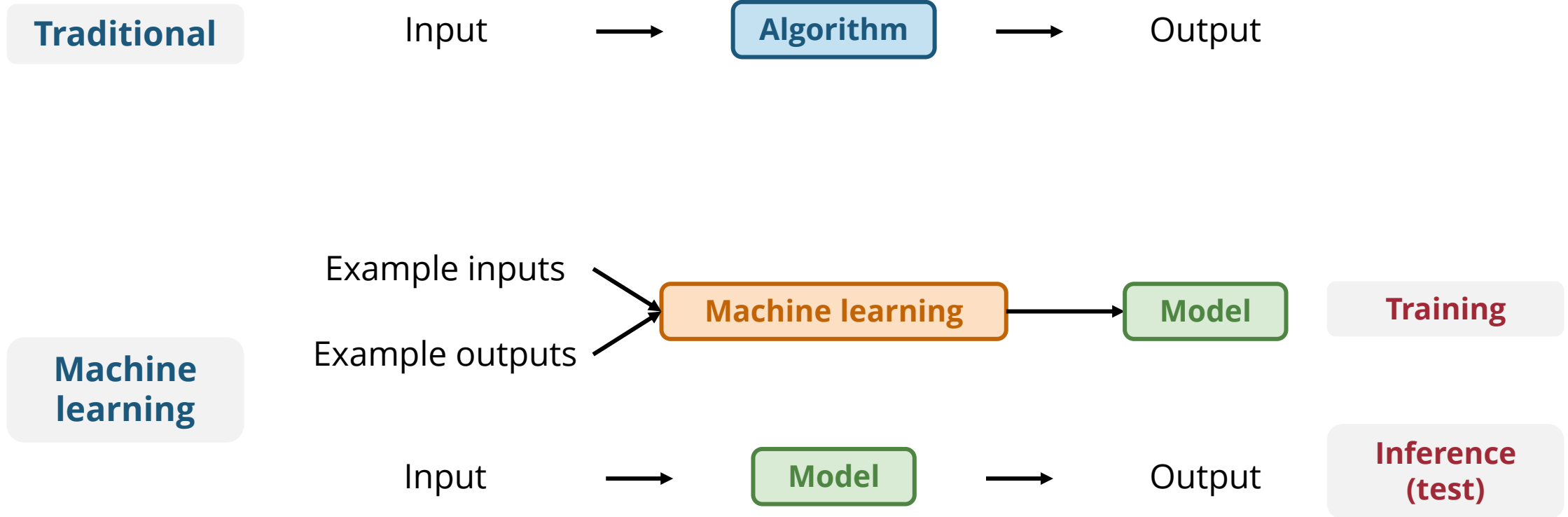
<https://www.theguardian.com/technology/2016/mar/15/alphago-what-does-google-advanced-software-go-next>

https://www.youtube.com/watch?v=PFMRDm_H9Sg

AI vs ML vs DL



Machine Learning



Components of a Machine Learning Model

Optimization

Defining inputs & outputs

Improve on task T,

with respect to performance metric P,

based on experience E

Loss function
(objective function)

Training data

Deep learning is almost the same as machine learning by this definition!

What's special about deep learning?

Types of Machine Learning

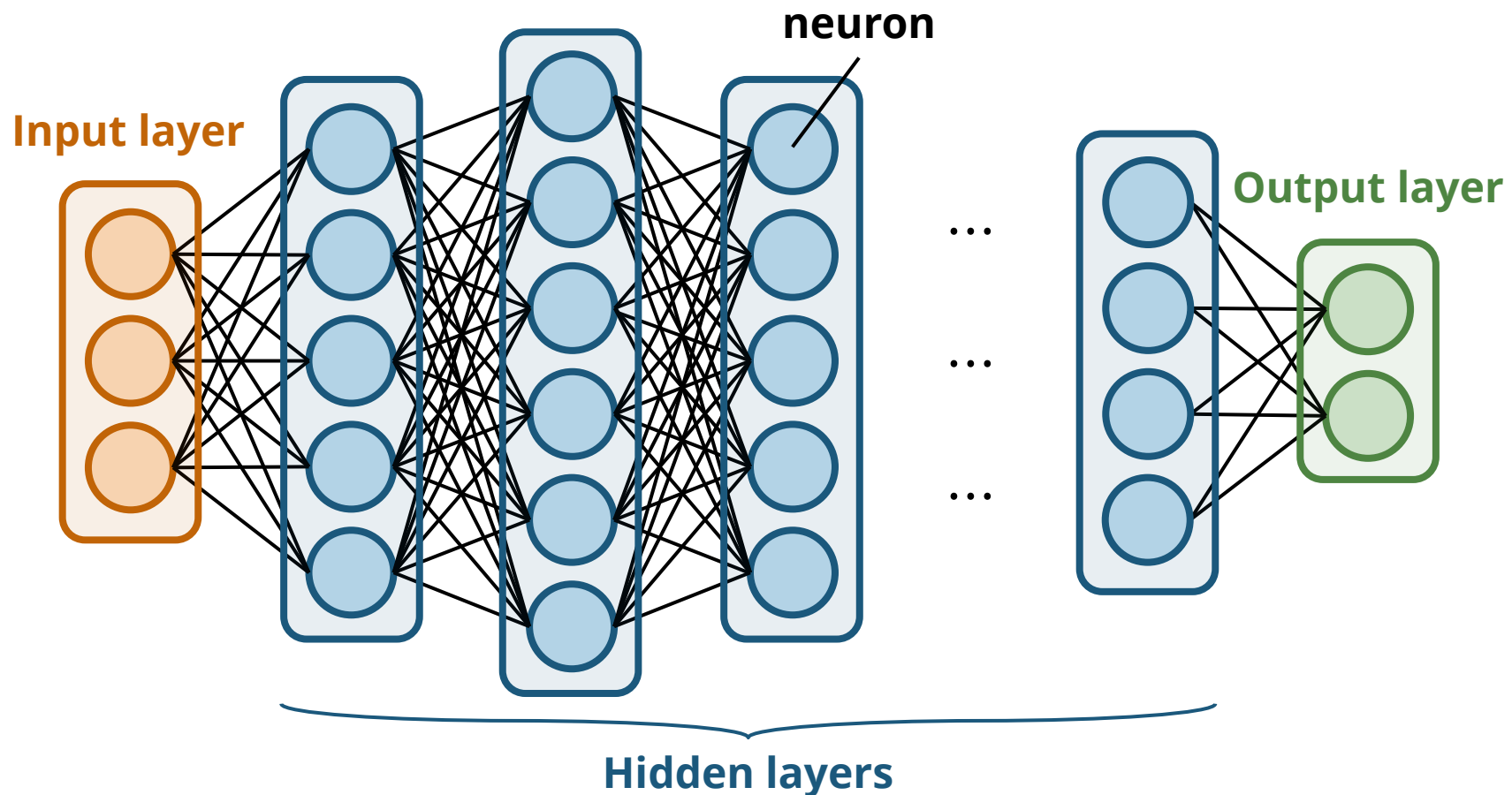
- **Supervised learning**
 - **Classification:** *discrete* outputs
 - **Regression:** *continuous* outputsGiven **pairs of example inputs and outputs**
- **Unsupervised learning**
 - **Self-supervised learning**Given **only example inputs**
- **Semi-supervised learning** Given **example inputs** and **a few example outputs**
- **Reinforcement learning** Given **scalar rewards** for **a sequence of actions**

Many generative AI models based on self-supervised learning!

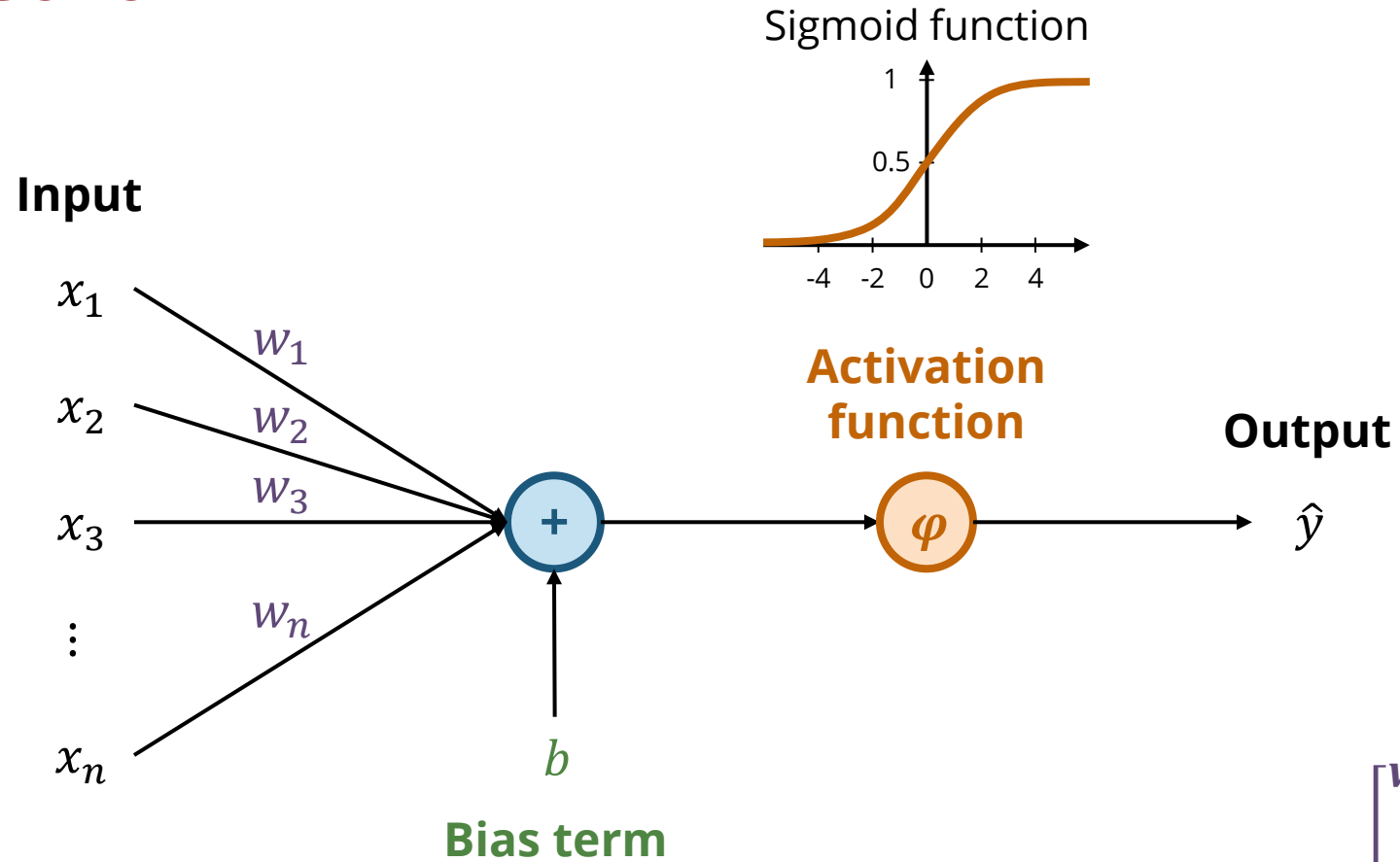
Review – Deep Learning Fundamentals

What is Deep Learning?

- A type of machine learning that uses **deep neural networks**



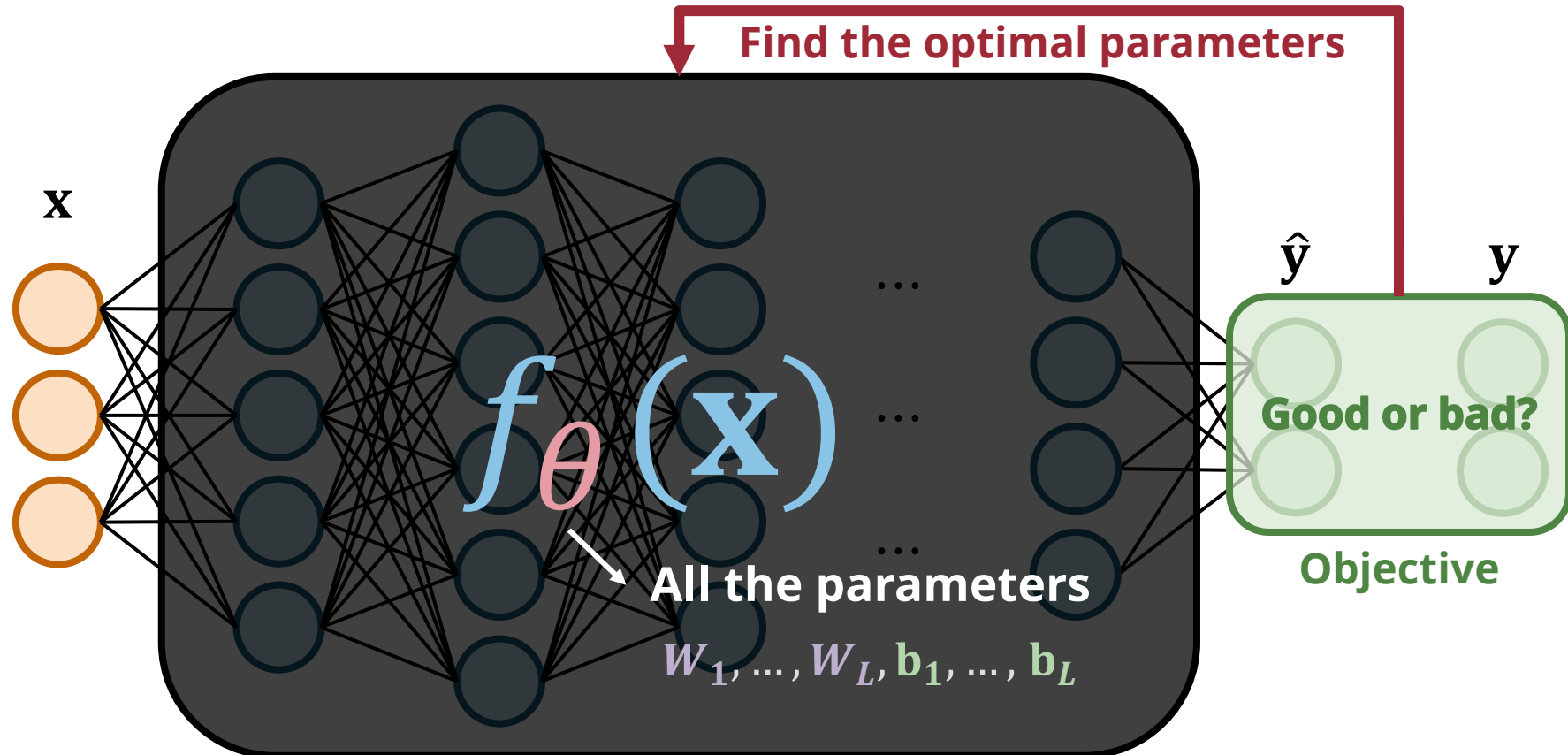
Inside a Neuron



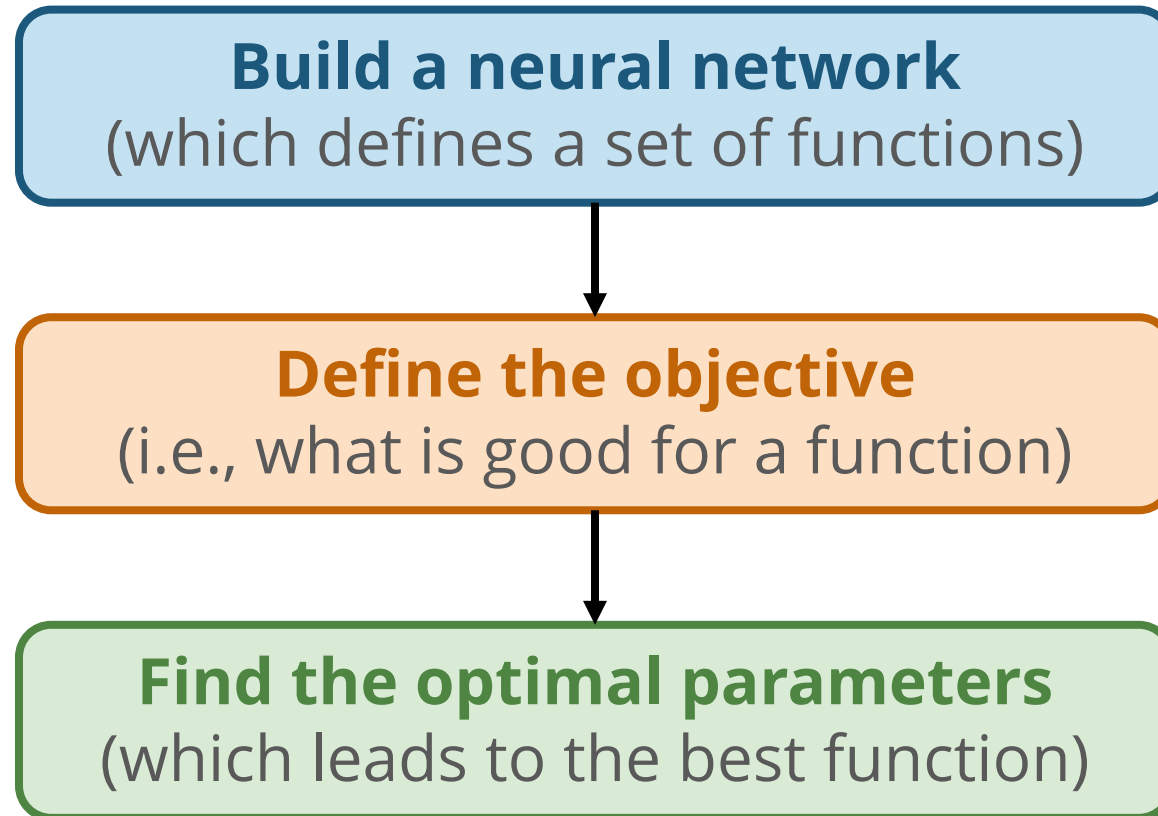
$$\hat{y} = \varphi(w_1x_1 + w_2x_2 + \dots + w_nx_n + b) = \varphi\left(\sum_{i=1}^n w_i x_i + b\right) = \varphi(\mathbf{w} \cdot \mathbf{x} + b)$$

Neural Networks are Parameterized Functions

- A neural network represents **a set of functions**



Training a Neural Network



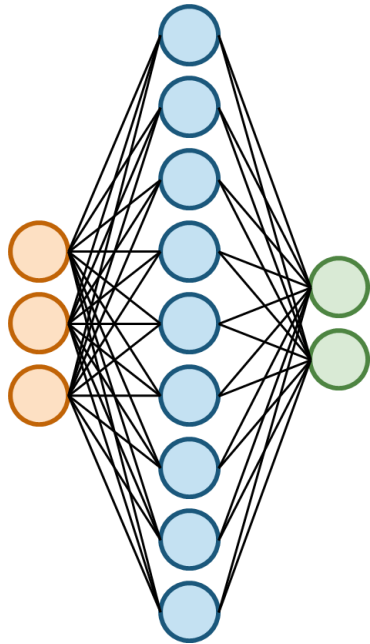
$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$$

$$Loss(\boldsymbol{\theta}) = \sum_k^N L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

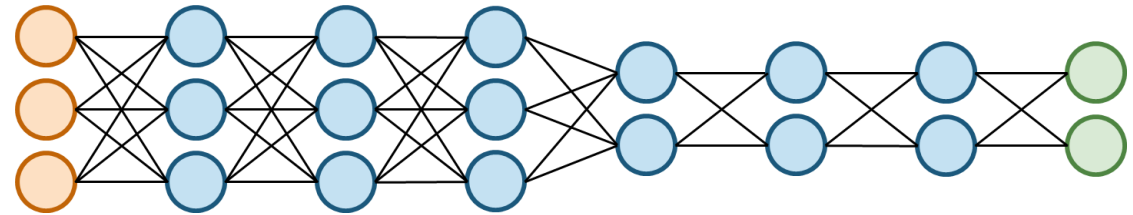
Shallow vs Deep Neural Networks – In Practice

Shallow neural nets



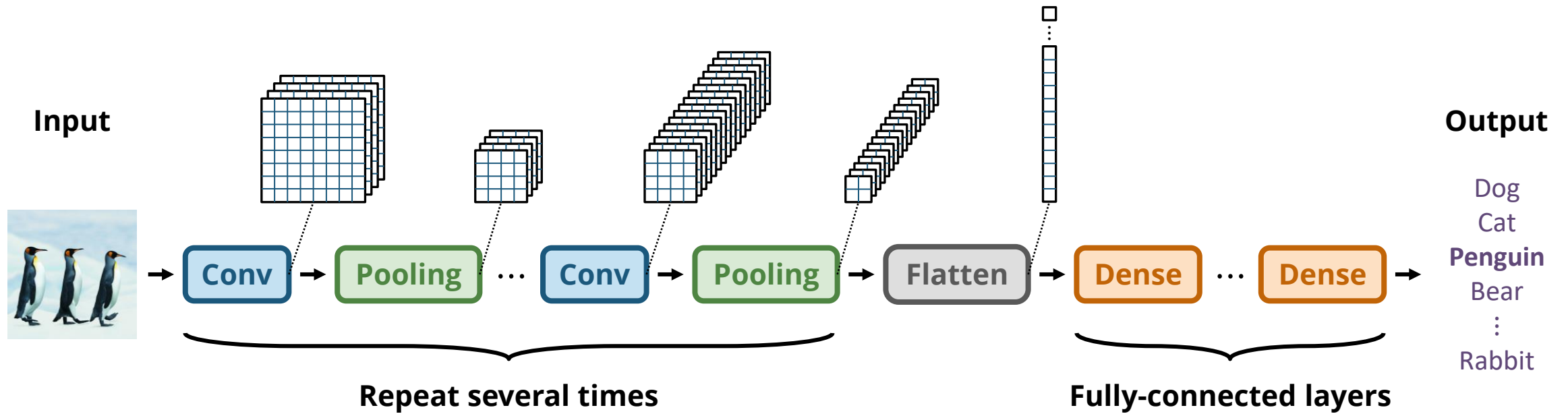
Less expressive
(less parameter efficient)

Deep neural nets



More expressive
(more parameter efficient)

Convolutional Neural Network (CNNs)



2D Convolution

Input

1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

*

Kernel

1	-1	-1
-1	1	-1
-1	-1	1

=

Output

9	
	9

High activation when the local pattern is close to the kernel

2D Convolution

Input

1	-1	-1	-1
-1	1	-1	-1
-1	-1	1	-1
-1	-1	-1	1

Kernel

1	-1	-1
-1	1	-1
-1	-1	1

*

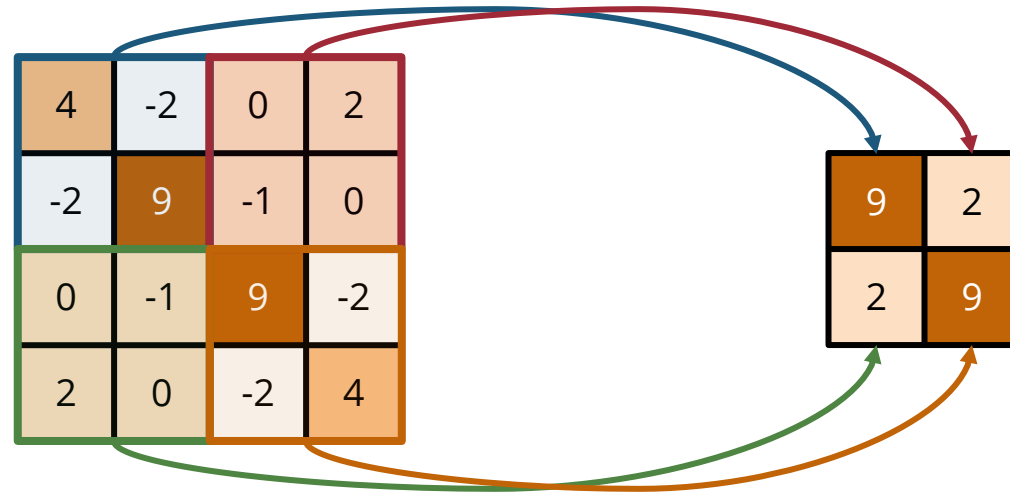
=

Output

	-1
-1	

**Low activation when
the local pattern
differs from the kernel**

Max Pooling Layer

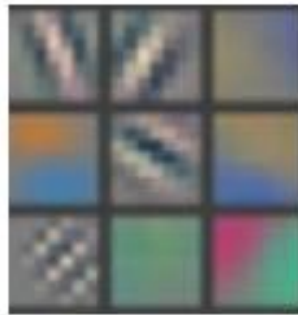


Downsample and keep the strongest activation in each block

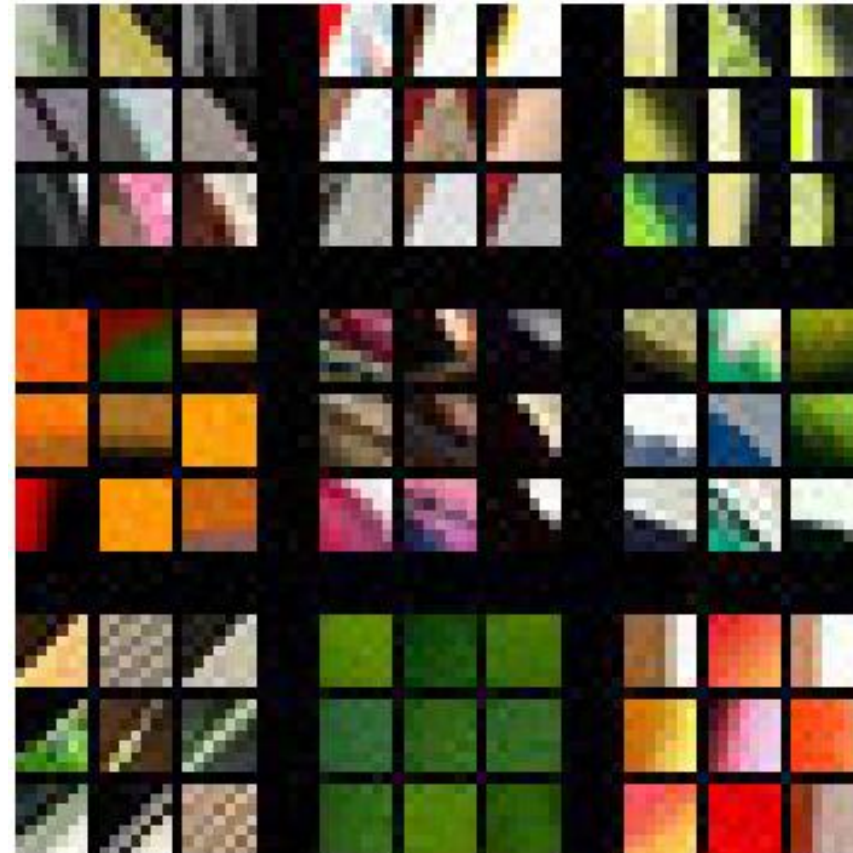
Learned CNN Kernels in a Trained AlexNet

Layer 1

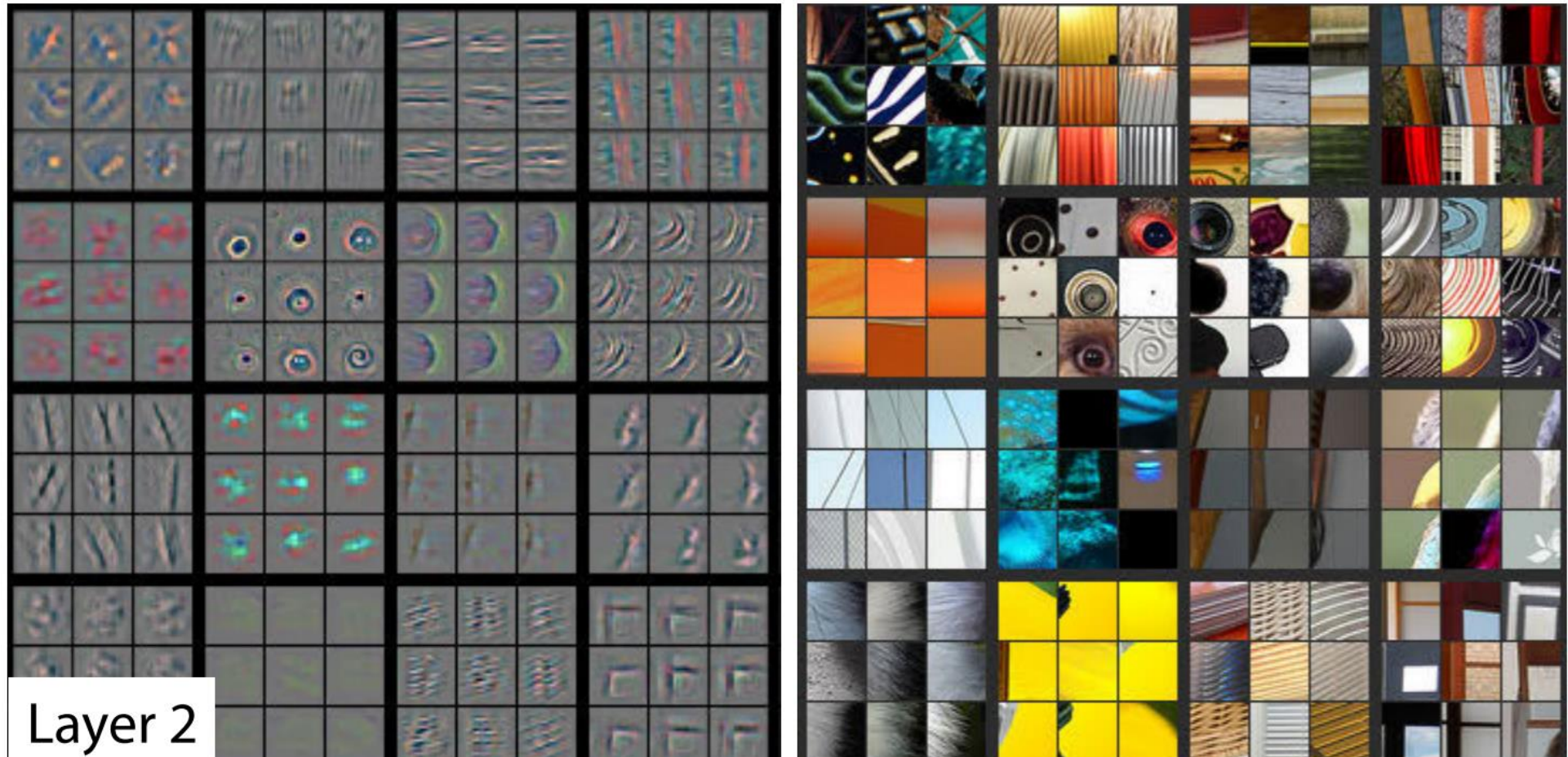
Learned CNN kernels



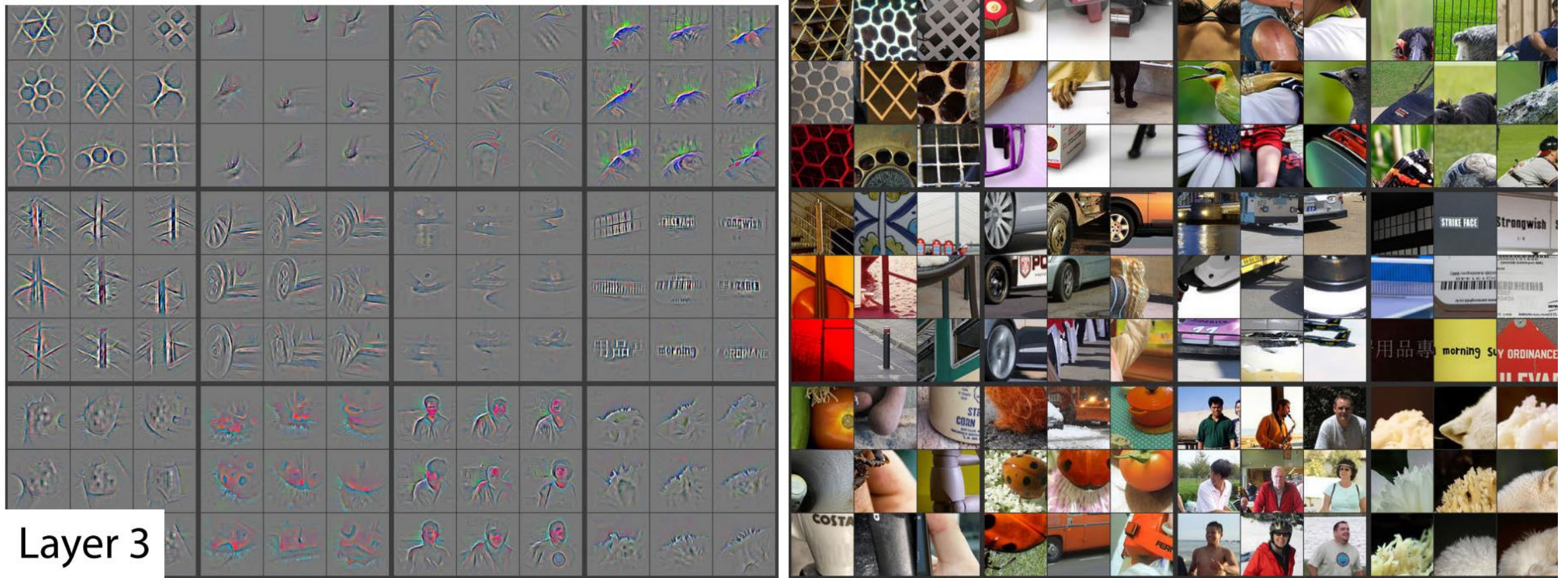
Top activations



Learned CNN Kernels in a Trained AlexNet

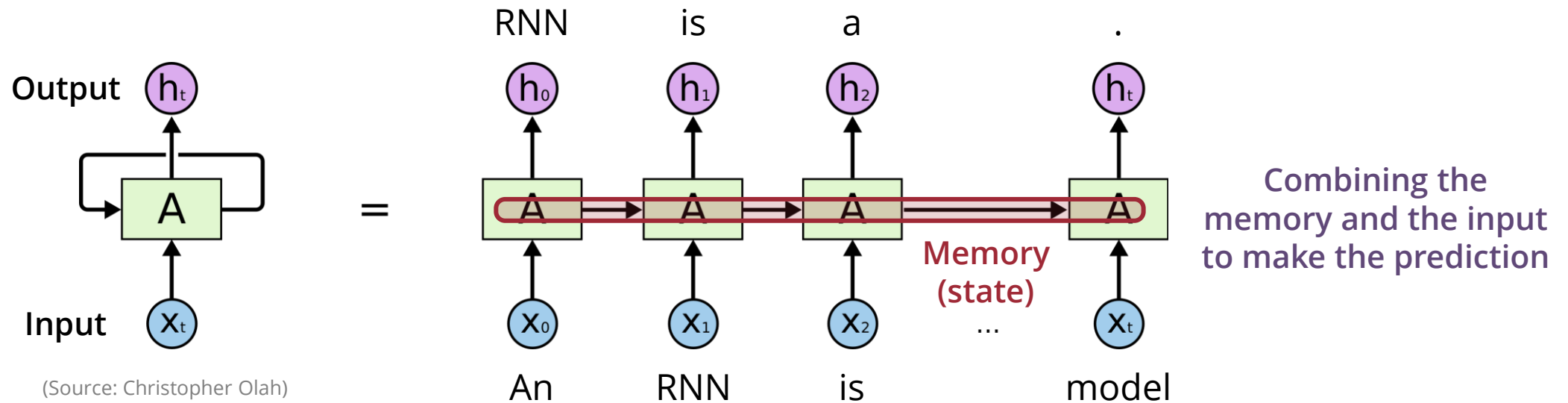


Learned CNN Kernels in a Trained AlexNet

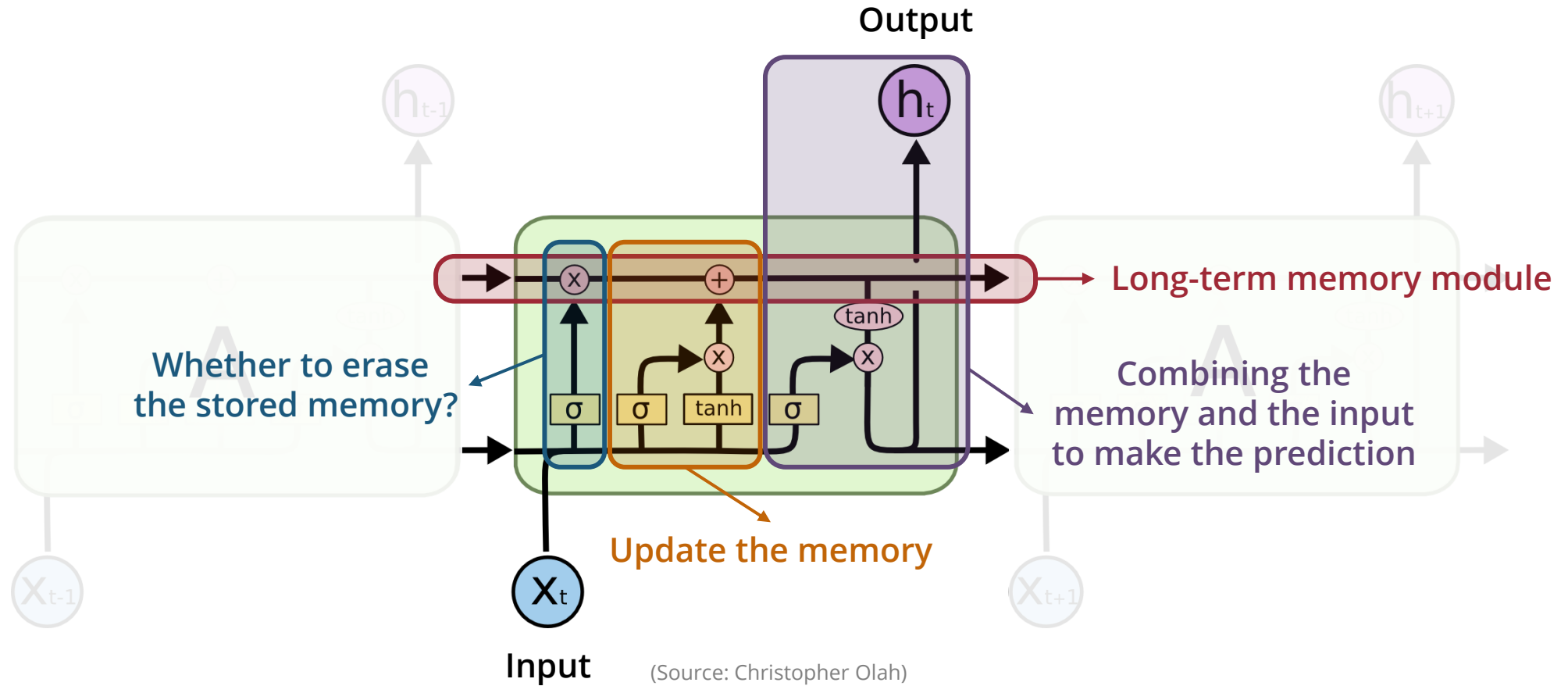


What is an RNN (Recurrent Neural Network)?

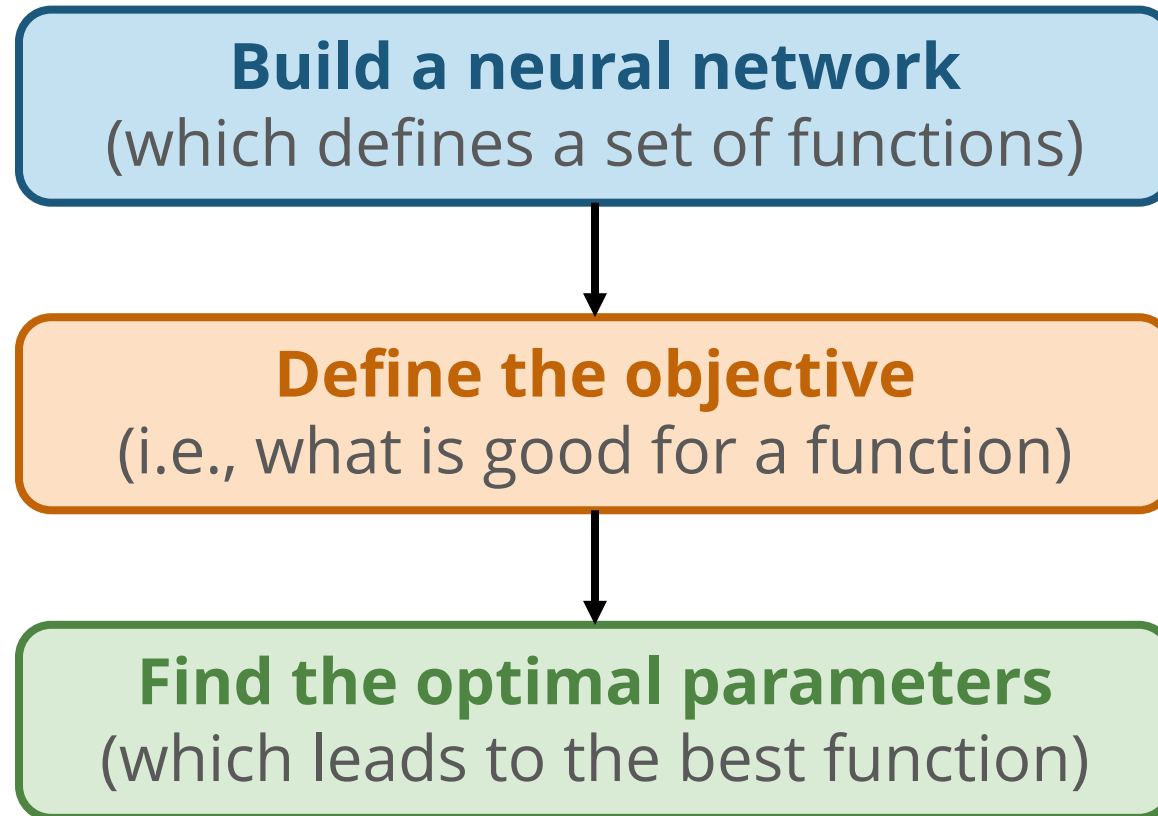
- A type of neural networks that have **loops**
- Widely used for **modeling sequences** (e.g., in natural language processing)



Demystifying LSTMs



Training a Neural Network

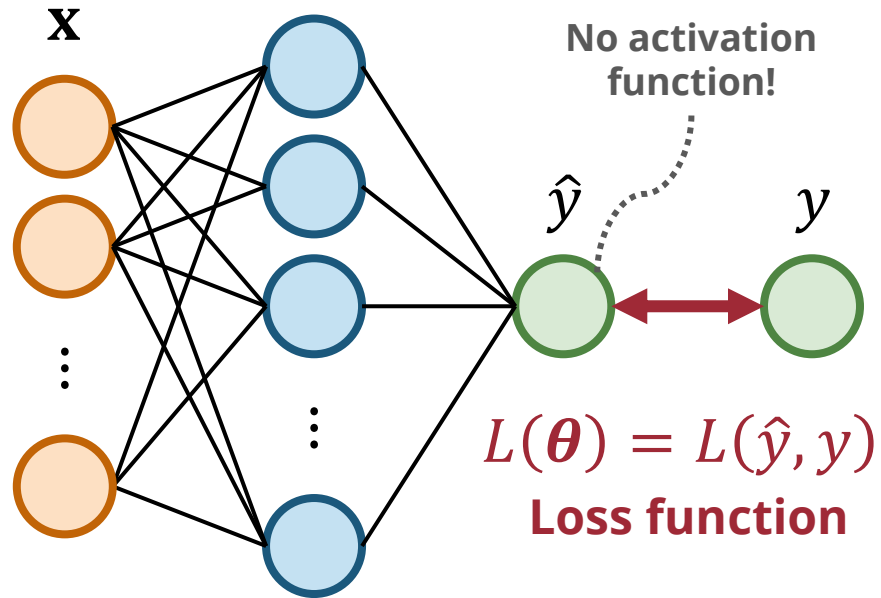


$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$$

$$Loss(\boldsymbol{\theta}) = \sum_k^N L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

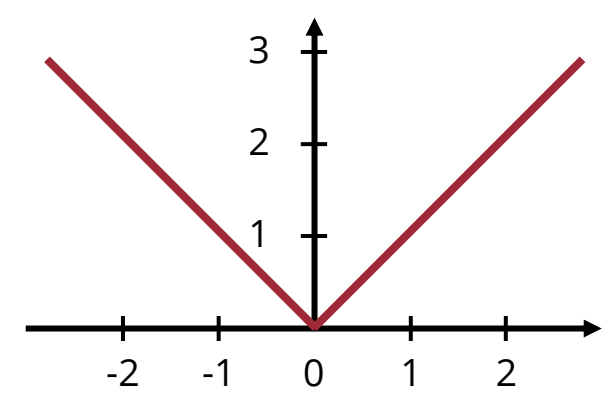
Common Loss Functions for Regression



$$Loss(\theta) = \sum_k^N L(\hat{y}_k, y_k)$$

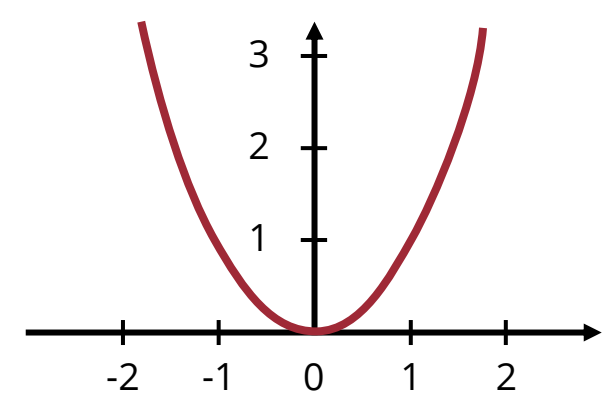
L1 loss

$$L(\hat{y}, y) = |\hat{y} - y|$$



L2 loss

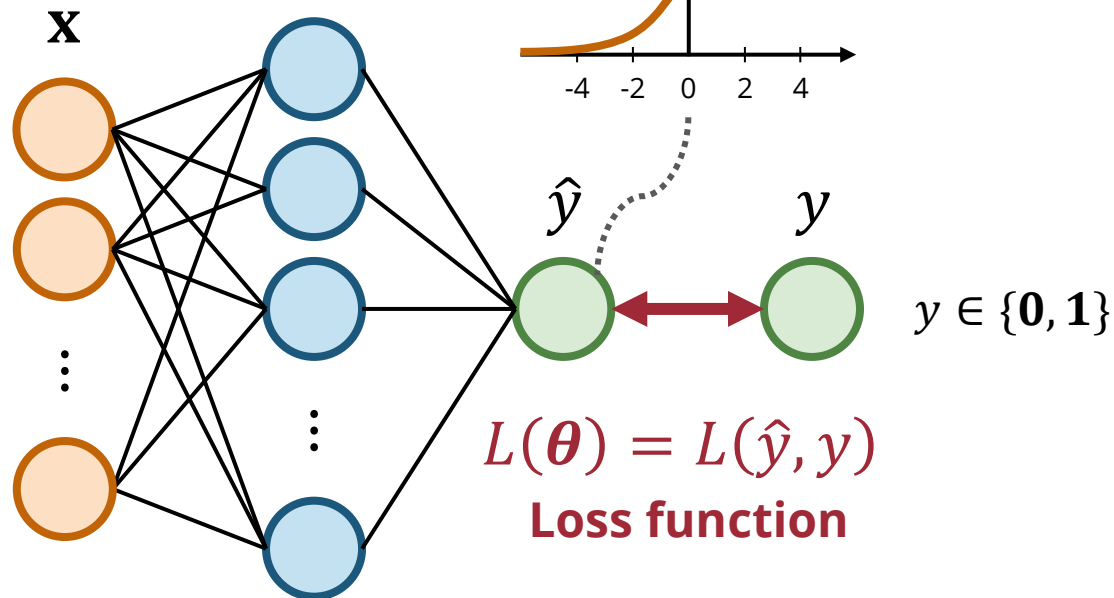
$$L(\hat{y}, y) = (\hat{y} - y)^2$$



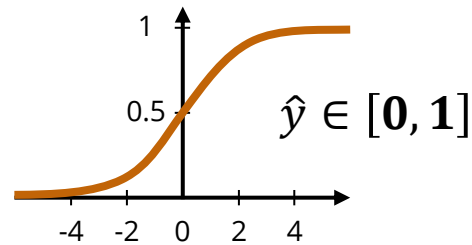
Binary Cross Entropy for Binary Classification

- **Logistic regression** approaches classification like regression

$$Loss(\theta) = \sum_k^N L(\hat{y}_k, y_k)$$



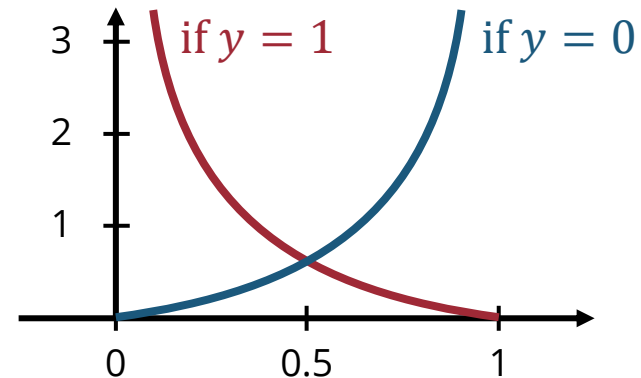
Sigmoid function



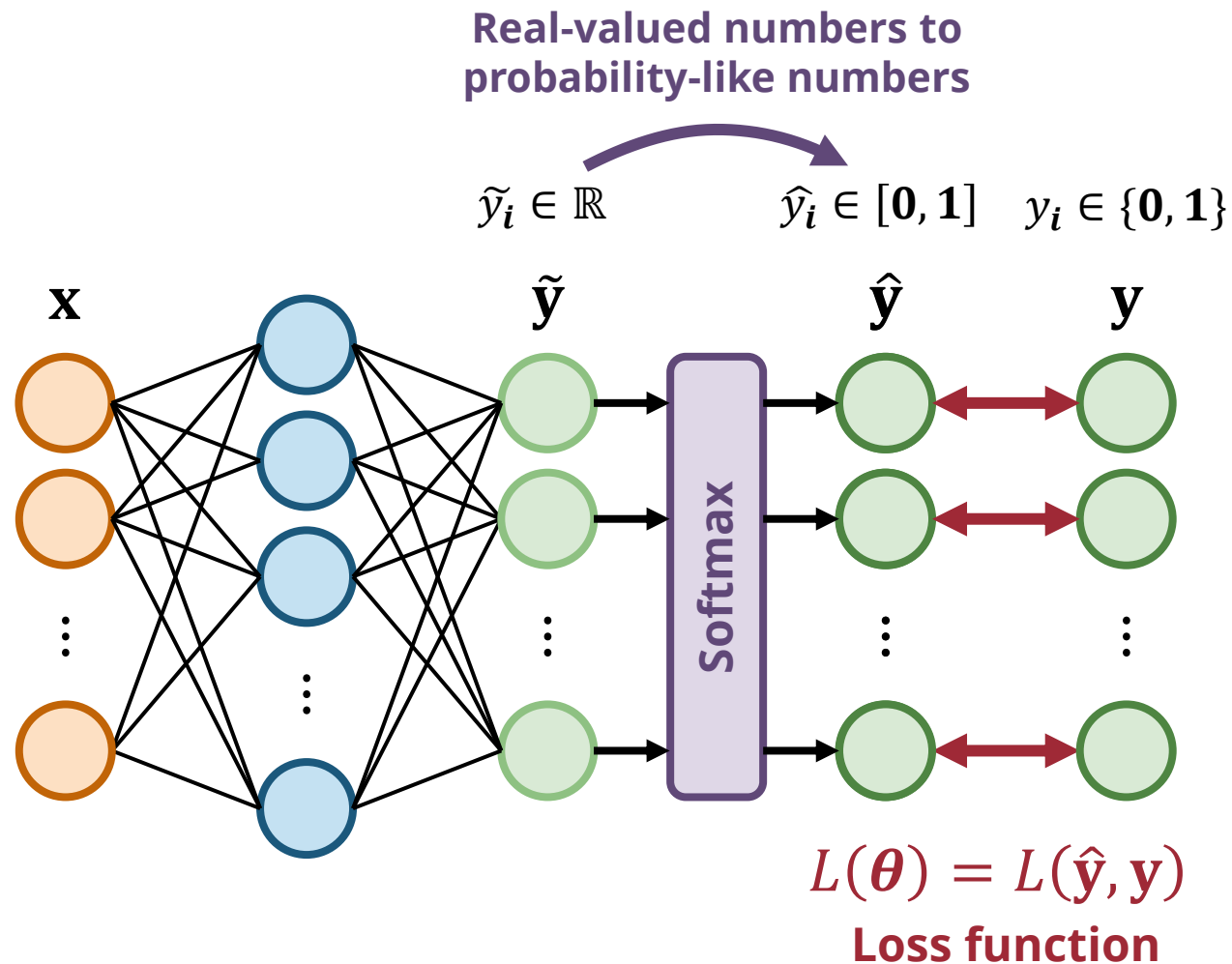
Binary cross entropy

(Also called log loss)

$$L(\hat{y}, y) = \begin{cases} -\log \hat{y}, & \text{if } y = 1 \\ -\log(1 - \hat{y}), & \text{if } y = 0 \end{cases}$$
$$= -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$



Cross Entropy for Multiclass Classification



Softmax

$$\hat{y}_i = \frac{e^{\tilde{y}_i}}{\sum_{j=1}^n e^{\tilde{y}_j}}$$

Cross entropy

$$L(\hat{\mathbf{y}}, \mathbf{y}) = - \sum_i^n y_i \log \hat{y}_i$$

$$Loss(\boldsymbol{\theta}) = \sum_k^N L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

Cross Entropy for Multiclass Classification

Binary Cross Entropy

Only one of them will be one!

$$L(\hat{y}, y) = -y \log \hat{y} - (1 - y) \log(1 - \hat{y})$$

Cross Entropy

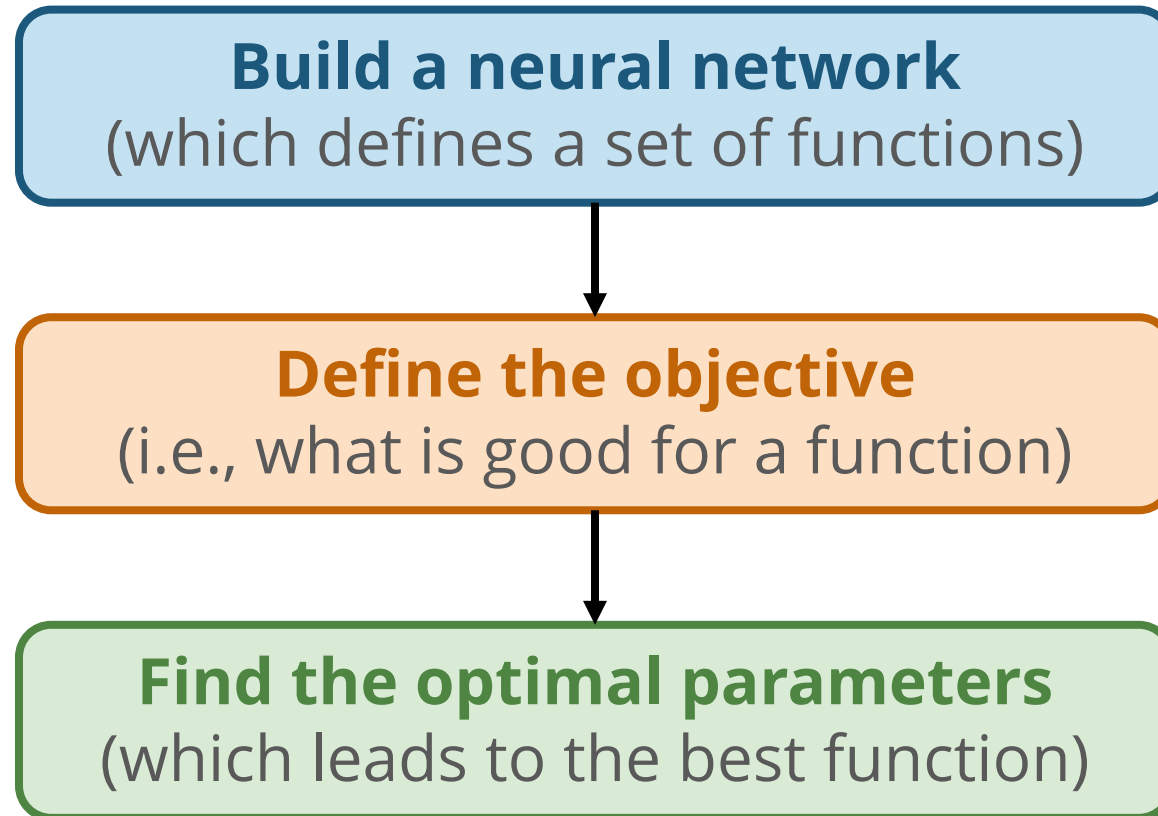
Only one of them will be one!

$$L(\hat{\mathbf{y}}, \mathbf{y}) = -y_1 \log \hat{y}_1 - y_2 \log \hat{y}_2 - \dots - y_i \log \hat{y}_n$$

$$= -\sum_i^n y_i \log \hat{y}_i$$

Log likelihood

Training a Neural Network



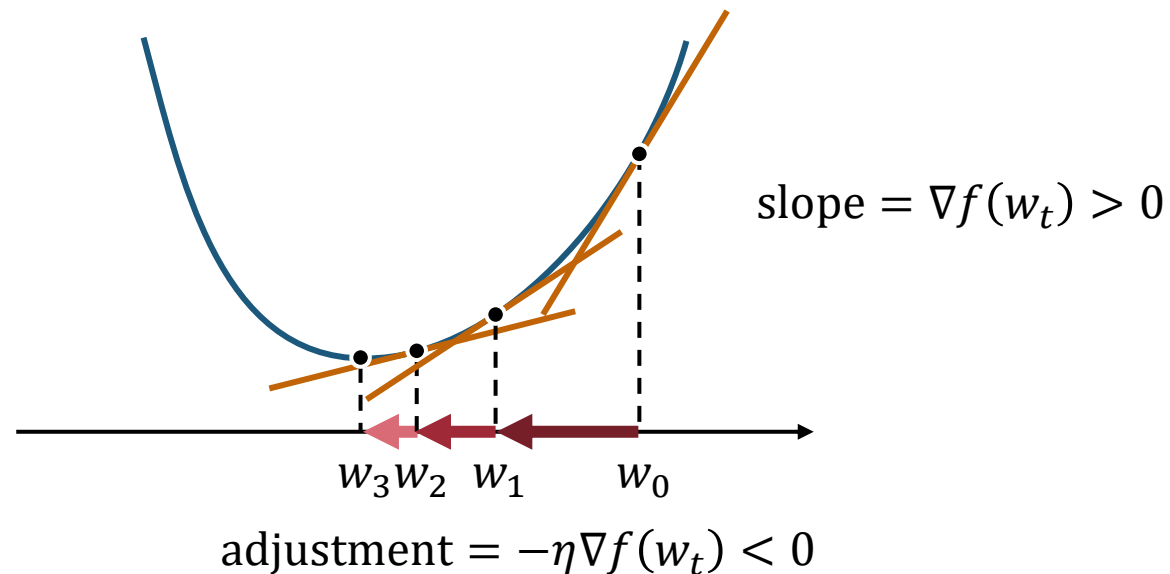
$$\hat{\mathbf{y}} = f_{\boldsymbol{\theta}}(\mathbf{x})$$

$$Loss(\boldsymbol{\theta}) = \sum_k^N L(\hat{\mathbf{y}}_k, \mathbf{y}_k)$$

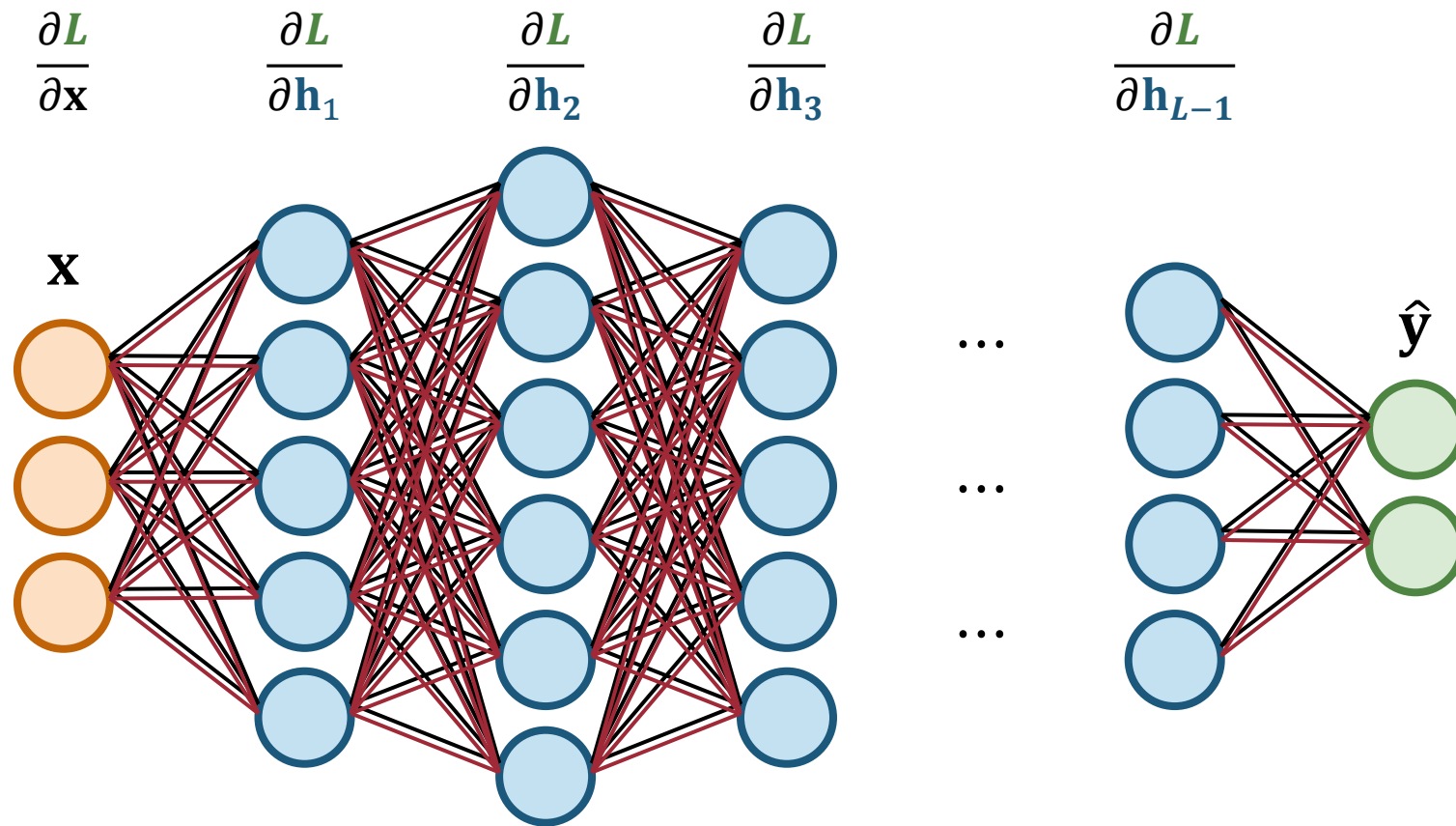
$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} L(\boldsymbol{\theta})$$

Gradient Descent – Pseudocode

- Pick an initial weight vector w_0 and learning rate η
- Repeat until convergence: $w_{t+1} = w_t - \eta \nabla f(w_t)$



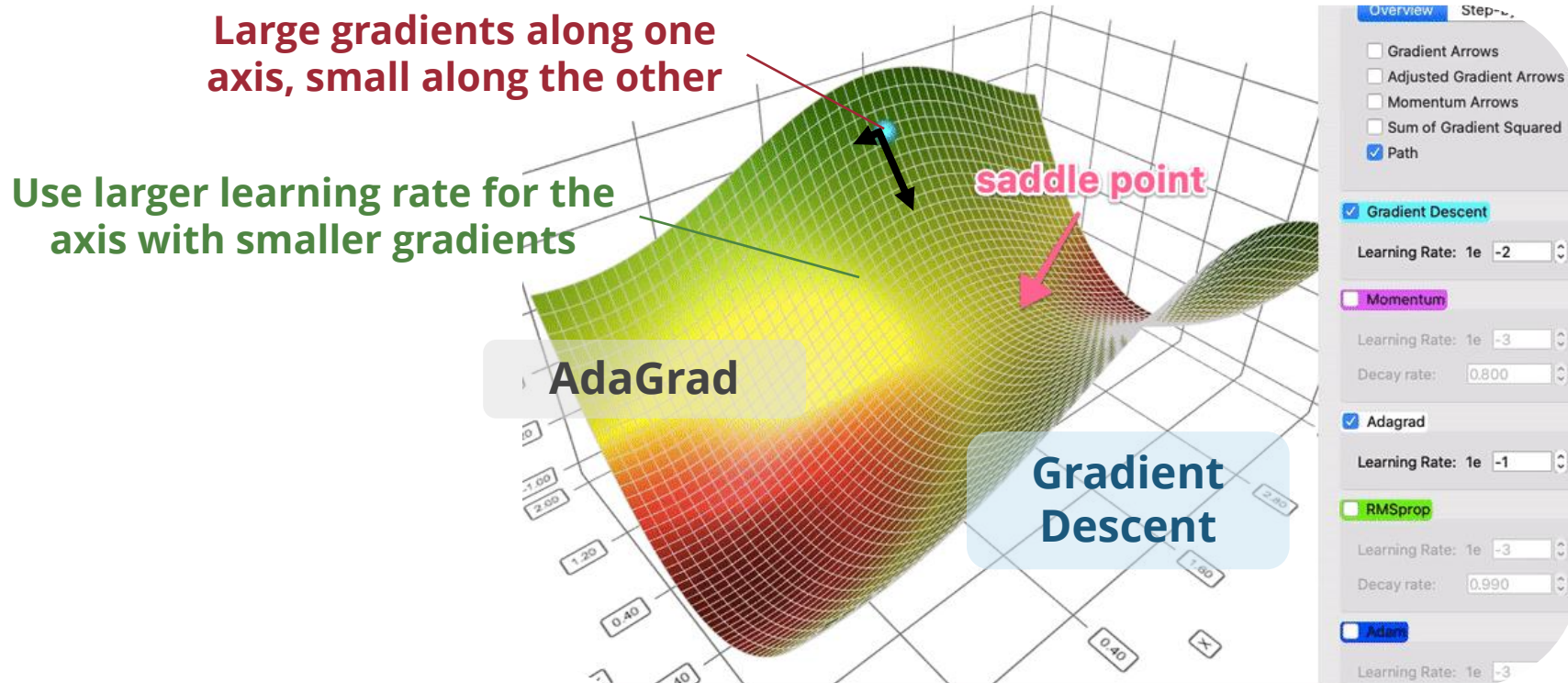
Forward Pass & Backward Pass



Backward pass
loss.backward()

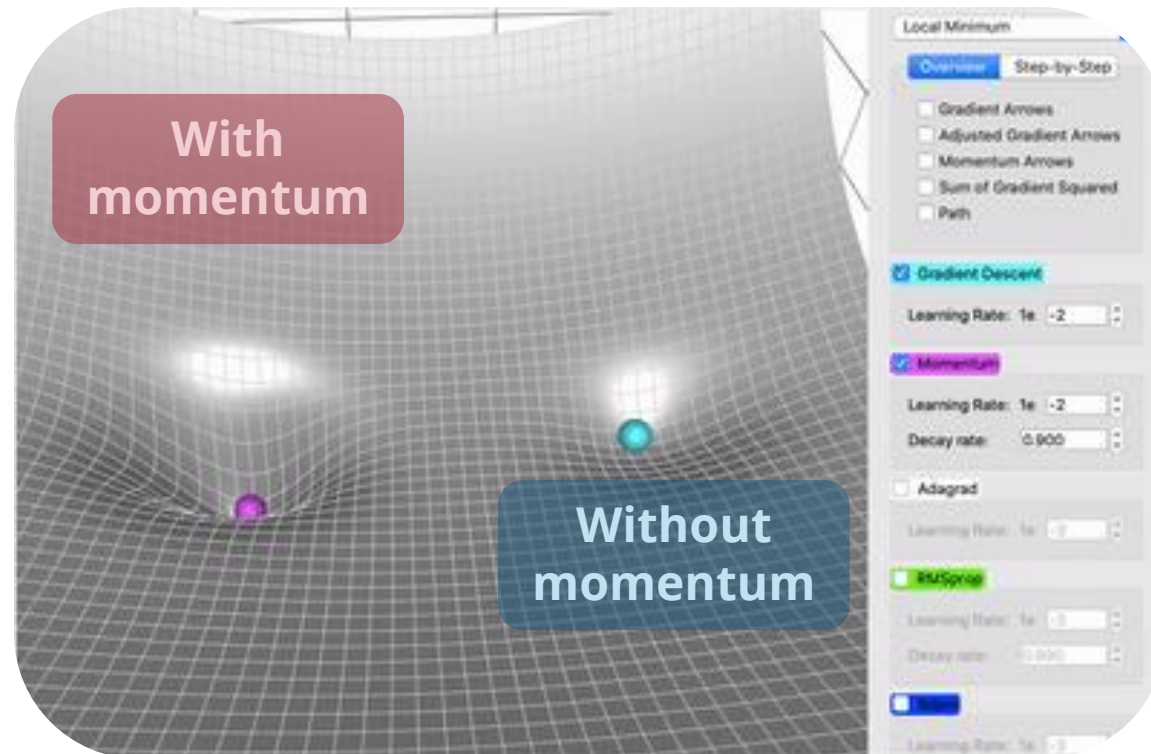
Gradient-based Adaptive Learning Rate

- **Intuition: Compensate axis that has little progress** by comparing the current gradients to the previous gradients



Momentum

- **Intuition:** Maintain the momentum to **escape from local minima**



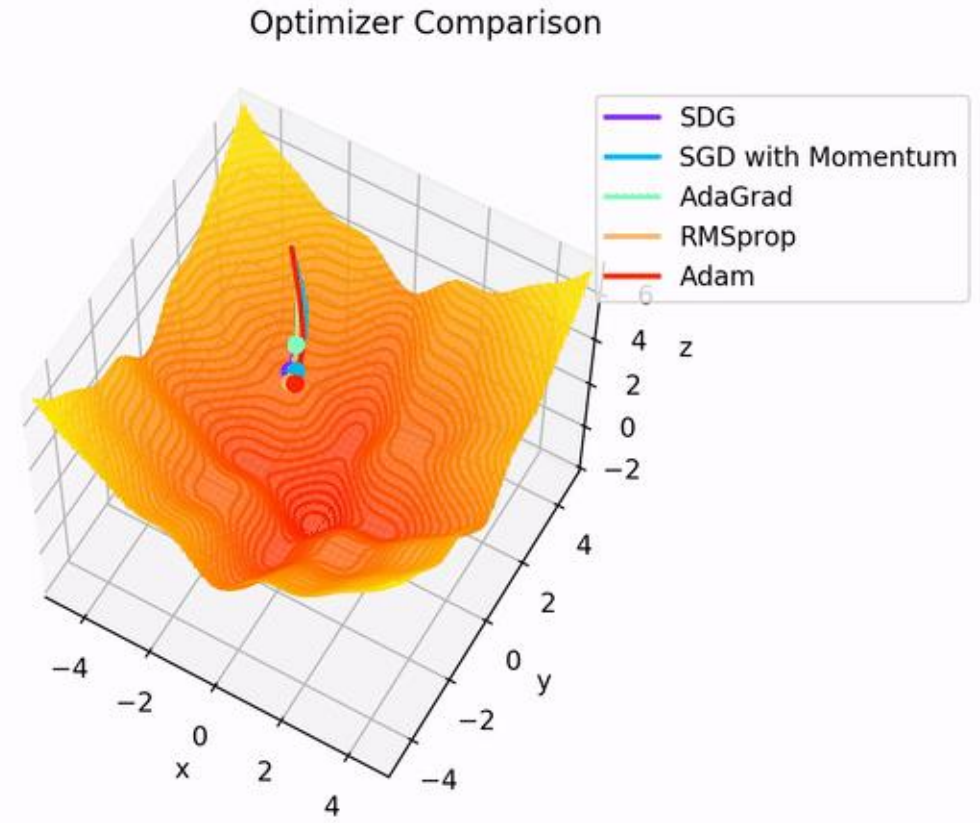
Comparison of Optimizers

- **Momentum**

- Gets you out of spurious local minima
- Allows the model to explore around

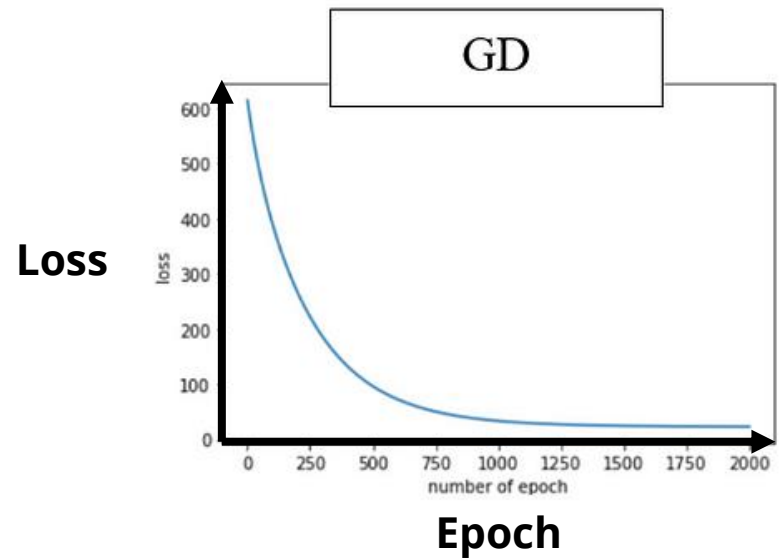
- **Gradient-based adaption**

- Maintains steady improvement
- Allows faster convergence

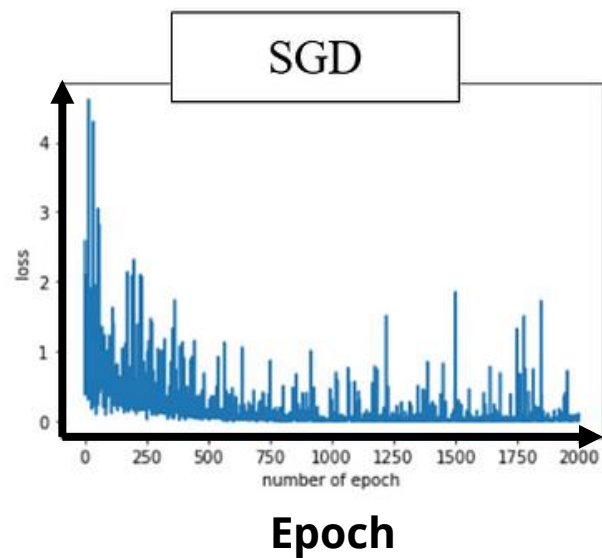


Mini-batch Gradient Descent

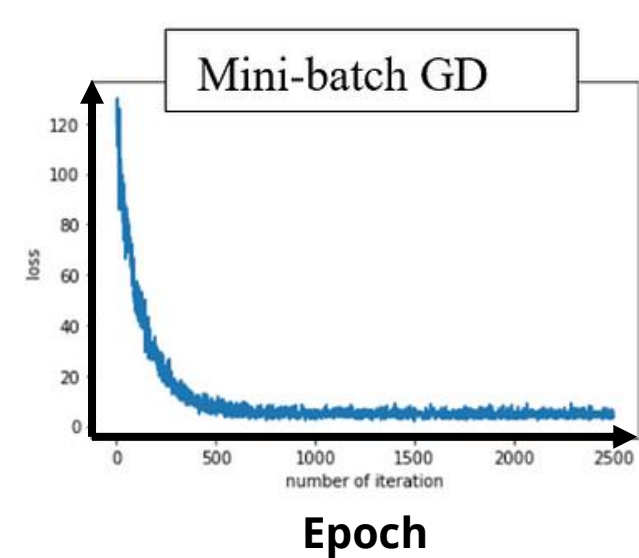
- **Intuition:** Estimate the gradient using **several random training samples**



batch size = N

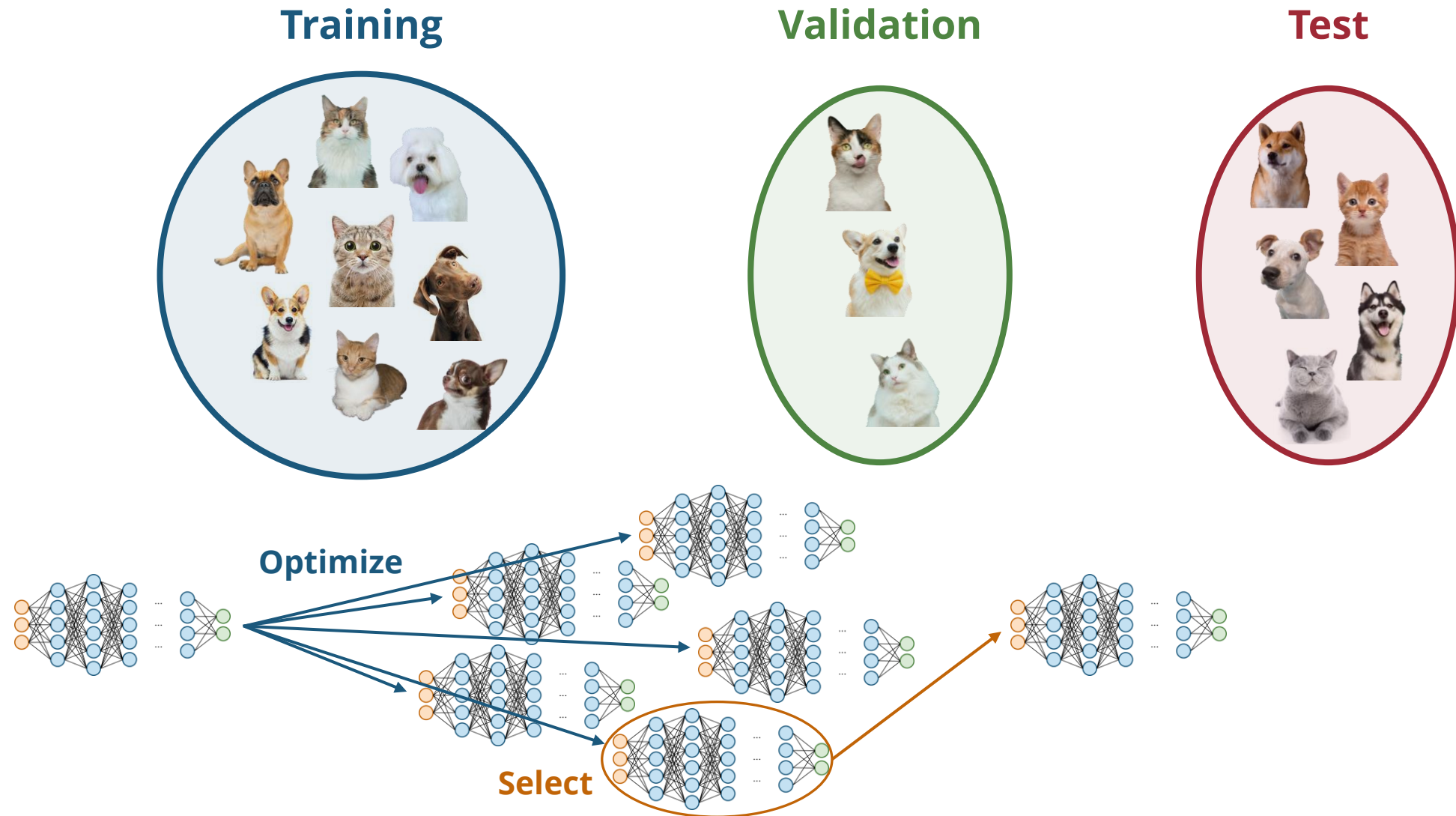


batch size = 1

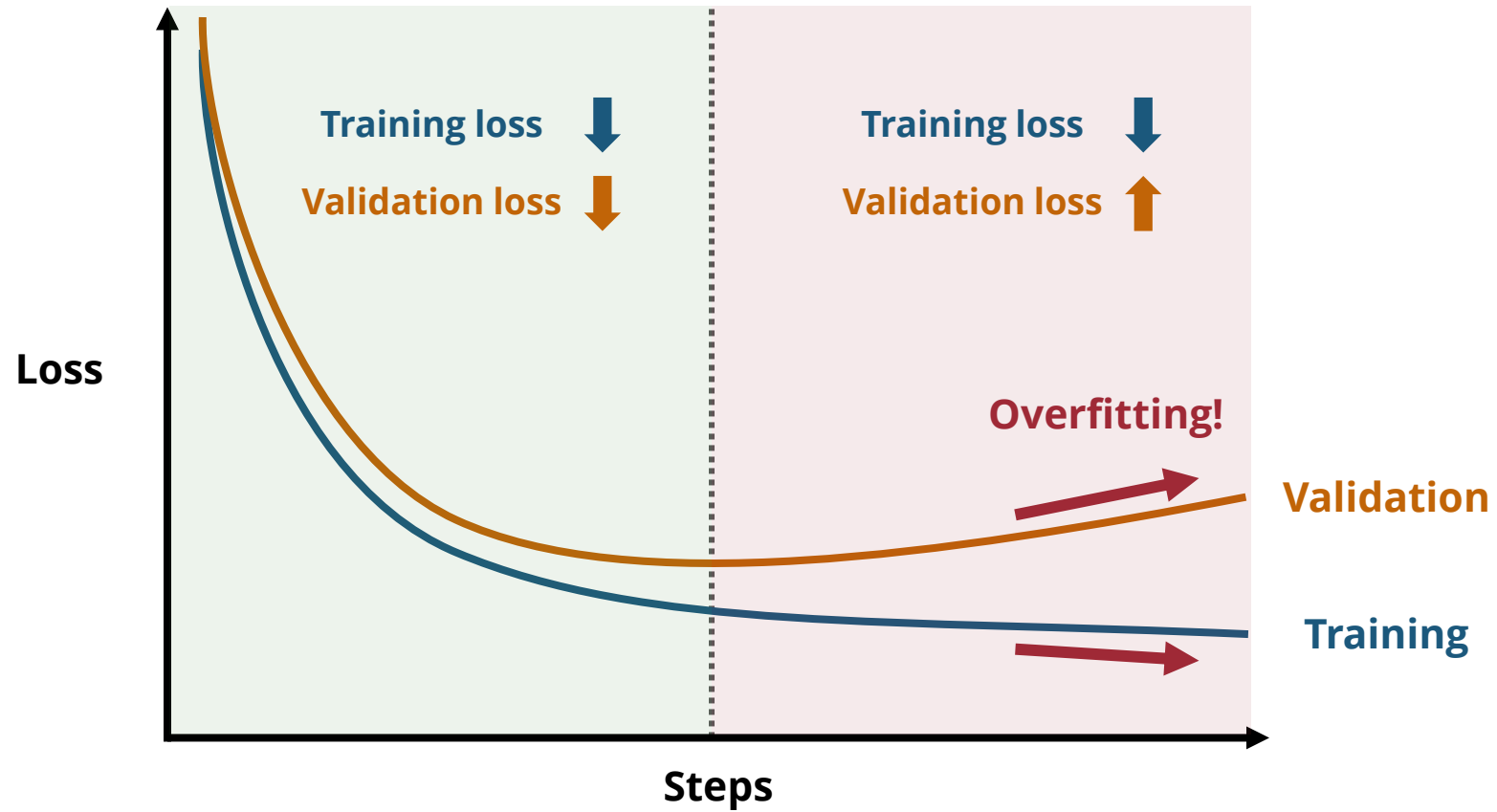


$1 < \text{batch size} < N$

Training-Validation-Test Pipeline



Training vs Validation Losses



Review – Deep Generative Models

Network Architectures vs Training Frameworks

Network architectures

Multilayer perceptron (MLP)

Convolutional neural networks (CNNs)

Recurrent neural networks (RNNs)

Transformers

ResNets

U-Nets

⋮

Training frameworks

Autoregressive

Autoencoders

Variational autoencoders (VAEs)

Generative adversarial networks (GANs)

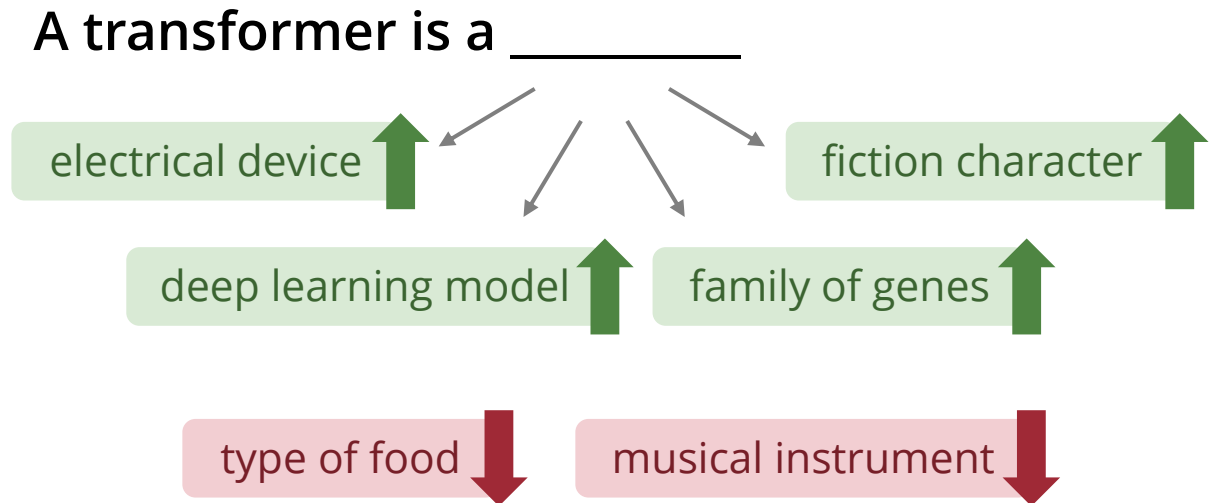
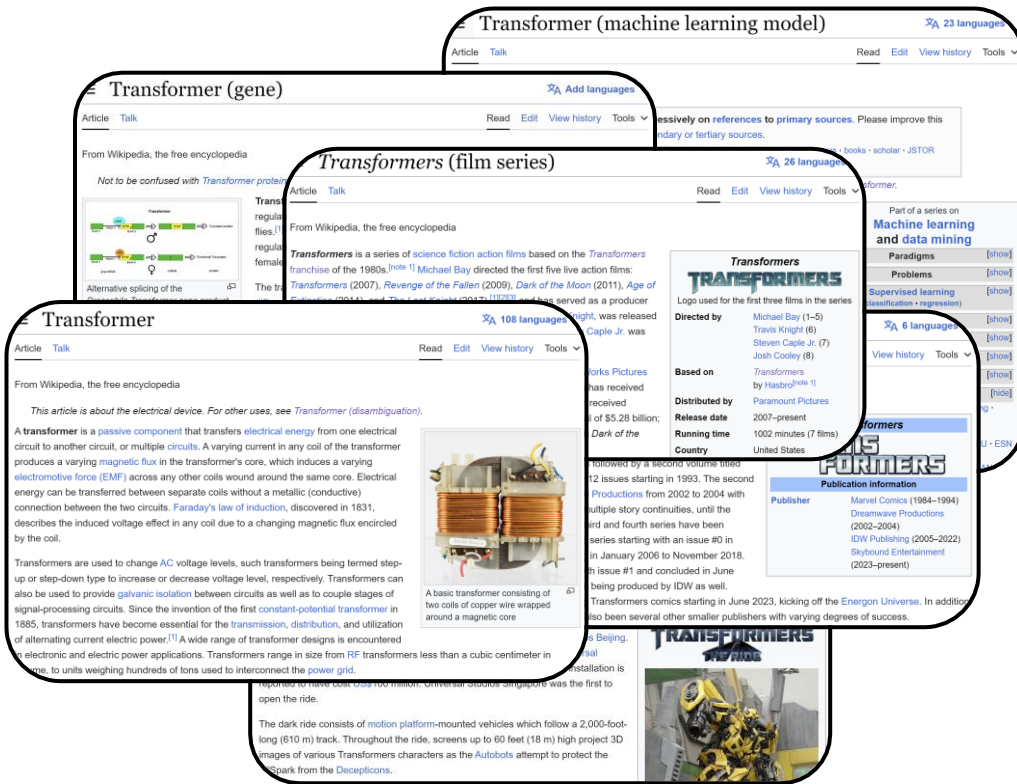
Diffusion models

Consistency models

⋮

Language Models

- Predicting the next word **given the past sequence of words**



Language Models (Mathematically)

- A class of machine learning models that **learn** the next word probability

$$P(x_i \mid x_1, x_2, \dots, x_{i-1})$$

Next word Previous words

$P(\text{electrical} \mid \text{A transformer is a})$	↑
$P(\text{character} \mid \text{A transformer is a})$	↑
$P(\text{gene} \mid \text{A transformer is a})$	↑
$P(\text{model} \mid \text{A transformer is a})$	↑
$P(\text{food} \mid \text{A transformer is a})$	↓
$P(\text{musical} \mid \text{A transformer is a})$	↓

Language Models – Generation

- How do we generate a new sentence using a trained language model?

A transformer is a

→ Model → deep

A transformer is a deep

→ Model → learning

A transformer is a deep learning

→ Model → model

A transformer is a deep learning model

→ Model → introduced

A transformer is a deep learning model introduced

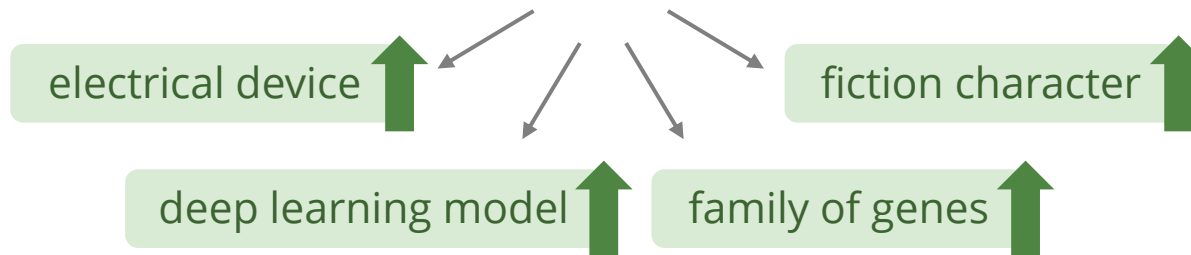
→ Model → in

A transformer is a deep learning model introduced in

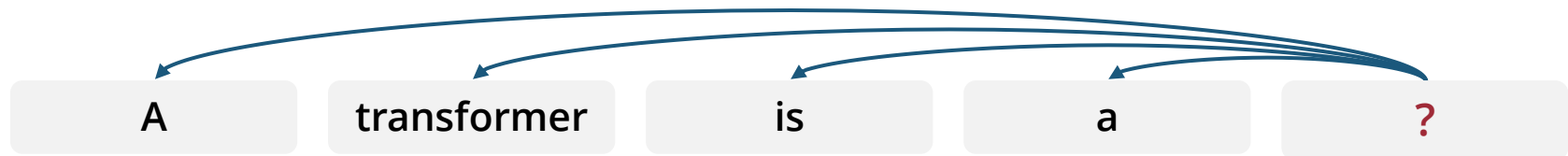
→ Model → 2017

Demystifying Transformers

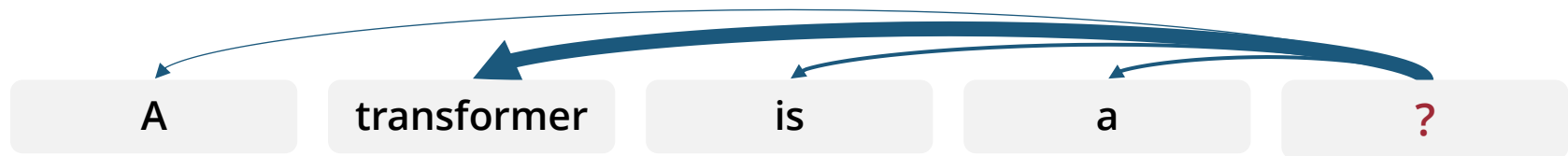
A transformer is a _____



Uniform attention



Variable attention

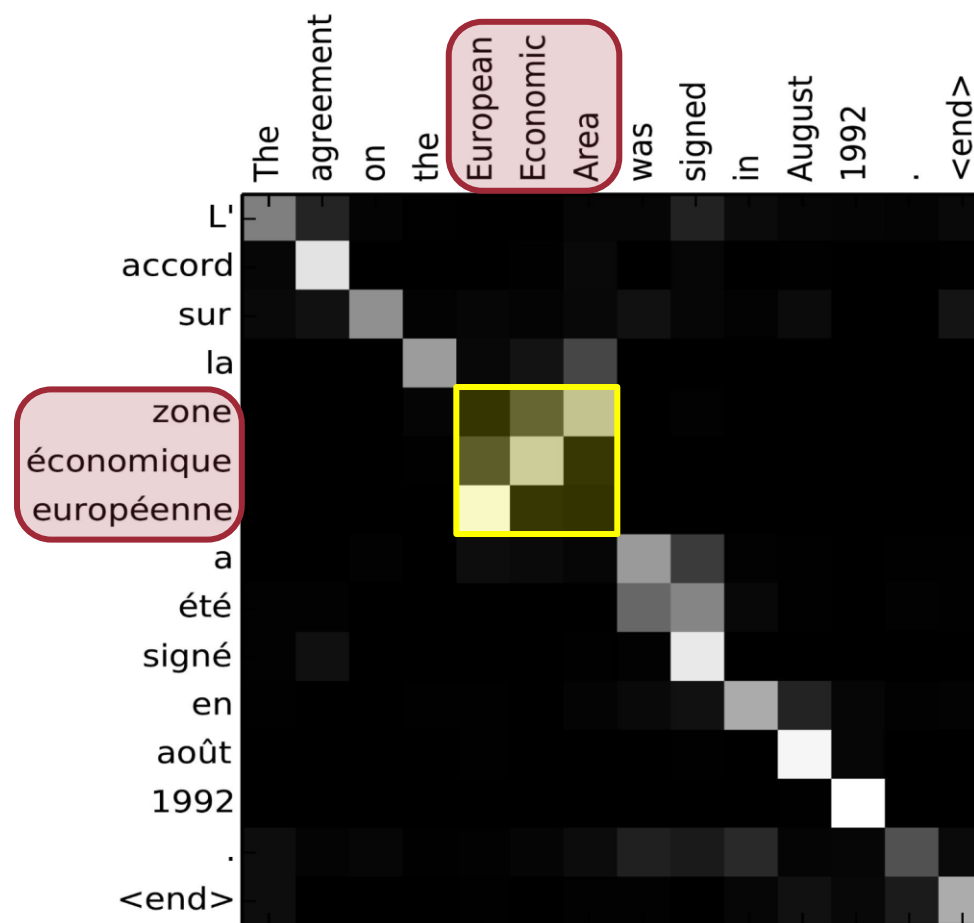


Transformers learn what to attend to from big data!

What does a Transformer Learn?

The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .
The FBI is chasing a criminal on the run .

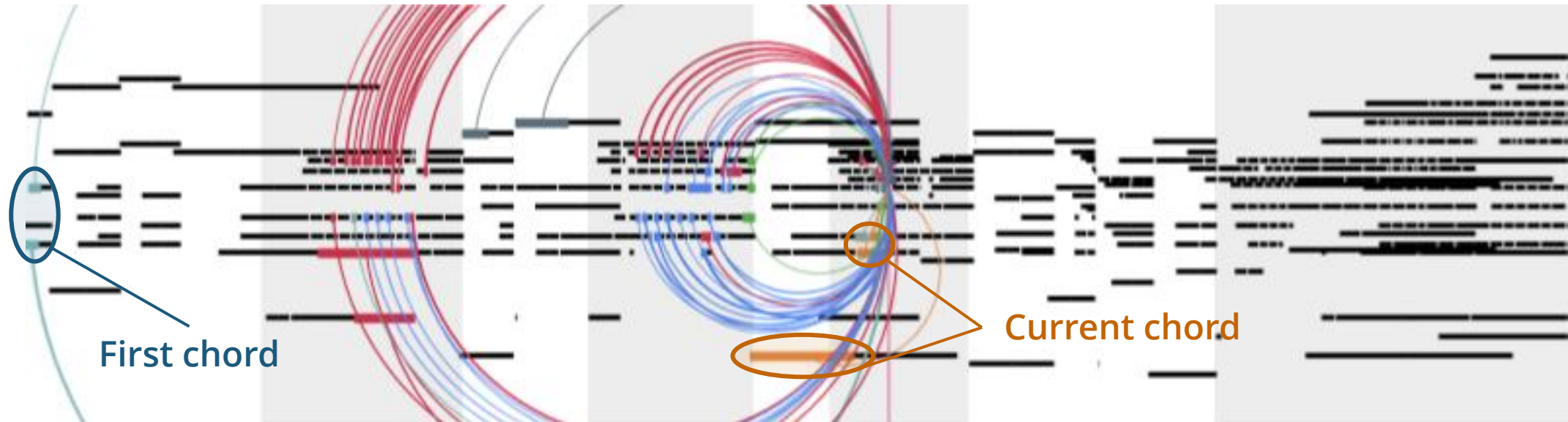
(Source: Cheng et al., 2016)



(Source: Bahdanau et al., 2015)

What does a Transformer Learn?

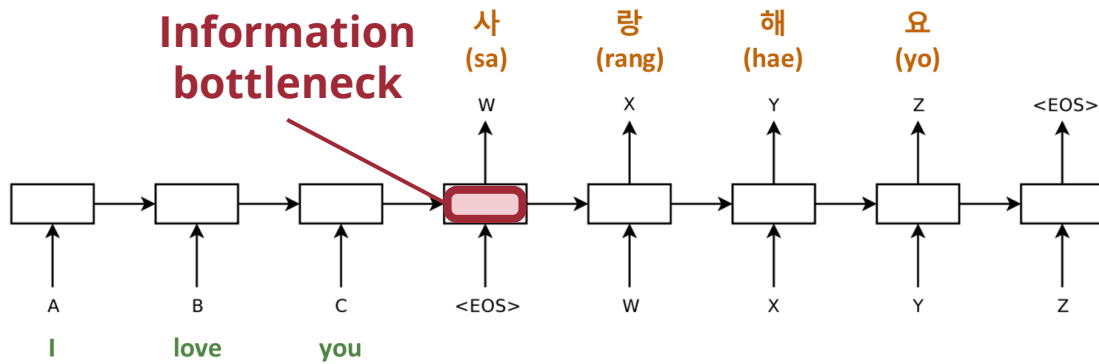
(Each color represents an attention head)



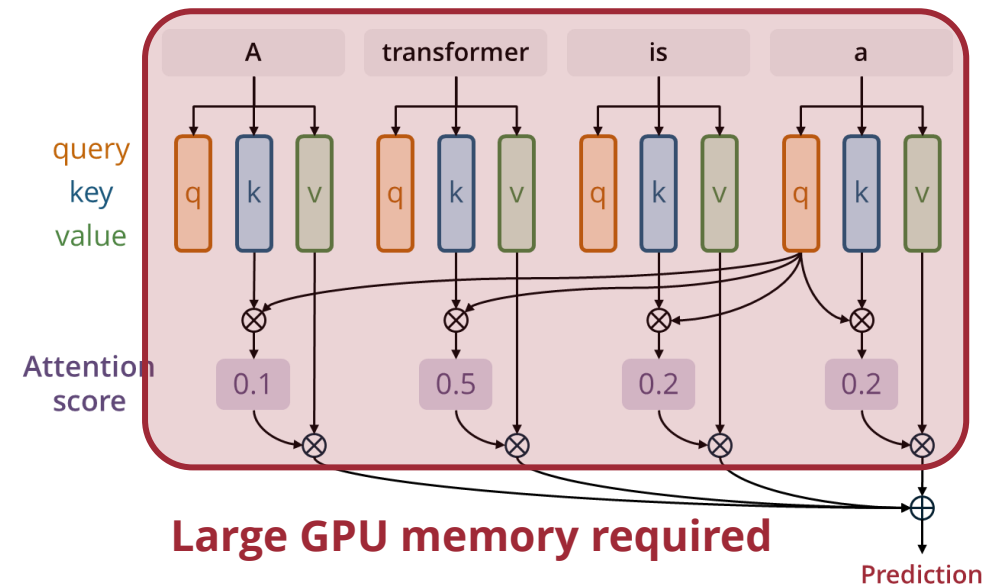
(Source: Huang et al., 2018)

Seq2seq vs Transformers

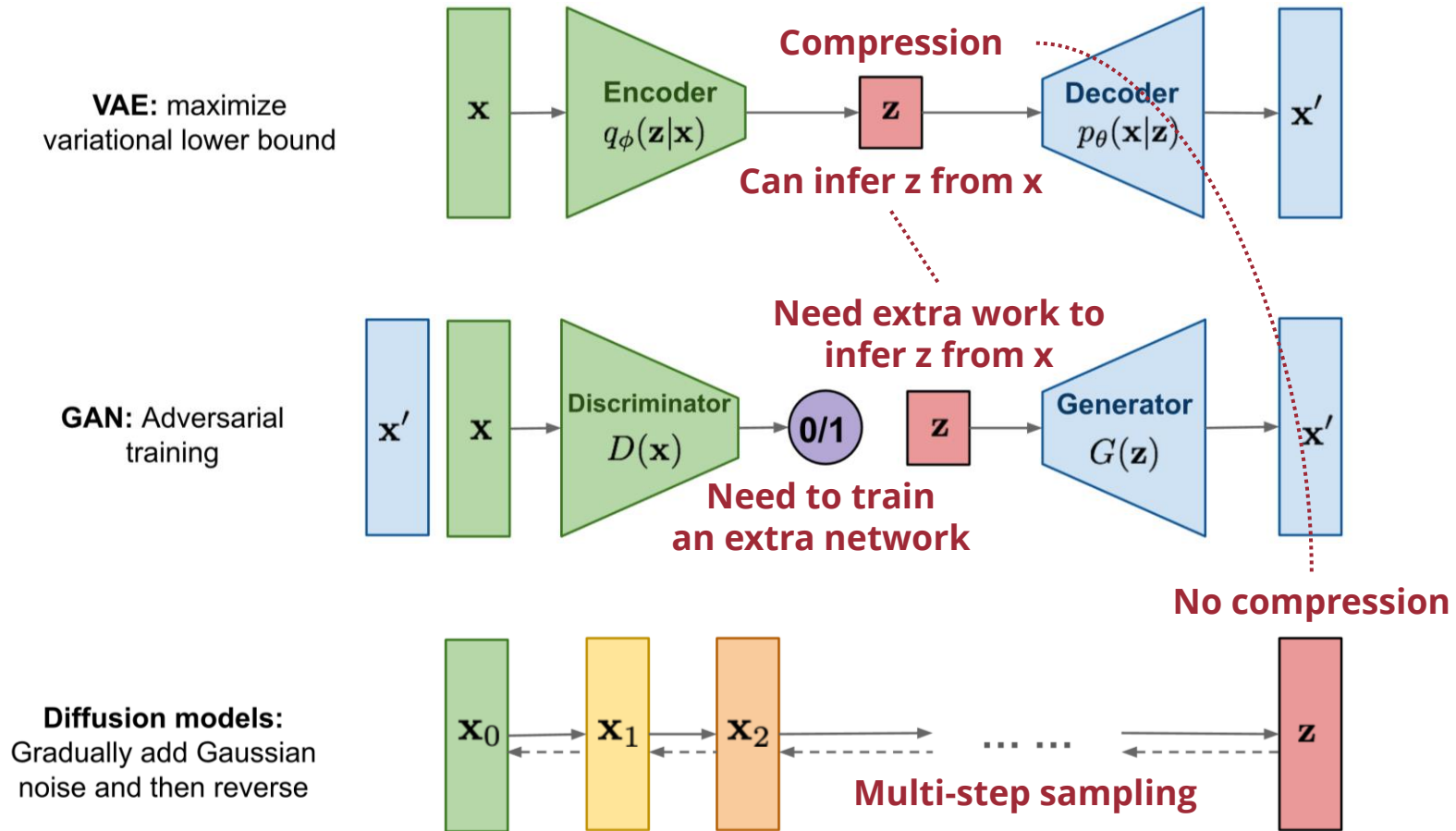
Seq2seq



Transformers



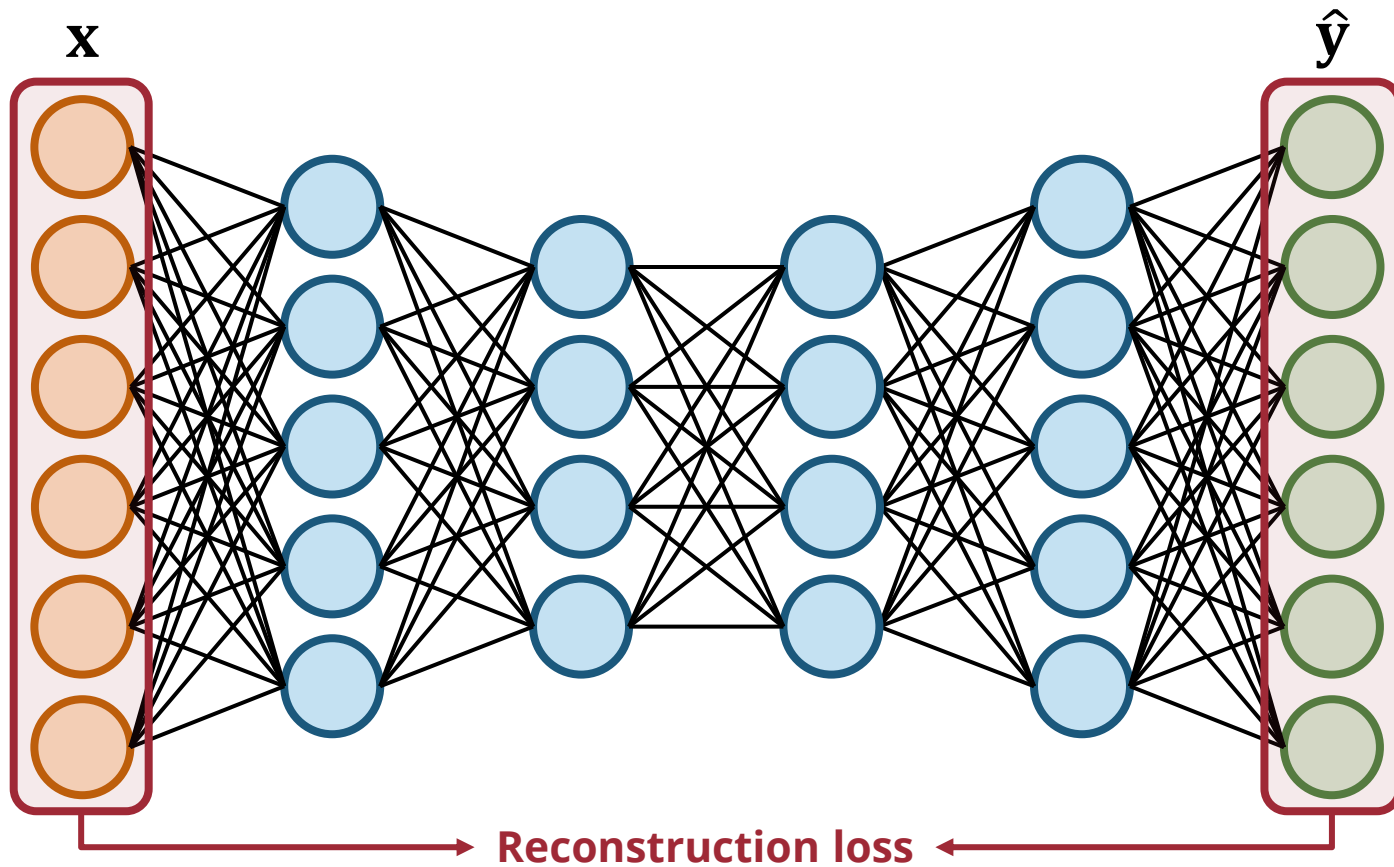
Comparison of Deep Generative Models



(Source: Weng, 2021)

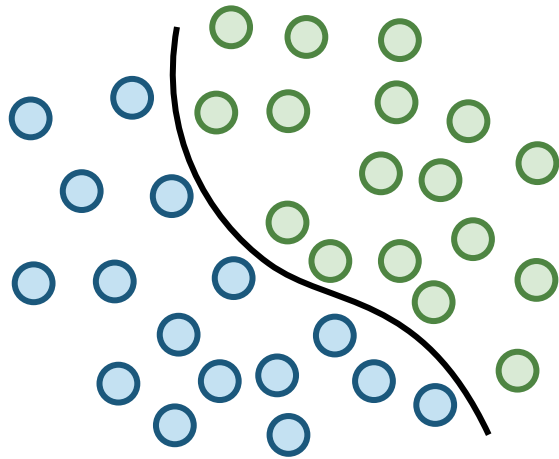
Autoencoders

- A neural network where the **input and output are the same**



Discriminative vs Generative Models

Discriminative



Discriminative models learn the decision boundary

$$P(y|x)$$

Generative

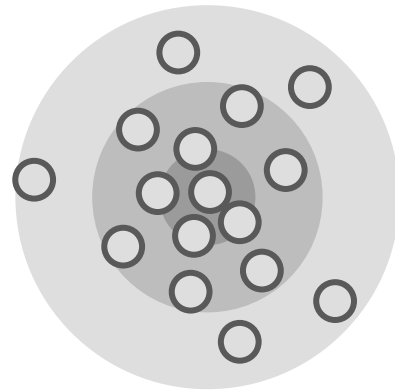


Generative models learn the underlying distribution

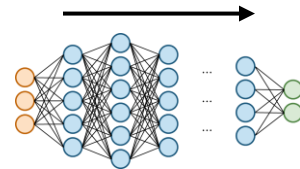
$$P(x) \text{ or } P(x|y)$$

Generating Data from a Random Distribution

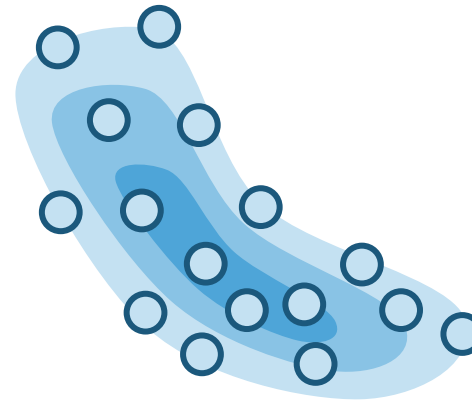
Random distribution



$P(z)$



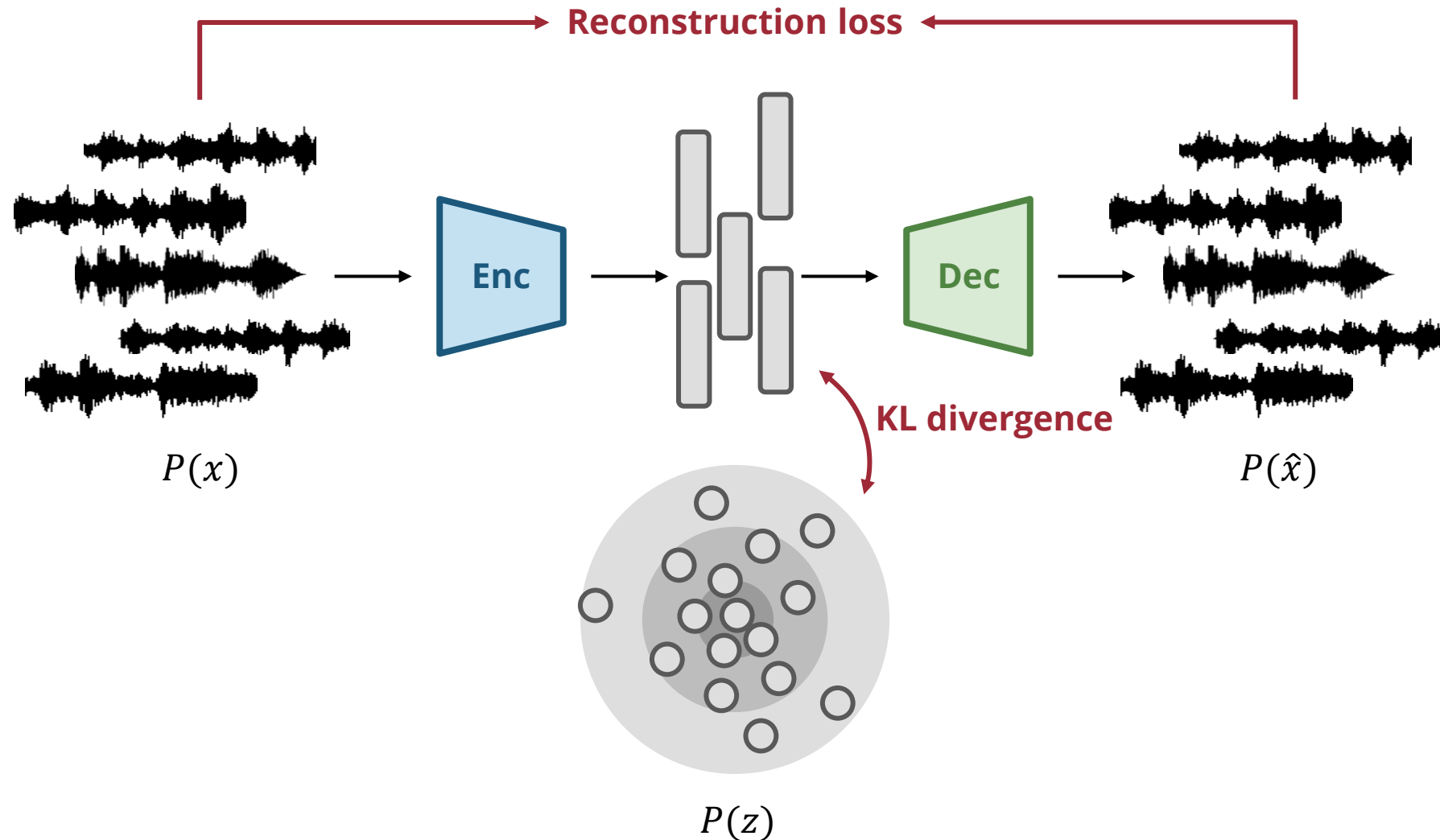
Data distribution



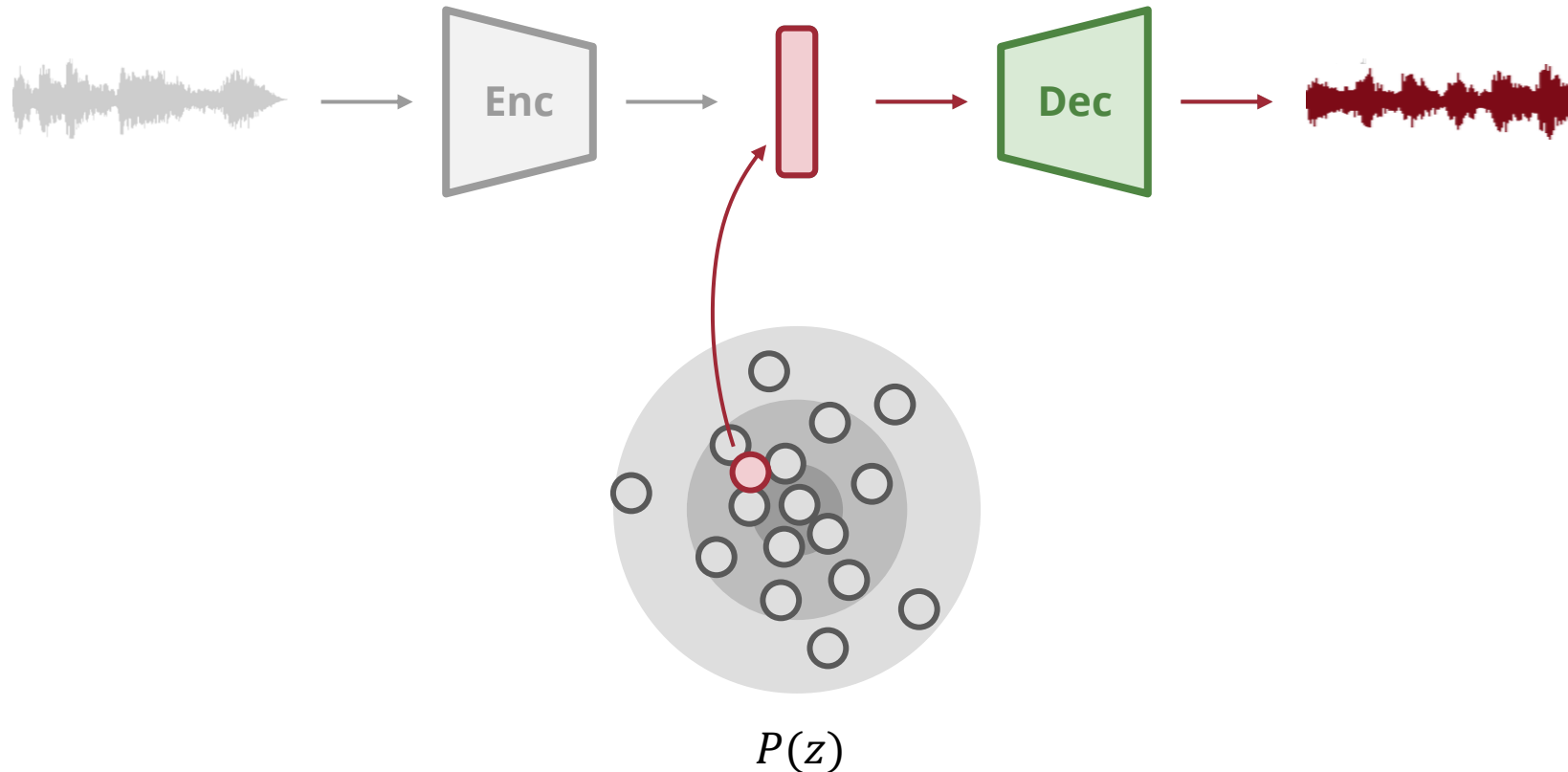
$P(x)$

If we can learn this mapping, we can easily generate new samples from the data distribution

Variational Autoencoders (VAEs) – Training

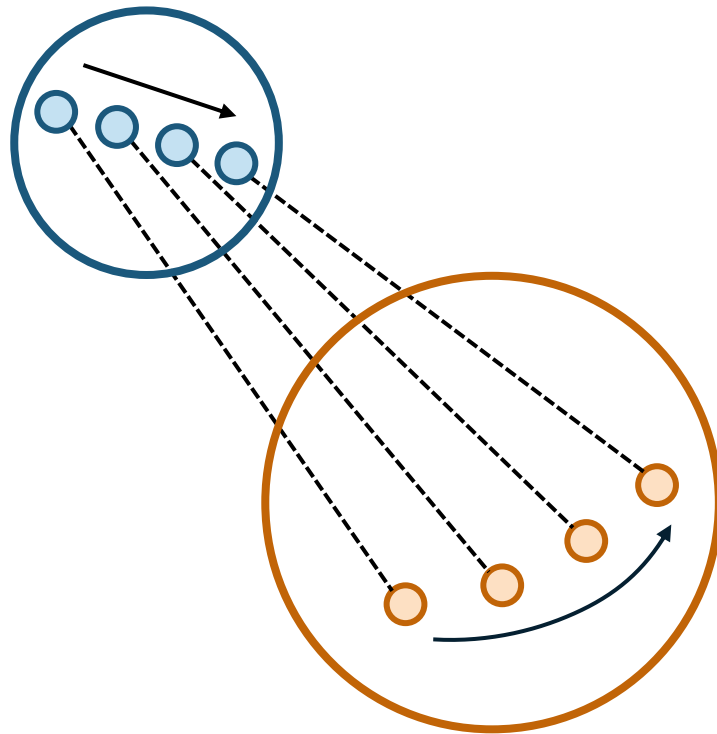


Variational Autoencoders (VAEs) – Generation

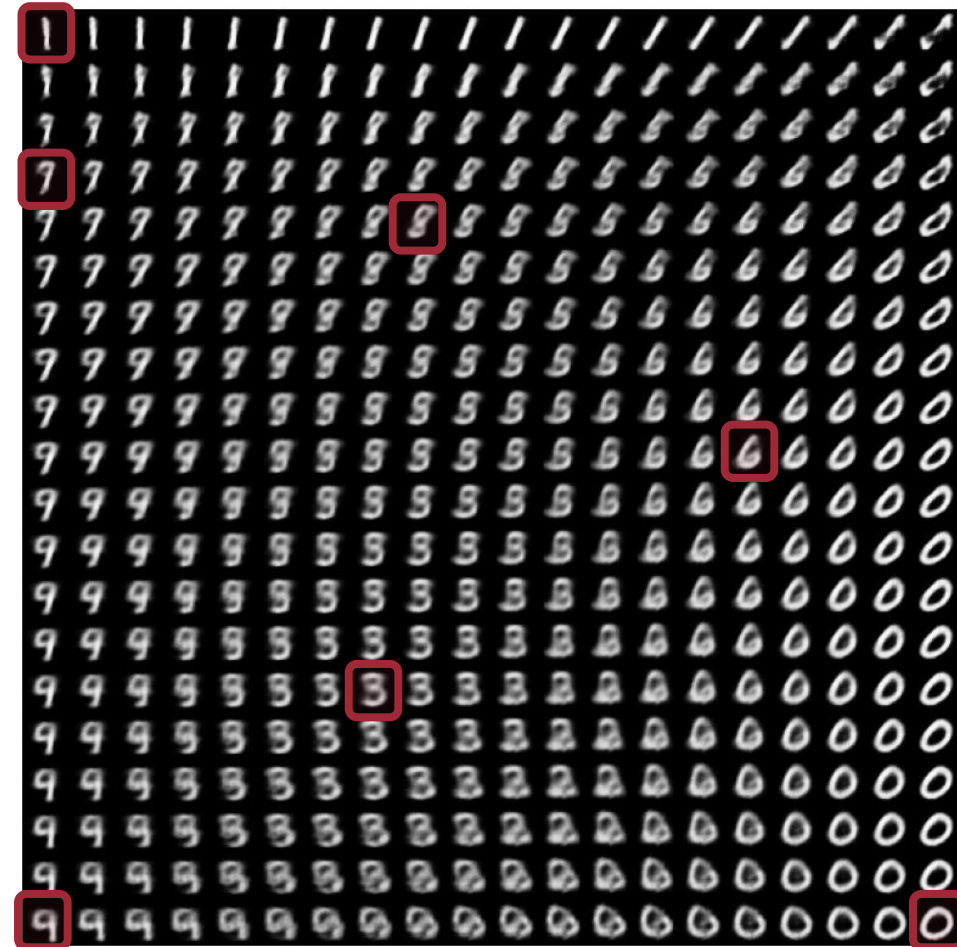


Decoding the Latent Space of a VAE

Latent space



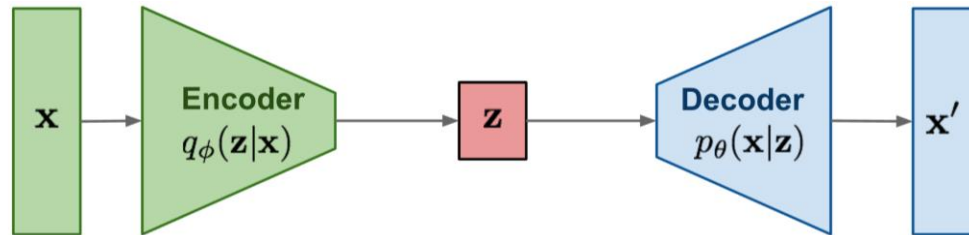
Data space



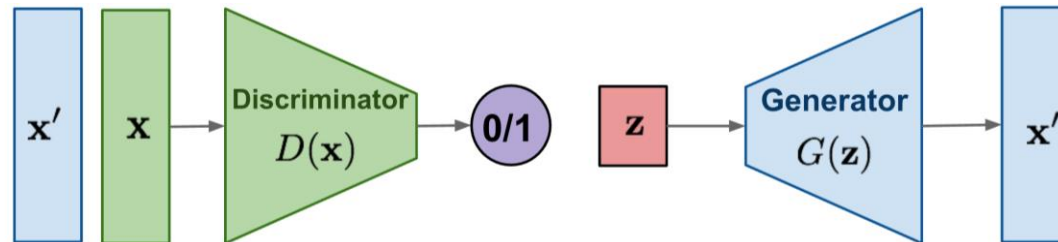
(Source: tensorflow.org)

Comparison of Deep Generative Models

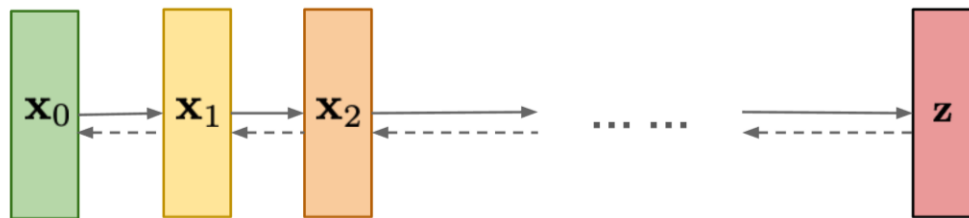
VAE: maximize variational lower bound



GAN: Adversarial training



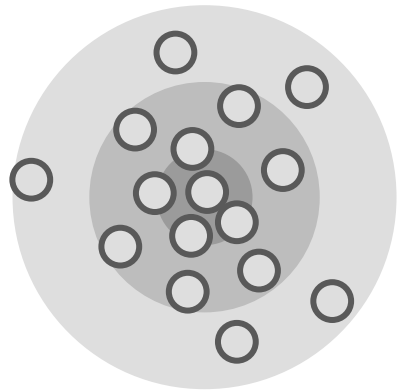
Diffusion models:
Gradually add Gaussian noise and then reverse



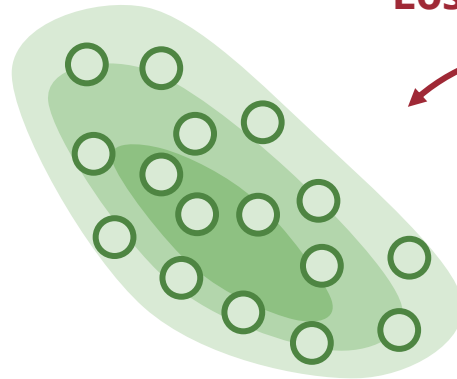
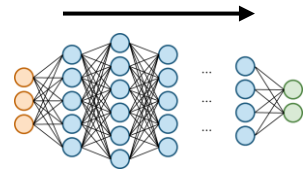
(Source: Weng, 2021)

A Loss Function for Distributions

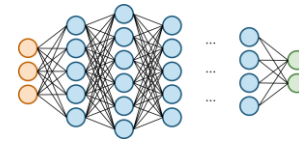
Random distribution



$P(z)$

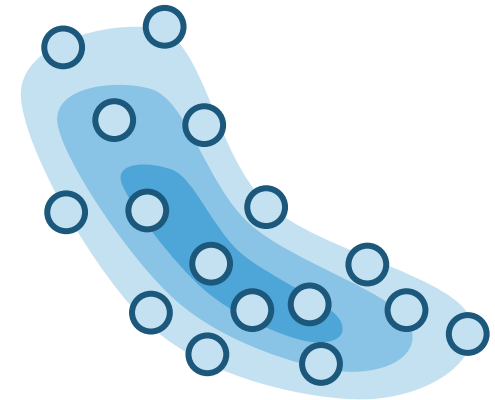


$P(\hat{x})$



Loss function?

Data distribution

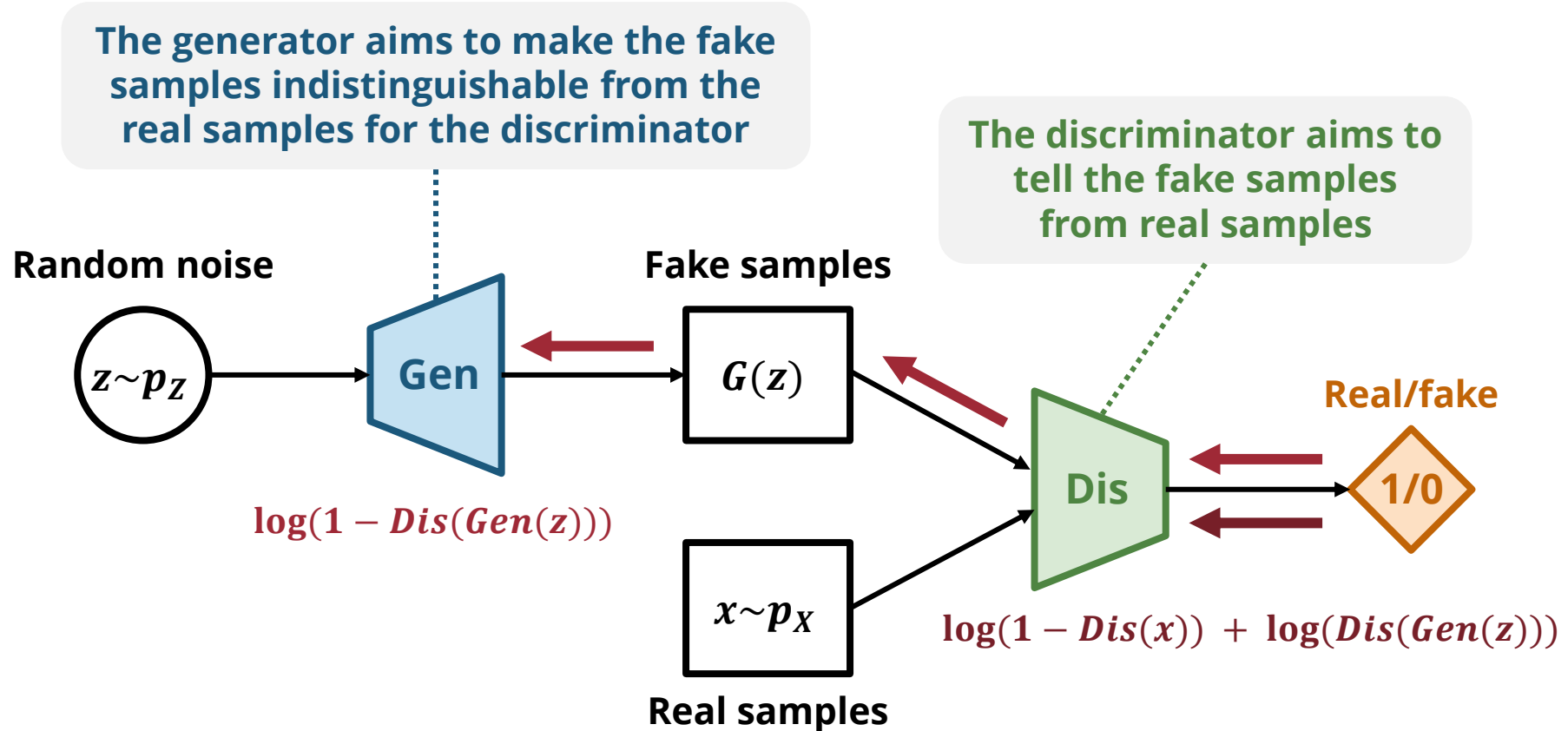


$P(x)$

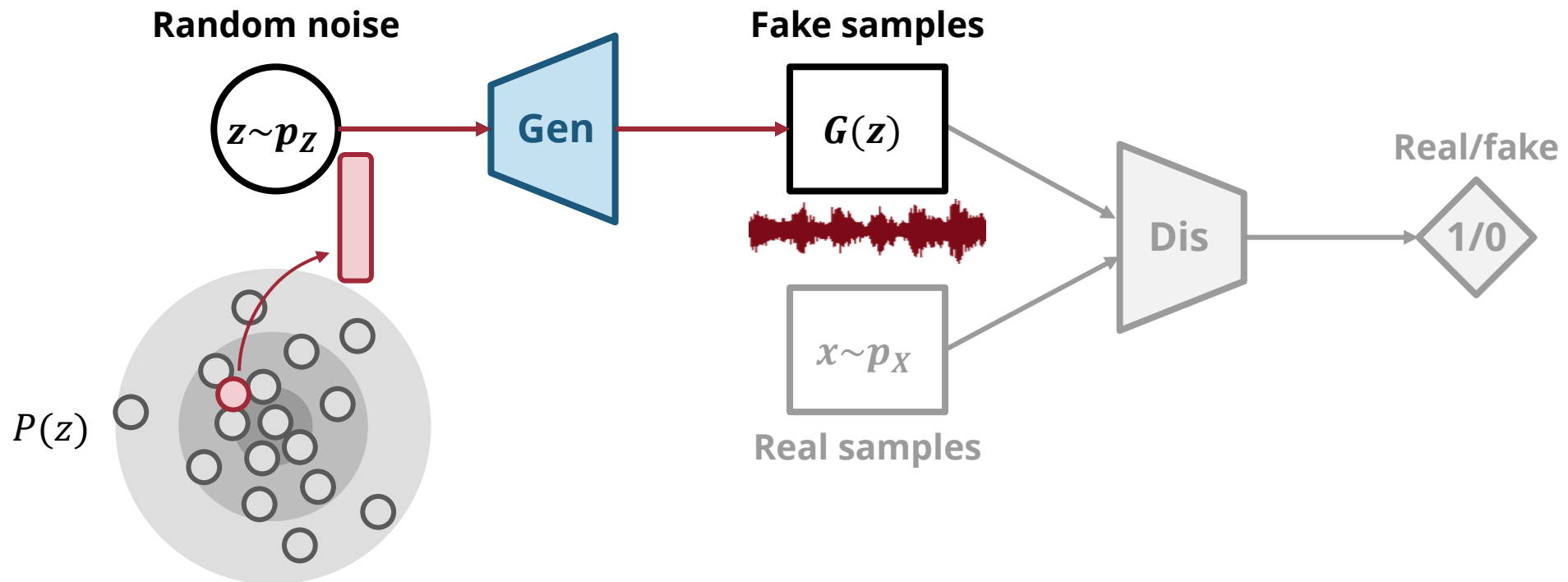
Unfortunately, no easy way to measure the difference between two distributions

But what about another neural network!?

Generative Adversarial Nets (GANs) – Training

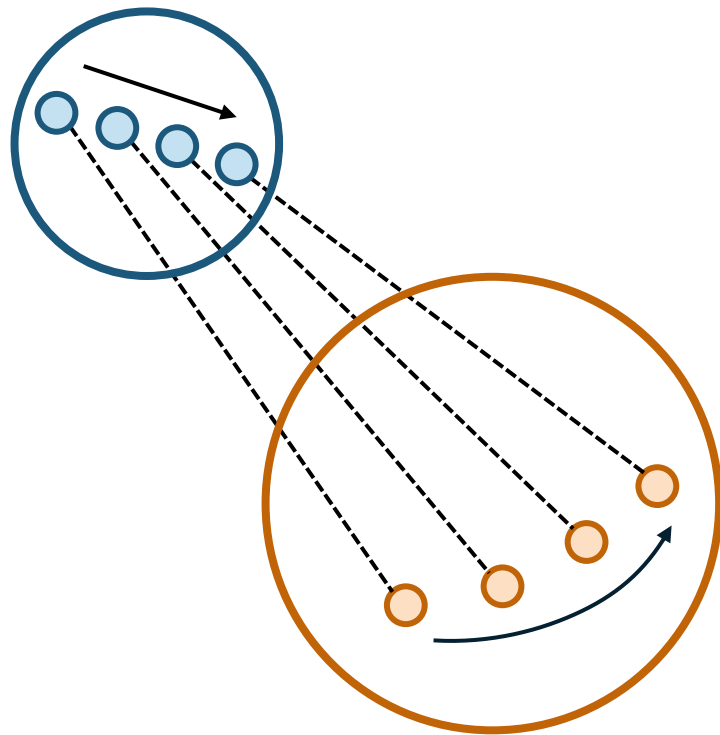


Generative Adversarial Nets (GANs) – Generation

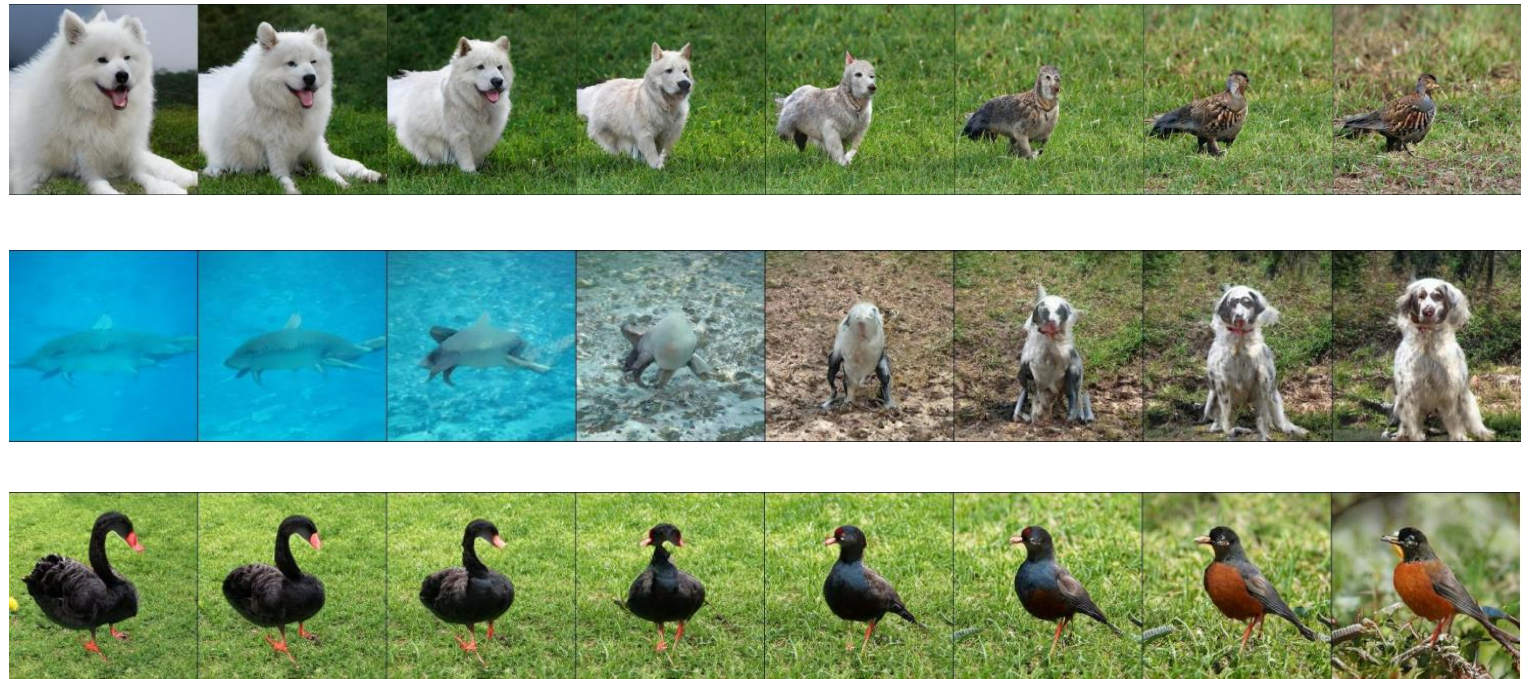


Interpolation on the Latent Space

Latent space



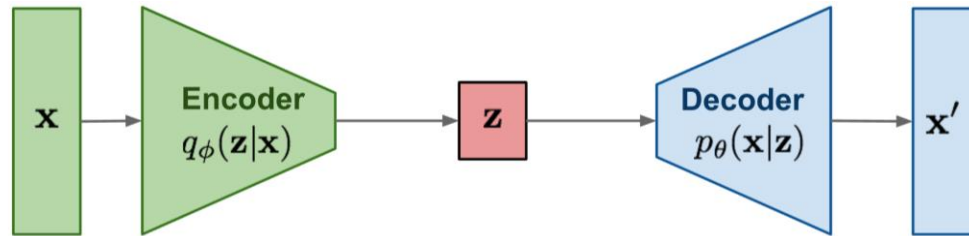
Data space



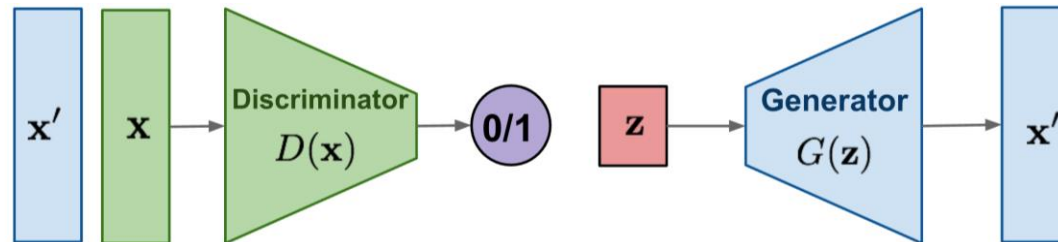
(Source: Brock et al., 2019)

Comparison of Deep Generative Models

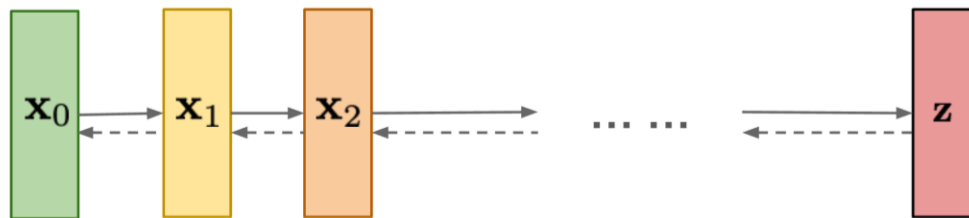
VAE: maximize variational lower bound



GAN: Adversarial training



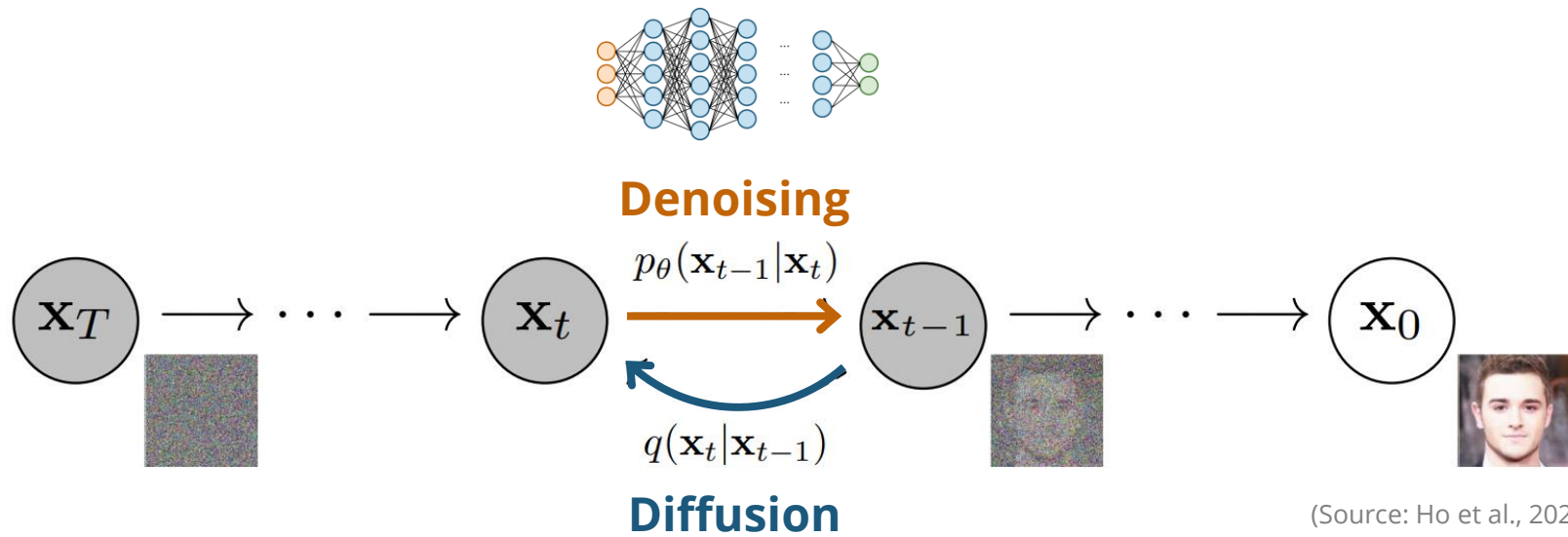
Diffusion models:
Gradually add Gaussian noise and then reverse



(Source: Weng, 2021)

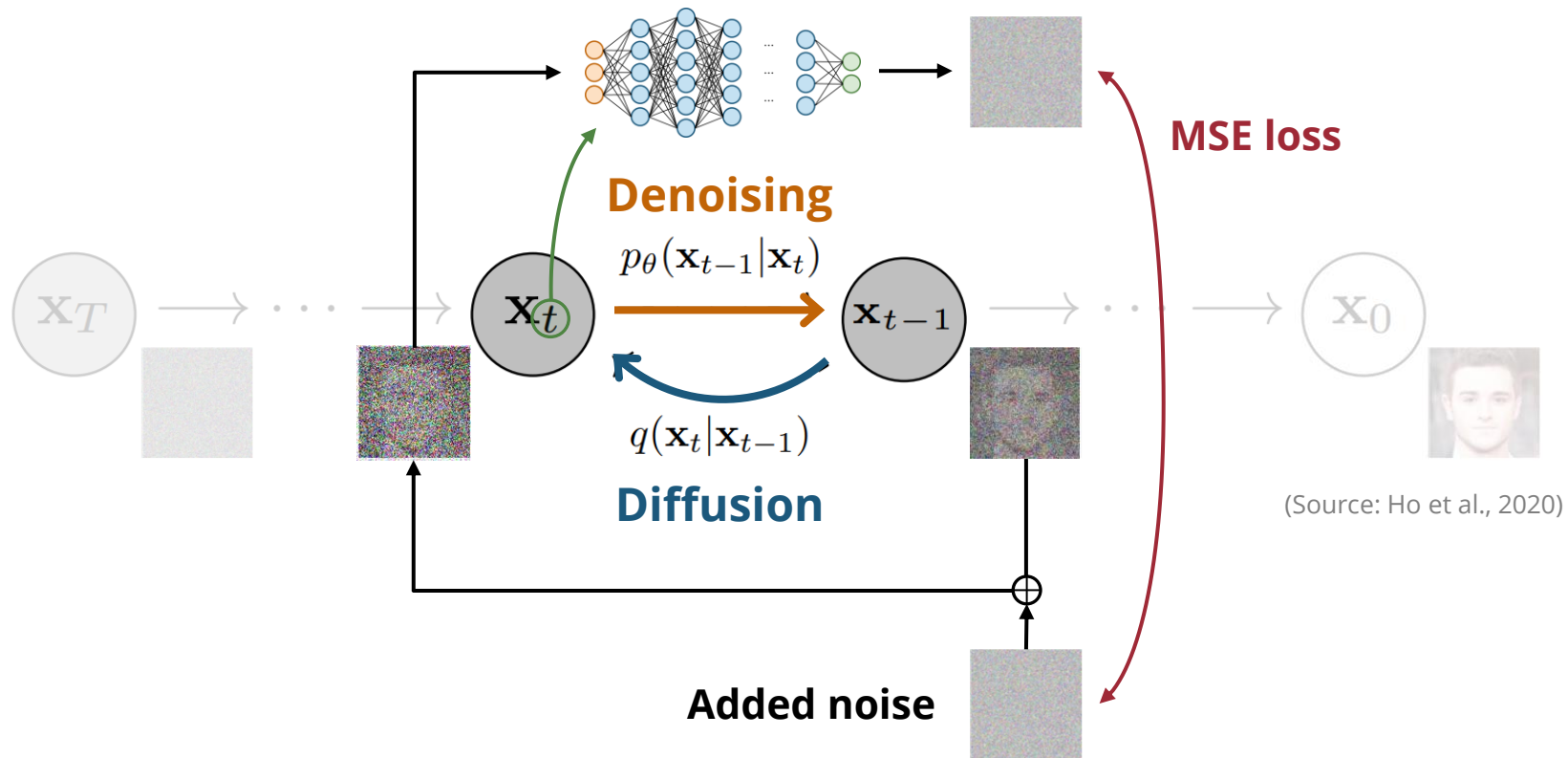
Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



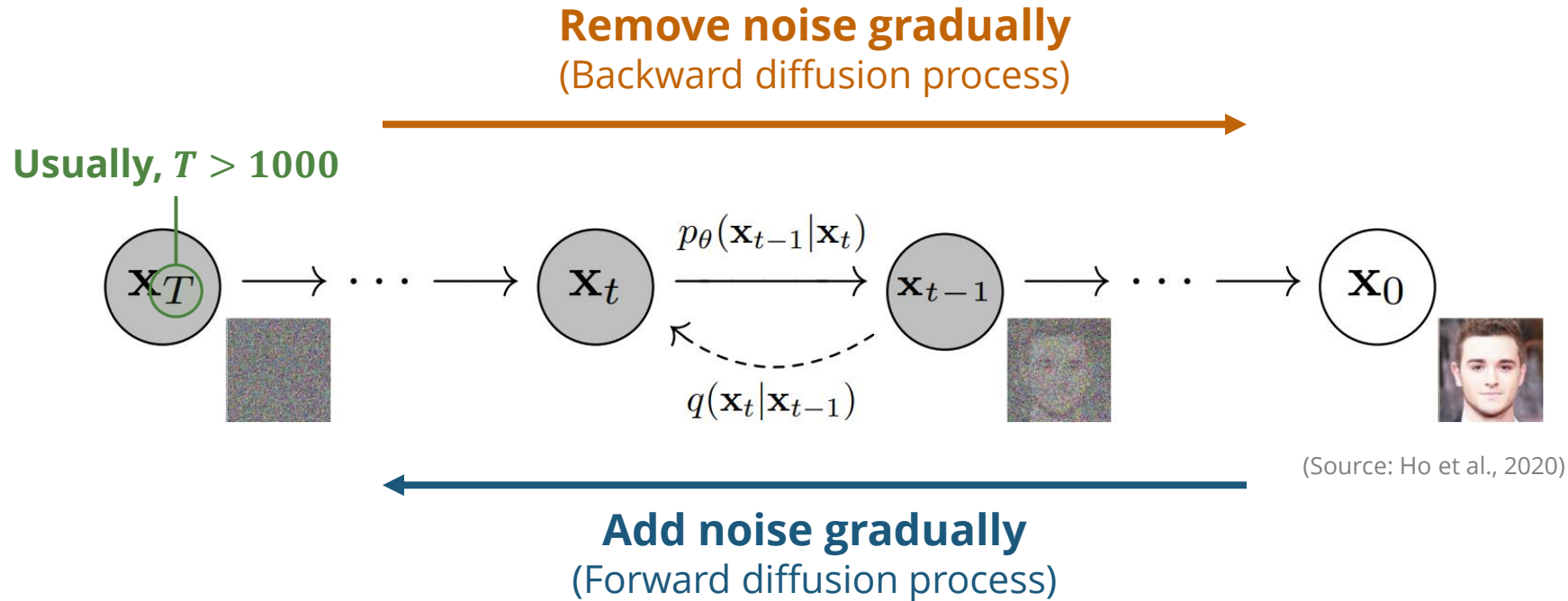
Diffusion Models – Training

- **Intuition:** Many denoising autoencoders stacked together

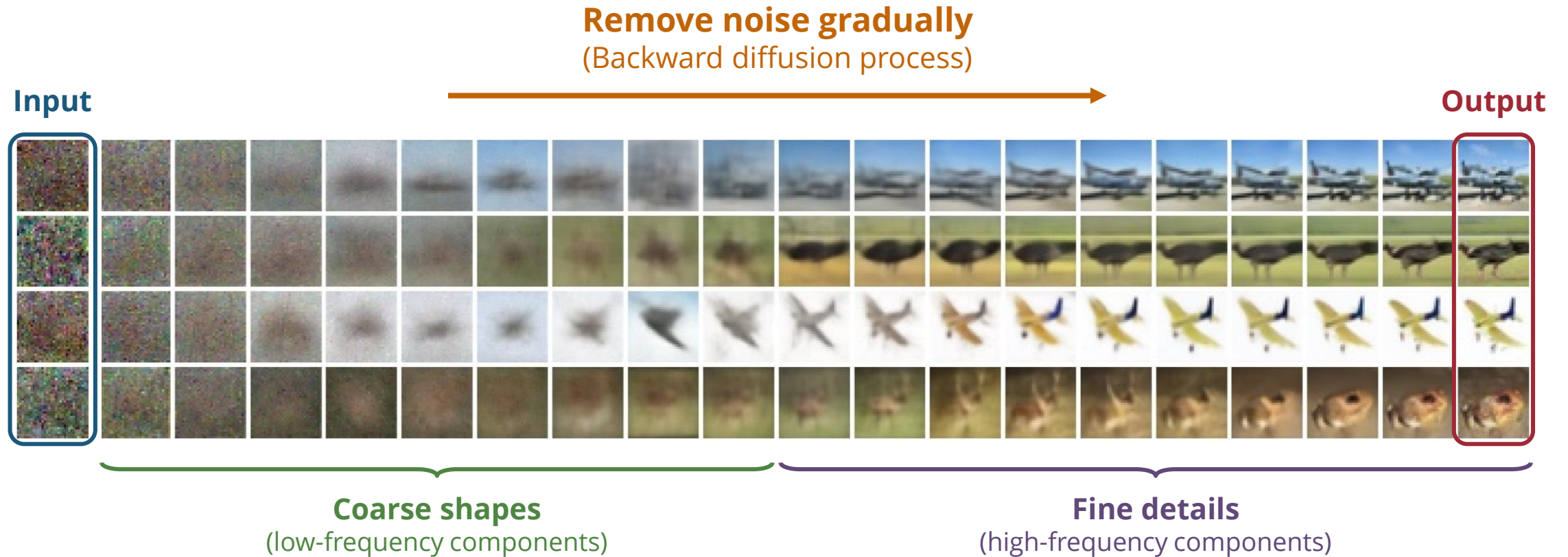


Diffusion Models

- **Intuition**: Many denoising autoencoders stacked together



Diffusion Models – Generation



(Source: Ho et al., 2020)

Discussions

Discussions

- **Why** do we need machine learning?
- **When** should we use deep learning?
- **When can't** we use deep learning?

- Should we choose the best model based on validation loss or accuracy?
- Which generative model works best?
- Is overfitting always an issue?

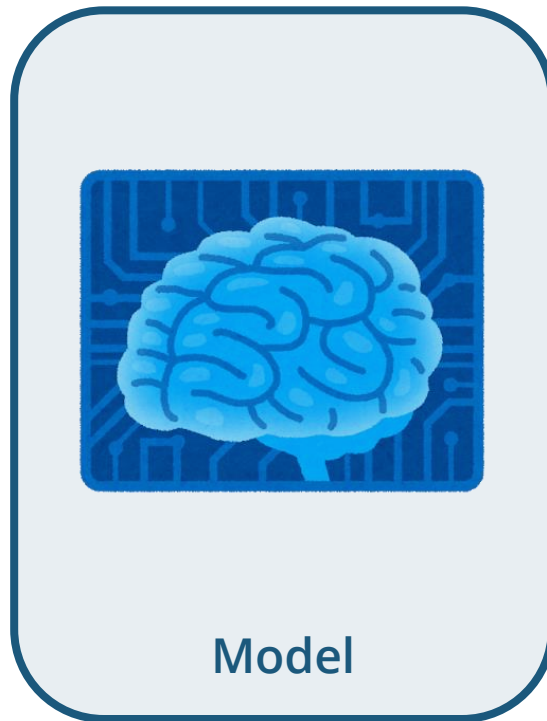
Building Blocks of Modern AI Systems



Data

What we'll talk about next!

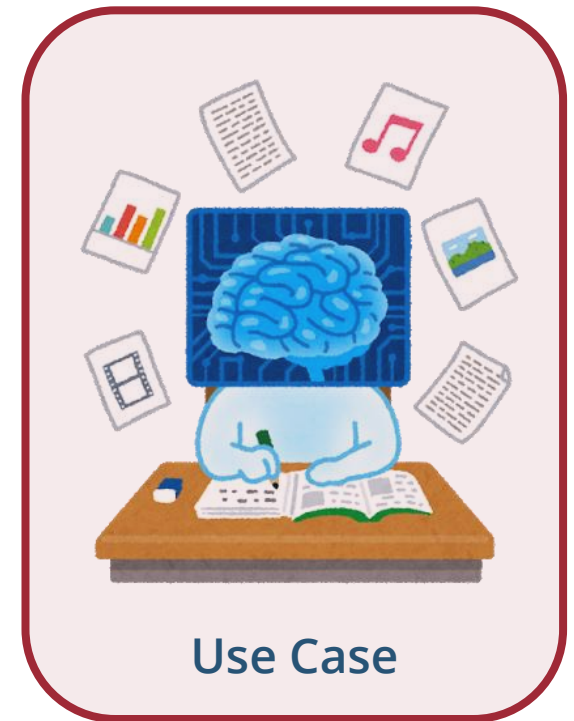
×



Model

What we've been focused on!

×



Use Case

What we'll talk about next!

That's It for the **First Half** of This Course!

