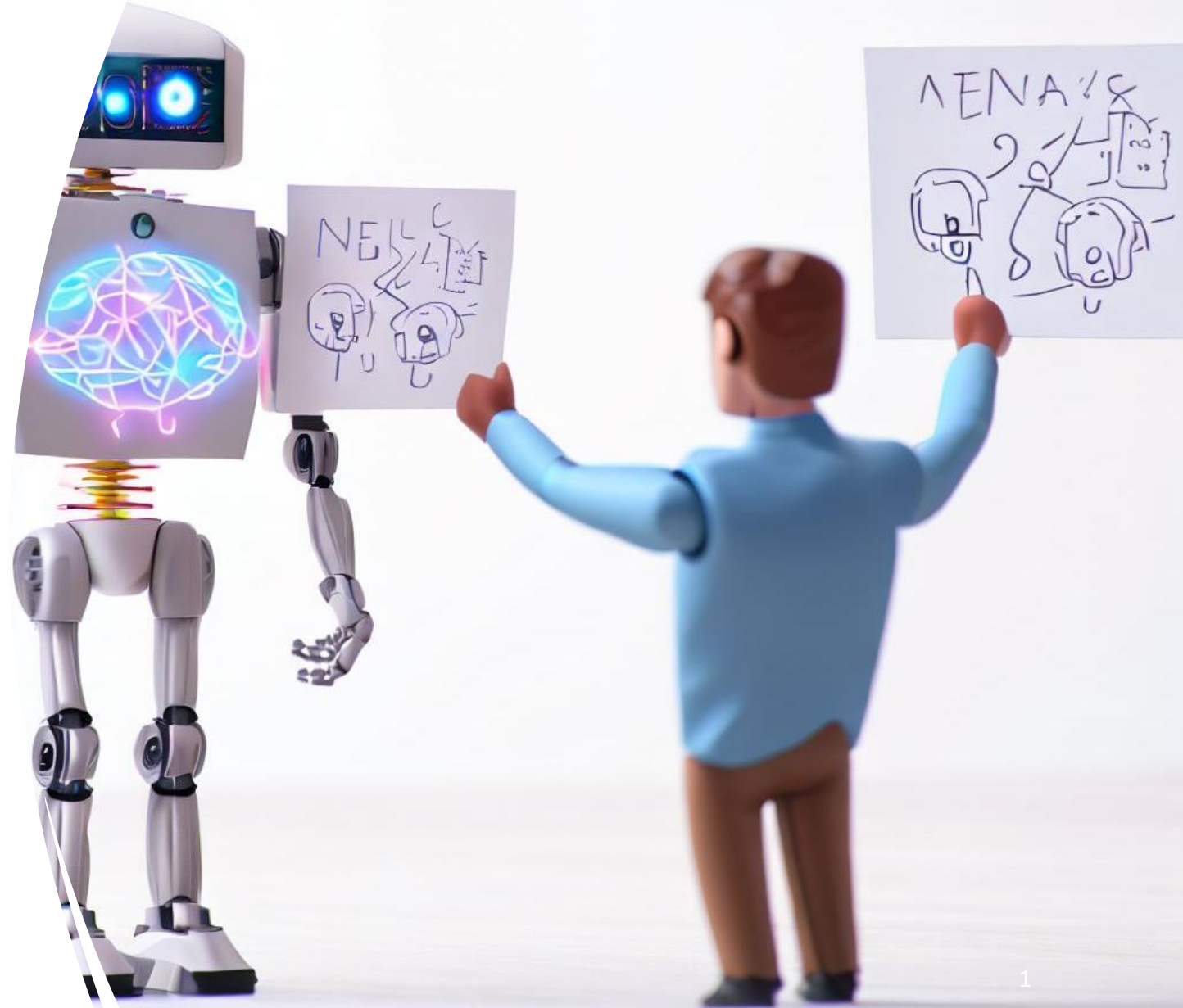
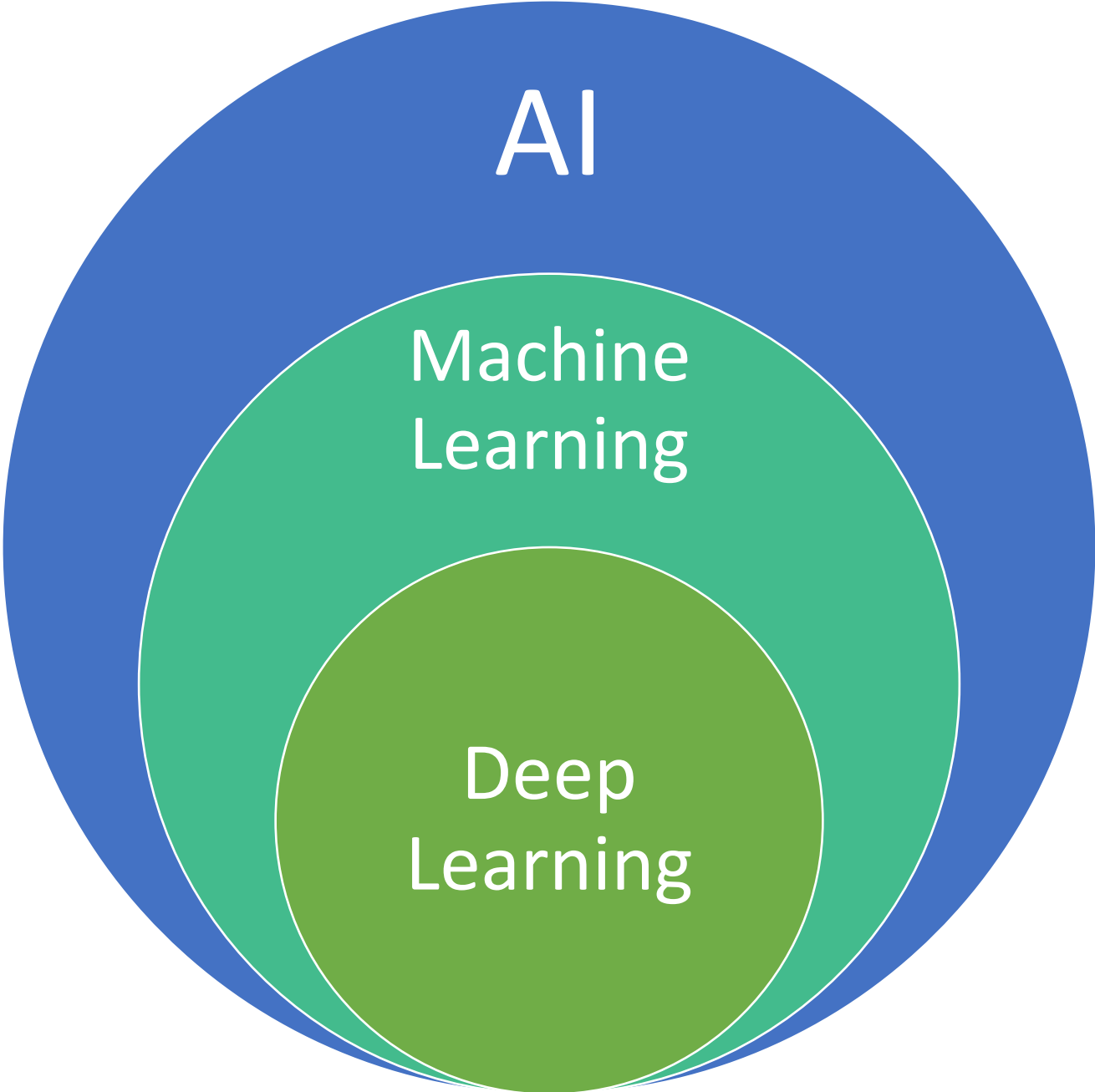


Introduction to Deep Learning

Prof. Kuan-Ting Lai
2023/9/19





AI

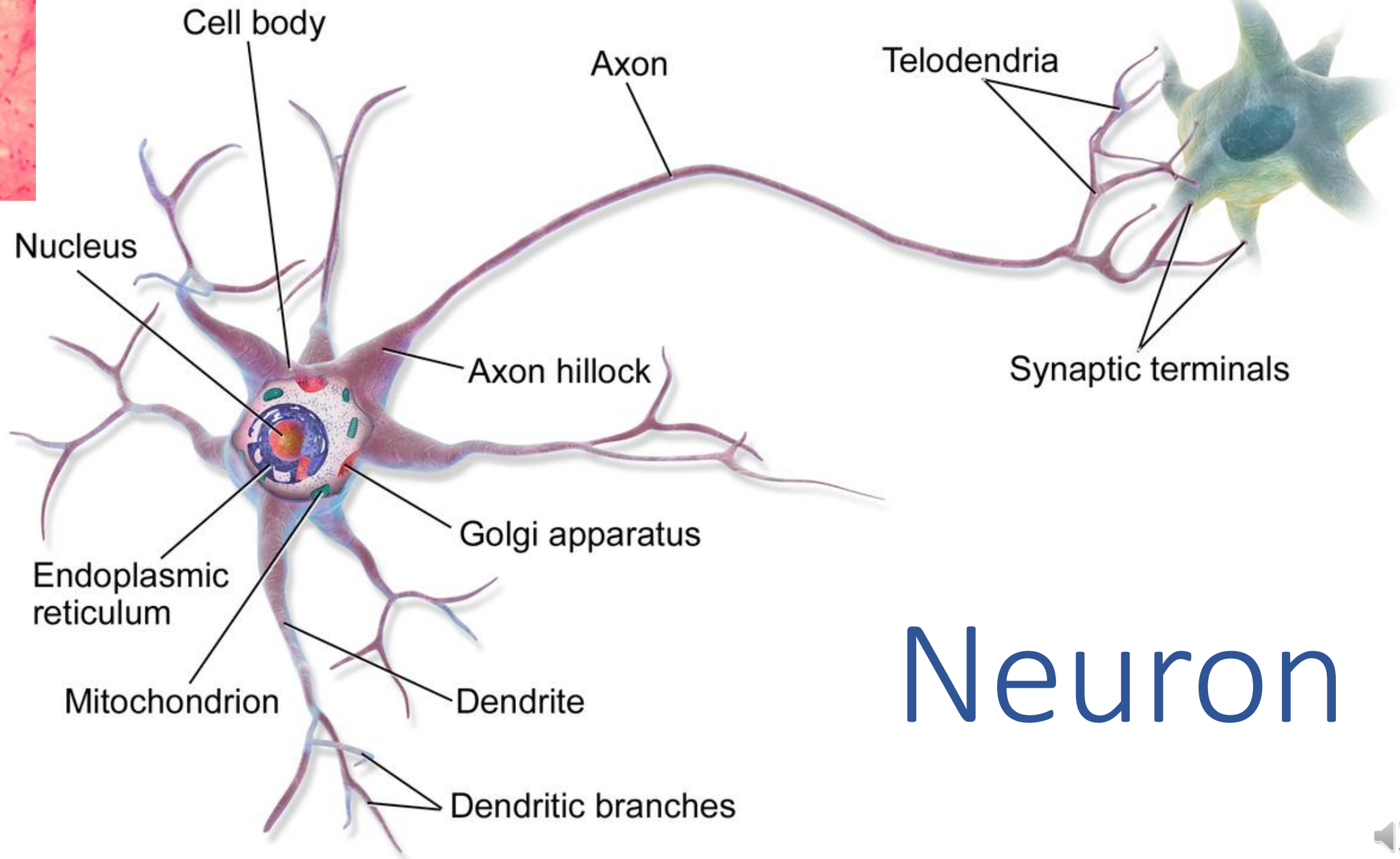
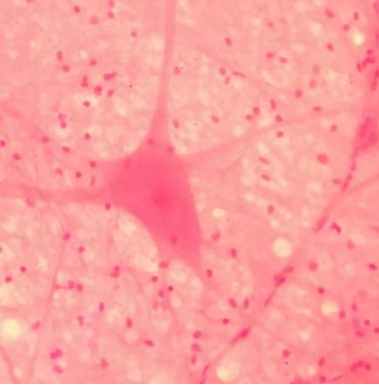
Machine
Learning

Deep
Learning



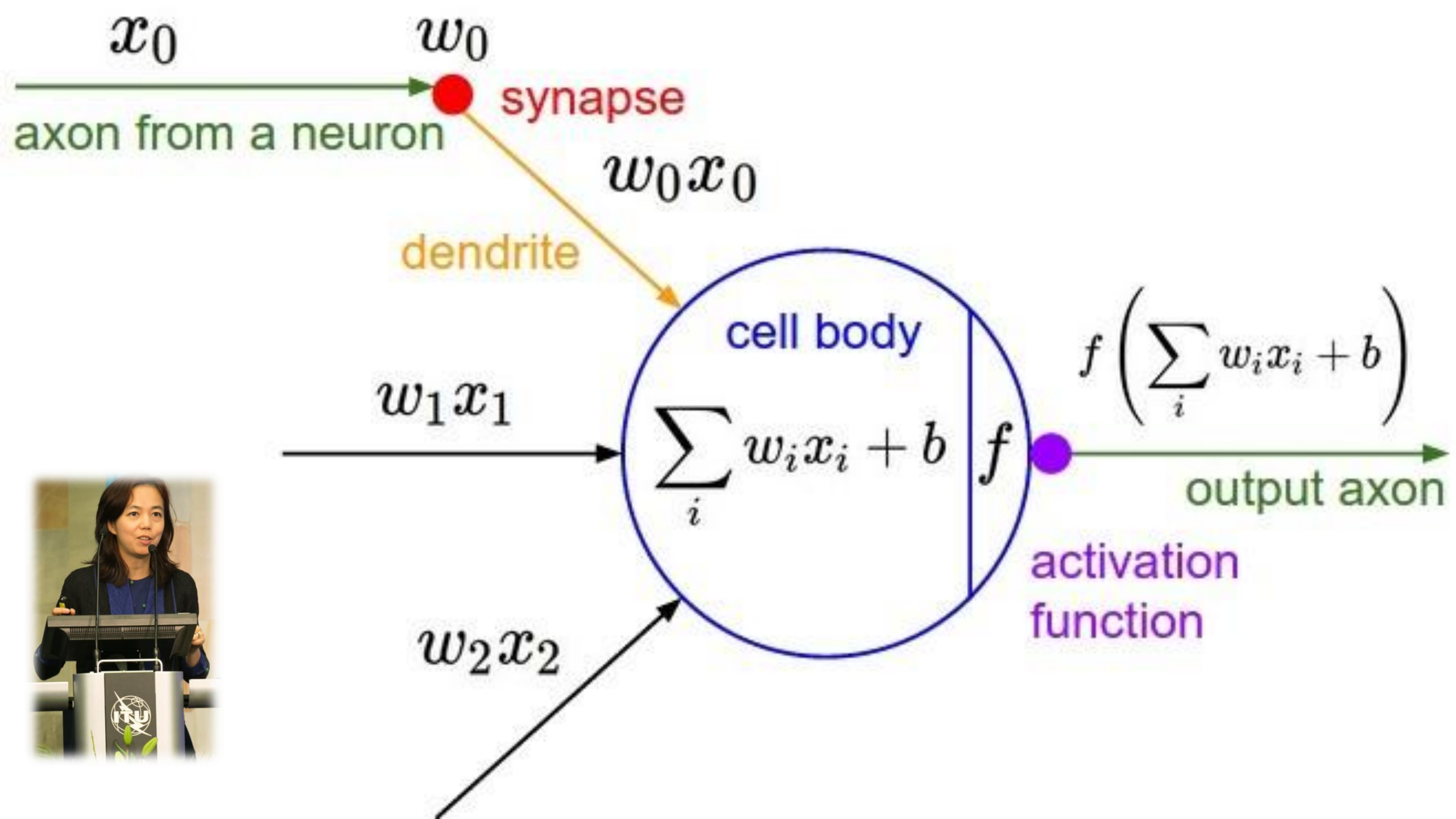
Neural Networks

A 3D visualization of a neural network. The background is a deep blue, filled with a complex web of glowing blue lines representing axons and dendrites. Several large, multi-lobed blue structures represent neurons, with smaller, similar structures scattered throughout. Interspersed among these are numerous small, bright orange-yellow spheres, some of which are larger and more prominent, suggesting active nodes or data points. The overall effect is a futuristic, digital representation of biological neural connectivity.



Neuron

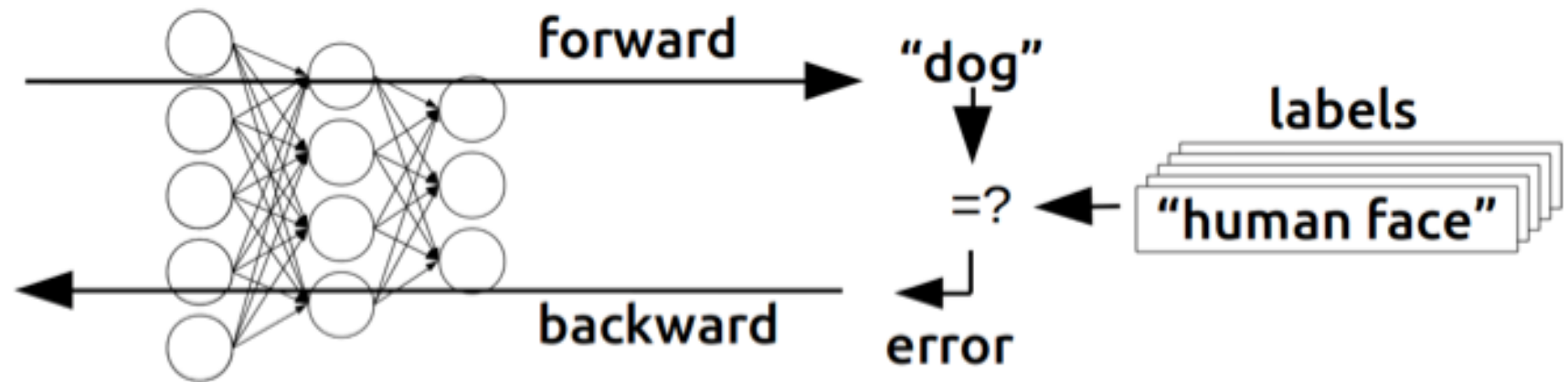
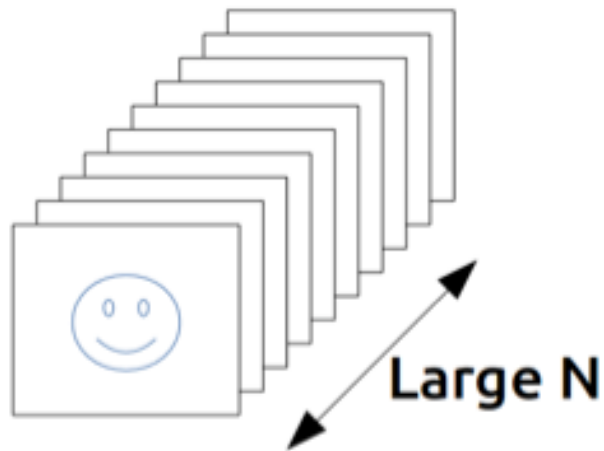




Learning (Training)

- Forward calculation + Backpropagation

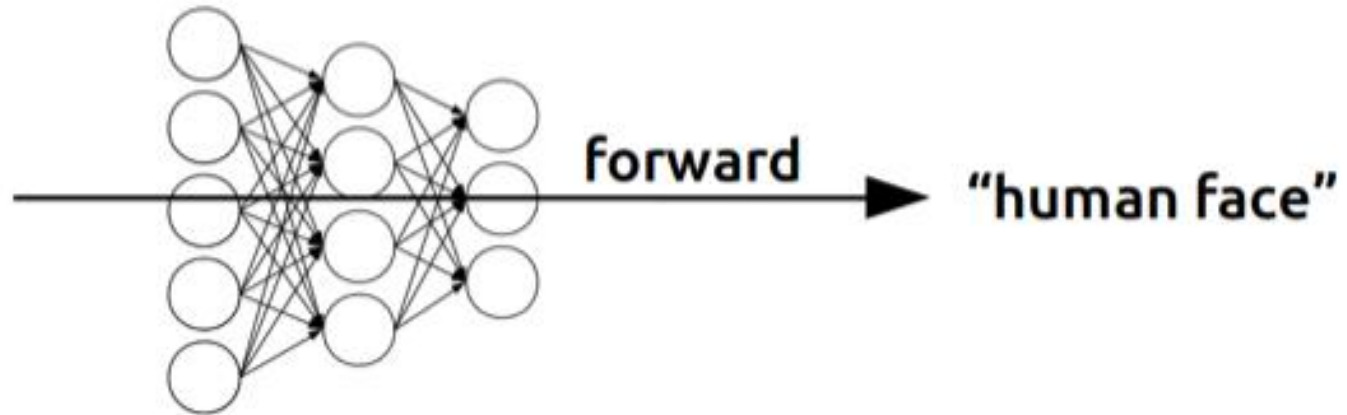
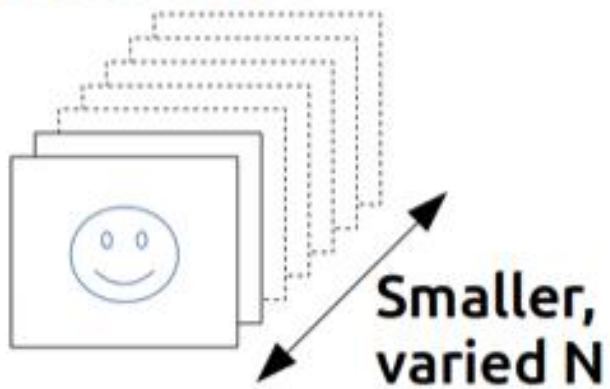
Training



Inference

- Forward calculation

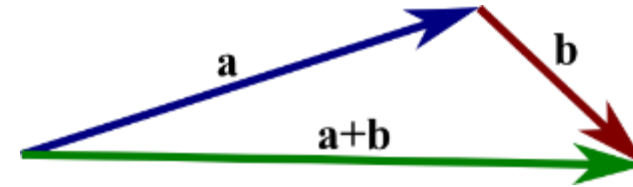
Inference



Single Variable vs. Multiple Variables

Linear Algebra

- Scalar
 - real numbers
- Vector (1D)
 - Has a magnitude & a direction
- Matrix (2D)
 - An array of numbers arranged in rows & columns
- Tensor ($\geq 3D$)
 - Multi-dimensional arrays of numbers

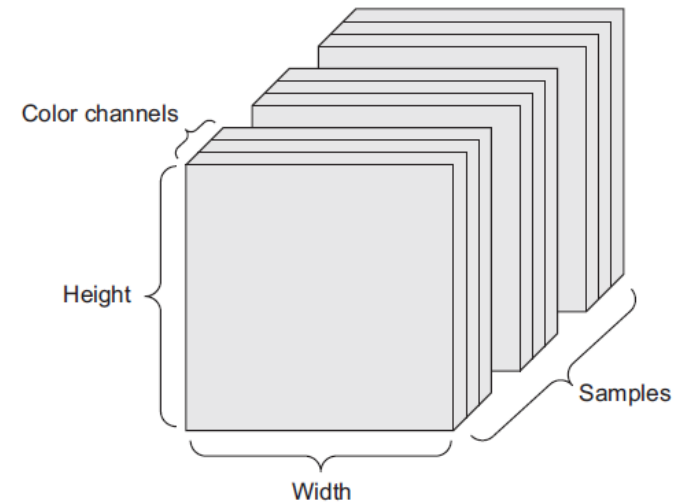
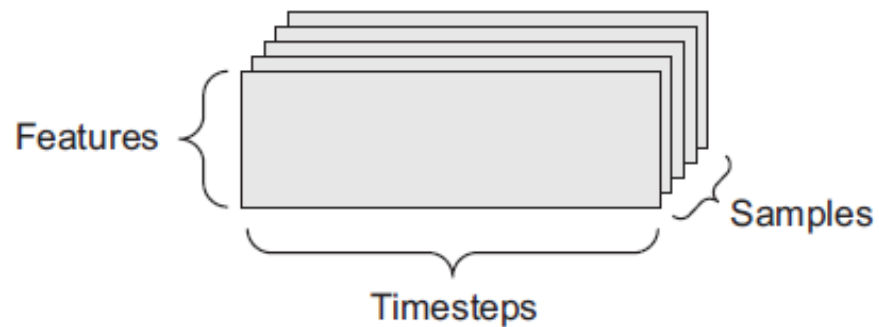


$$\mathbf{A} = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix}$$



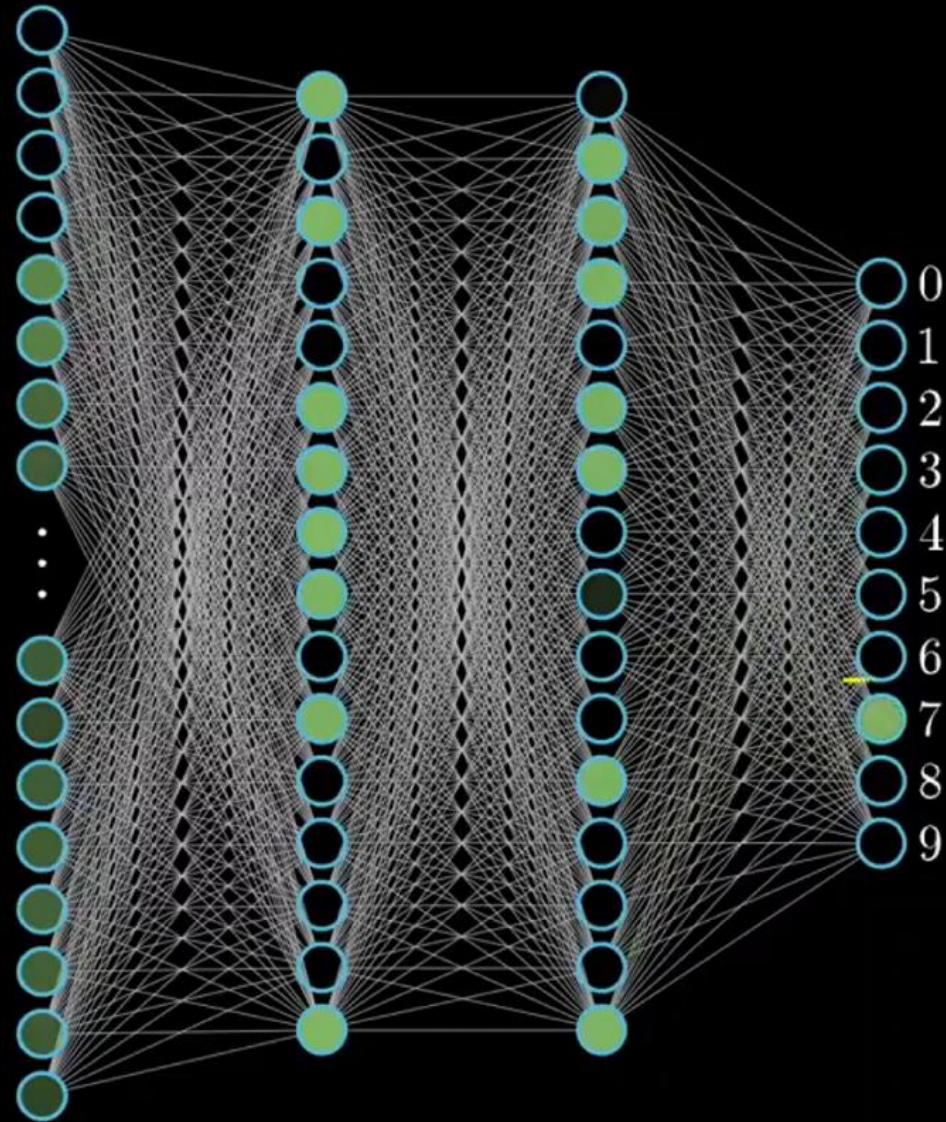
Real-world examples of Data Tensors

- Timeseries Data – 3D (samples, timesteps, features)
- Images – 4D (samples, height, width, channels)
- Video – 5D (samples, frames, height, width, channels)



Fully Connected Network

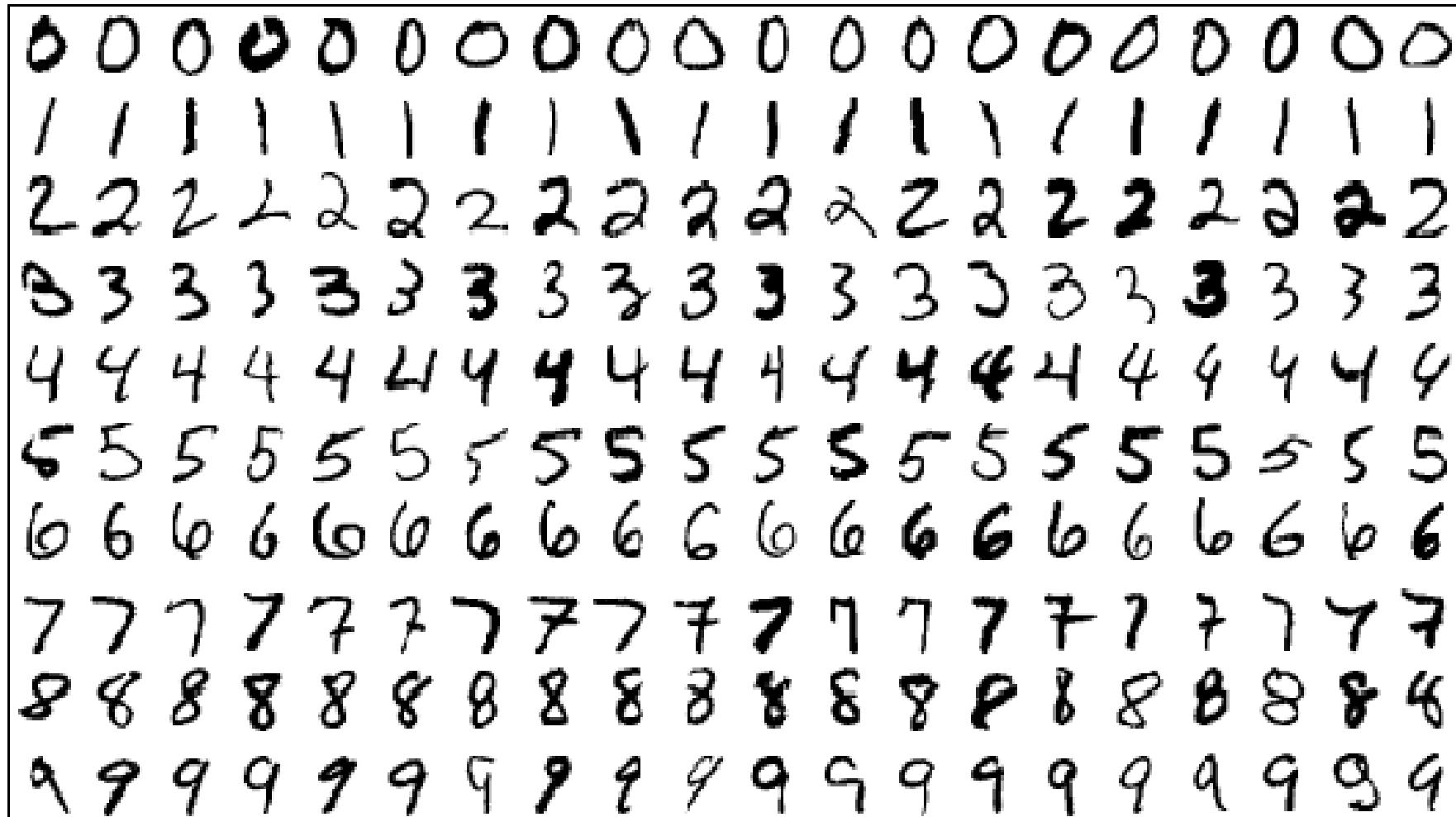
- Each neuron is connected to every neuron in the next layer

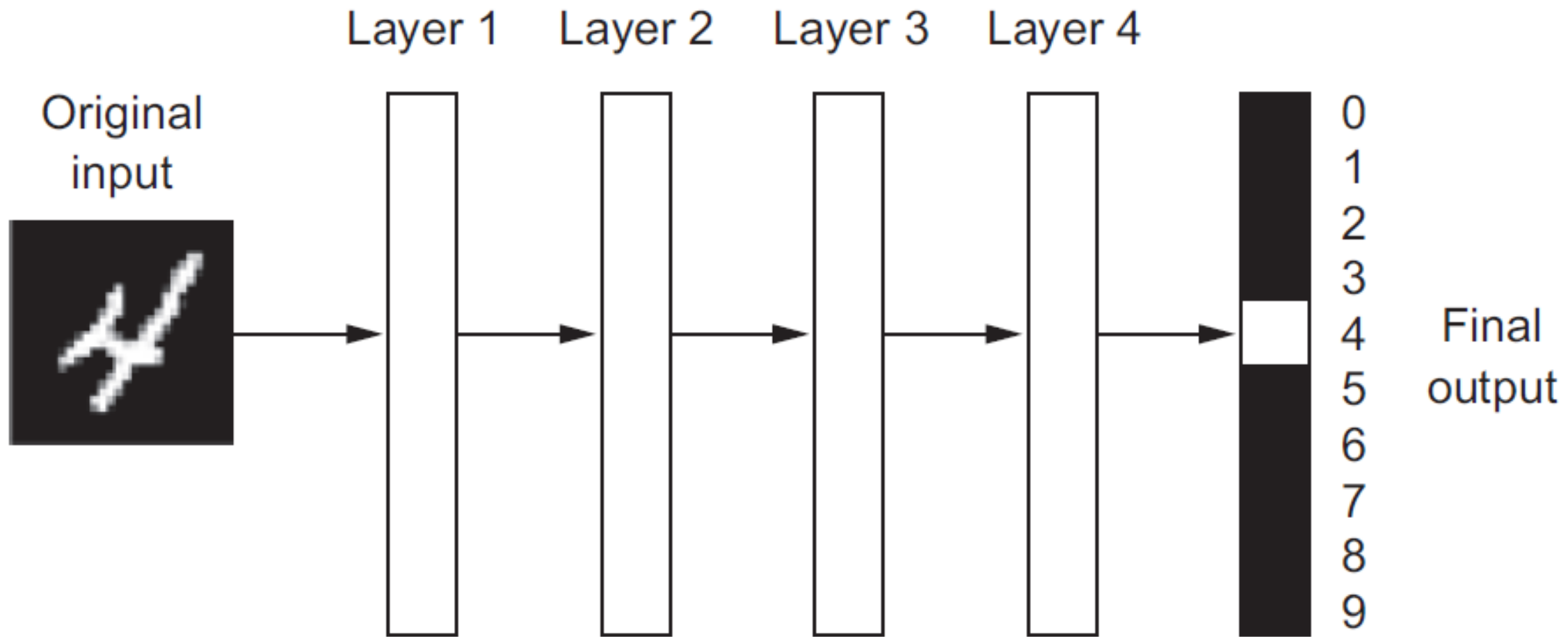


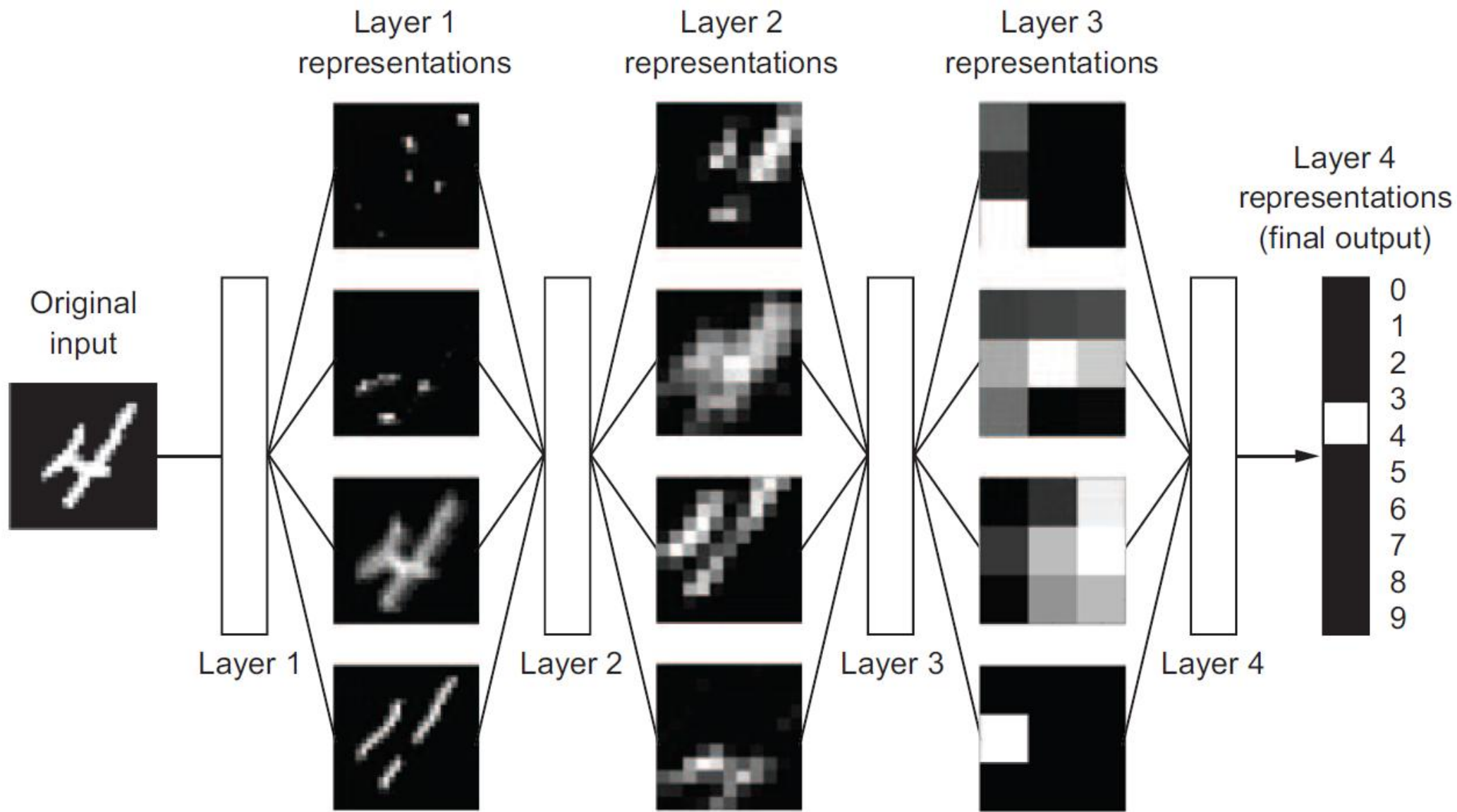
3Blue1Brown

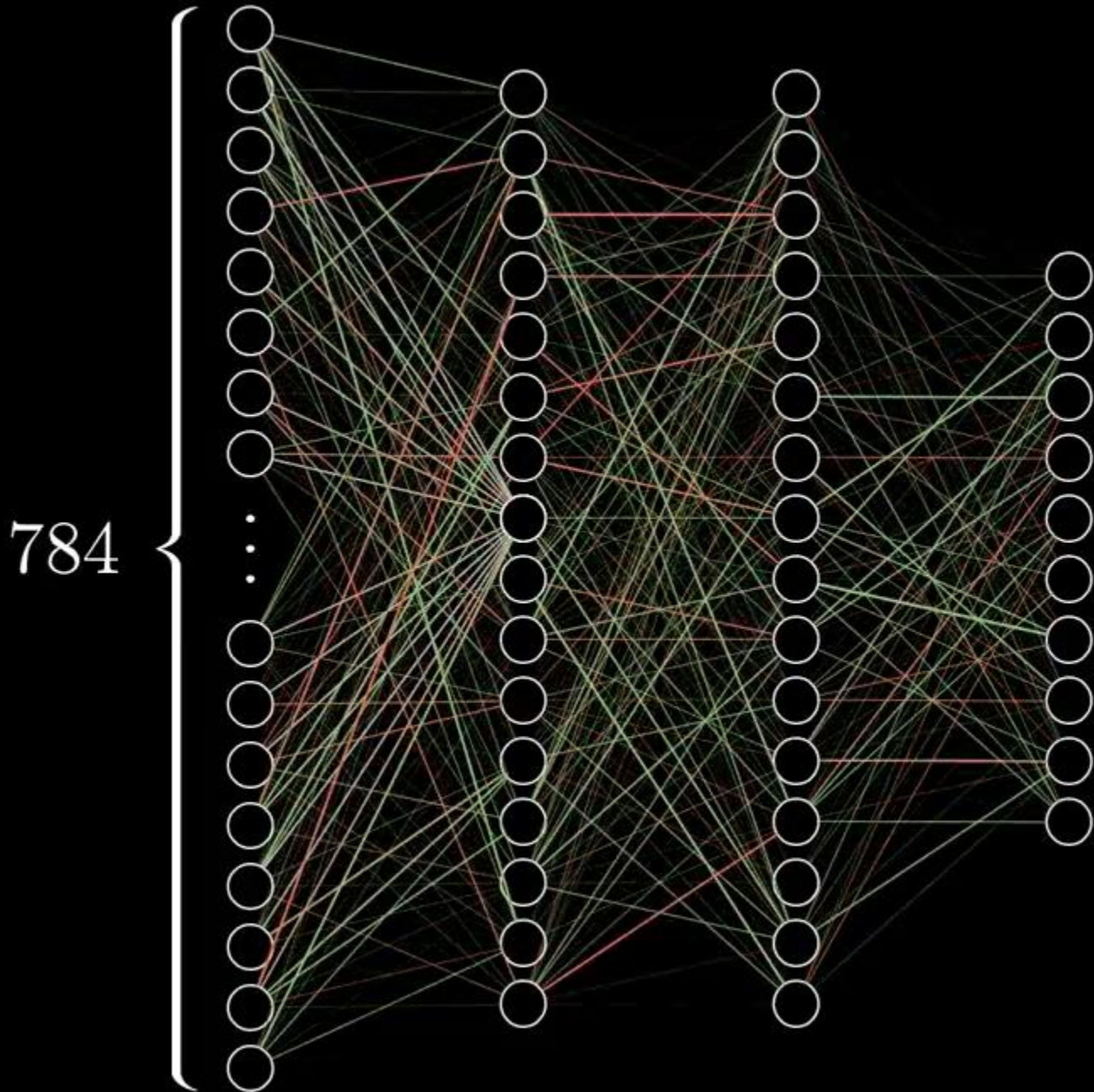
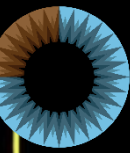
Example: Recognizing Handwritten Digits

- MNIST dataset









$$784 \times 16 + 16 \times 16 + 16 \times 10$$

weights

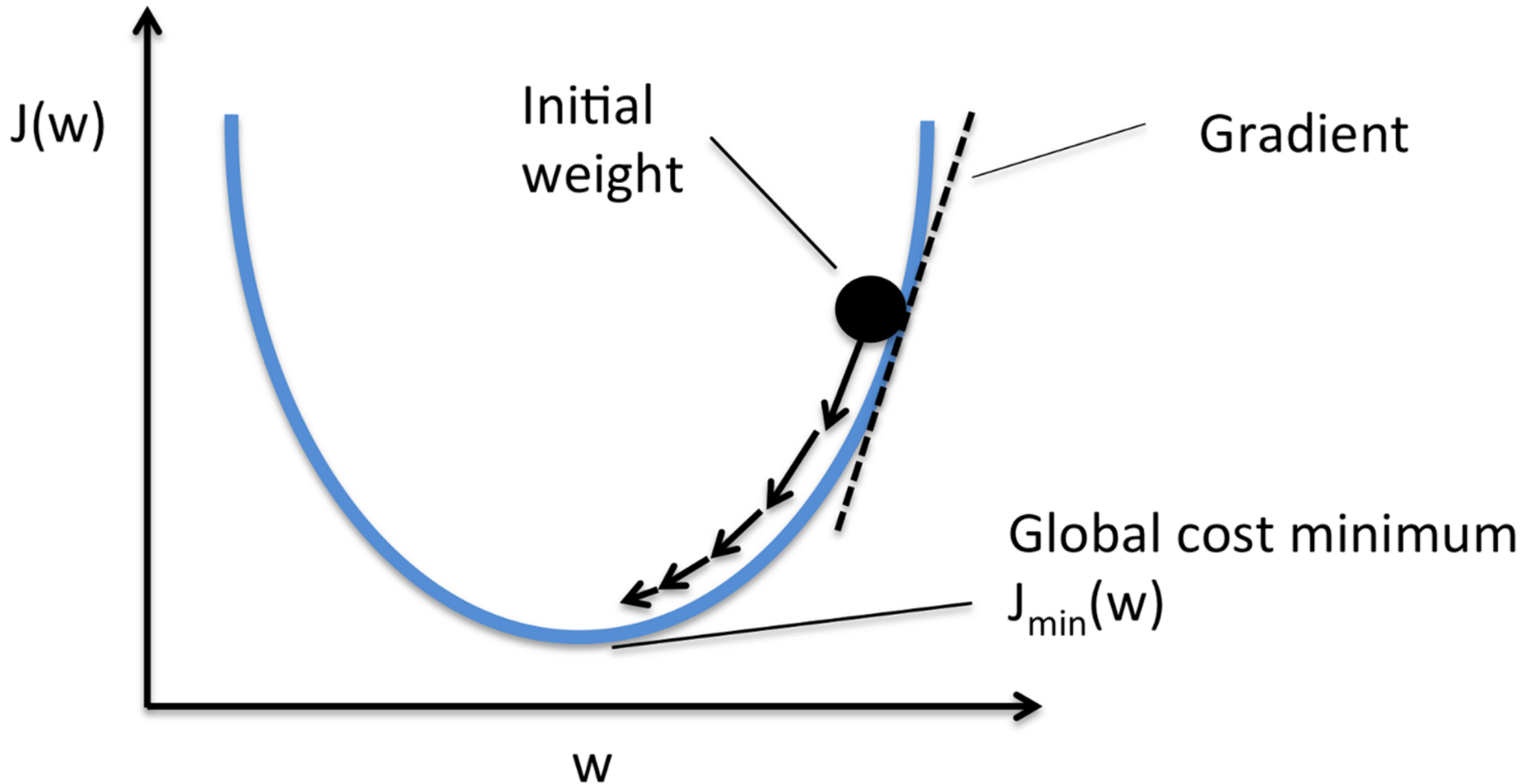
$$16 + 16 + 10$$

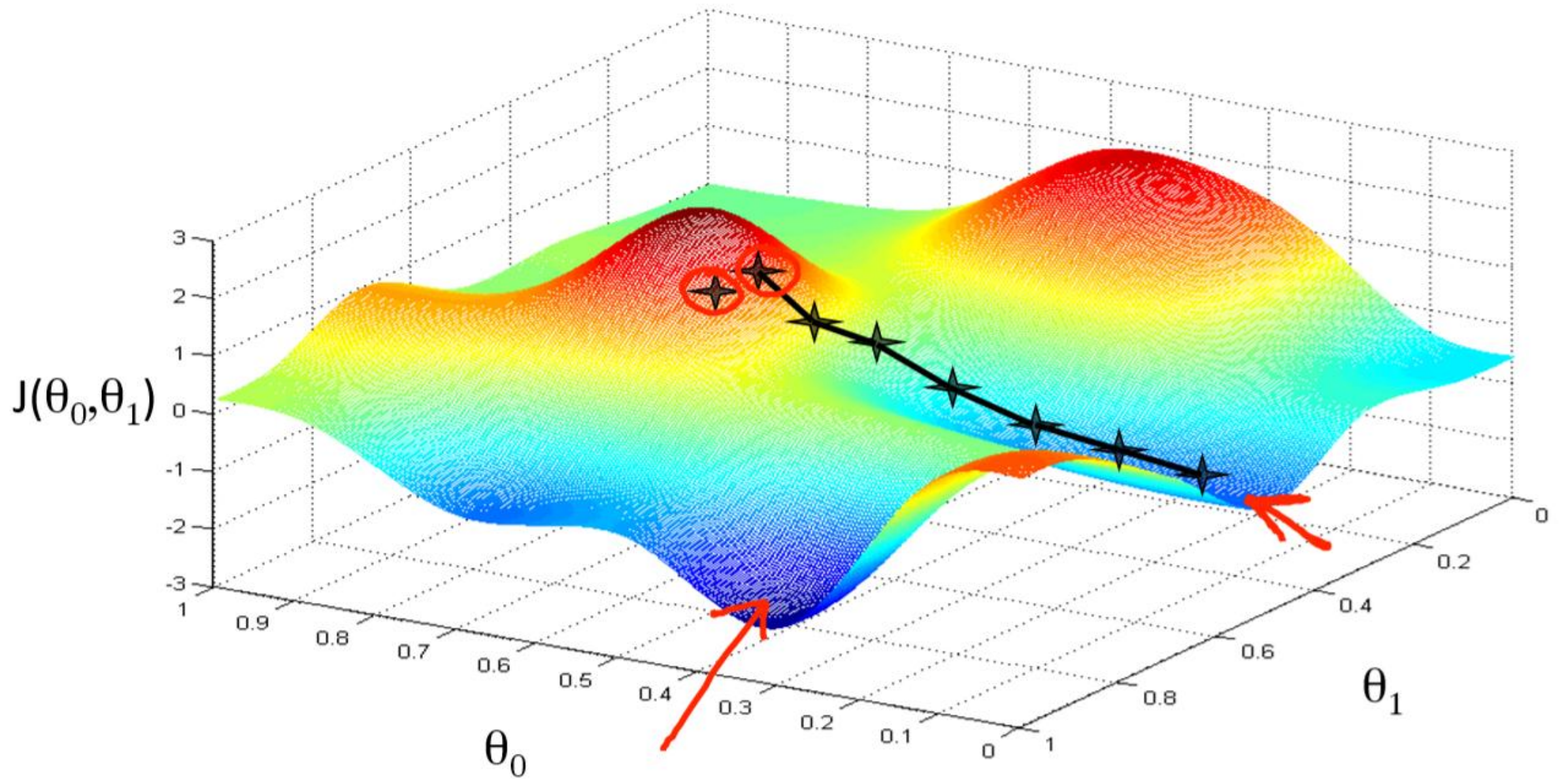
biases

13,002

Learning \rightarrow Finding the right weights and biases

Gradient Descent





<https://hackernoon.com/gradient-descent-aynk-7cbe95a778da>



Cost Function

- Mean-Squared Error

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N (f_{\theta}(x_i) - y_i)^2$$



Gradient Descent of MSE

- Gradient of MSE

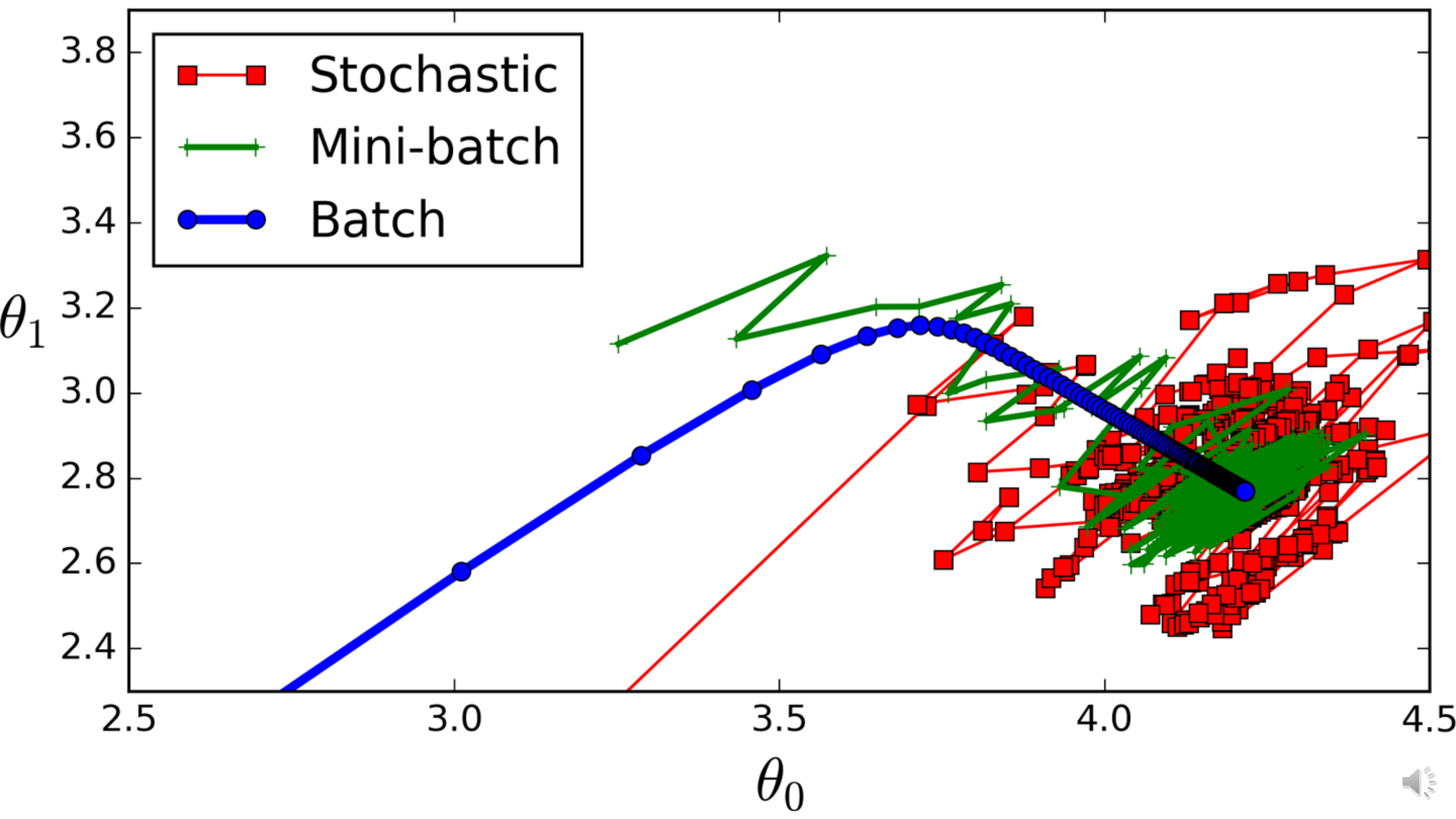
$$\frac{\partial J(\theta)}{\partial \theta} = \frac{2}{N} \sum_{i=1}^N (f_{\theta}(x_i) - y_i) f'_{\theta}(x_i)$$

- Update

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial J(\theta)}{\partial \theta_j}$$

- Repeat until Convergence





Cost function

$$C(w_1, w_2, \dots, w_{13,002})$$

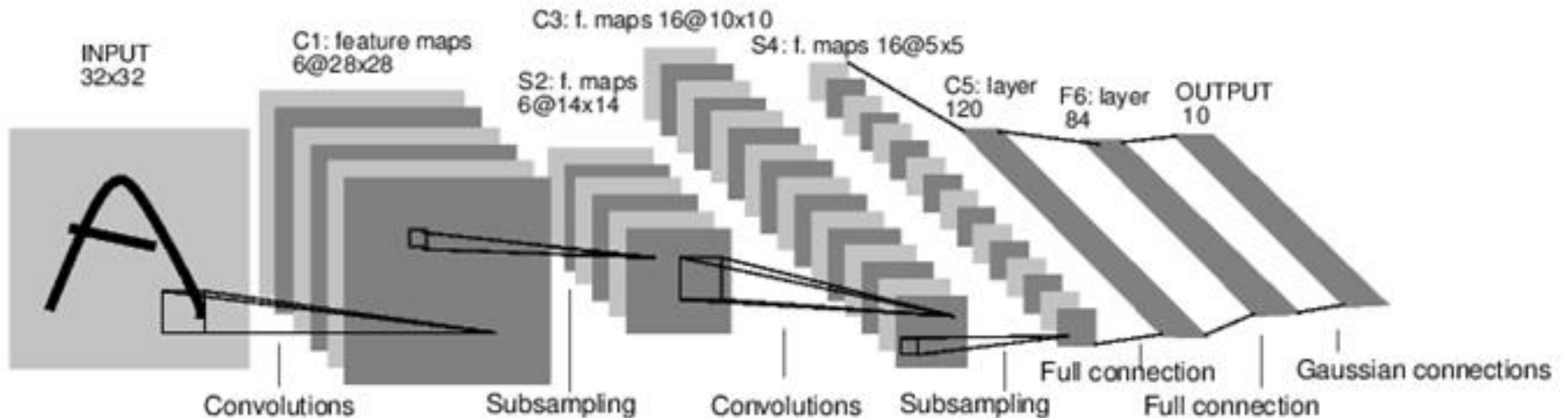
Weights and biases

Convolutional Neural Networks (CNNs)



Convolutional Neural Networks (CNNs)

- <https://medium.com/@sh.tsang/paper-brief-review-of-lenet-1-lenet-4-lenet-5-boosted-lenet-4-image-classification-1f5f809dbf17>



A Full Convolutional Neural Network (LeNet)



14,197,122 images, 21841 classes
(2021/9/21)

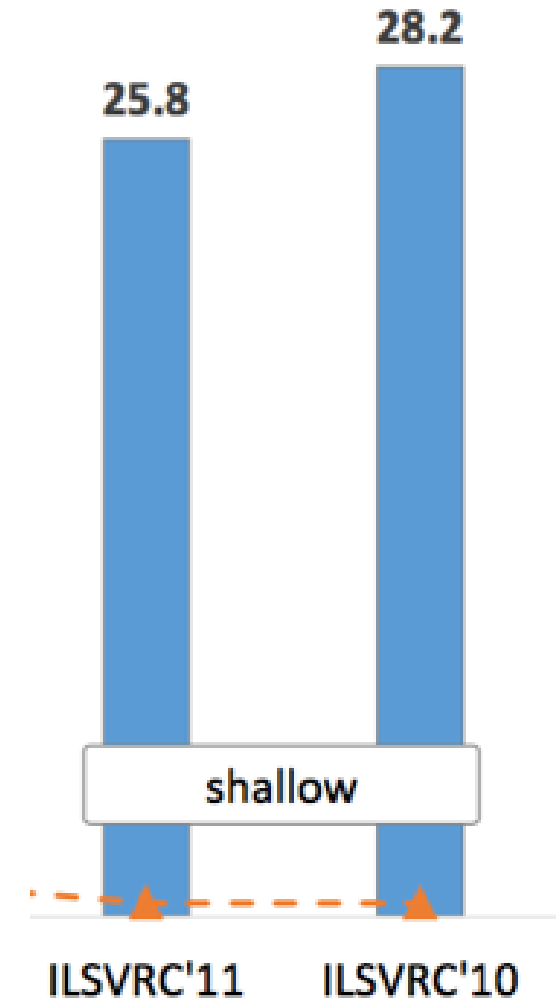


ImageNet Large Scale Visual Object Recognition Challenge (ILSVRC)

- 1000 categories
- For ILSVRC 2017
 - **Training images** for each category ranges from 732 to 1300
 - 50,000 validation **images** and 100,000 test **images**.
- Total number of images in ILSVRC 2017 is around 1,150,000

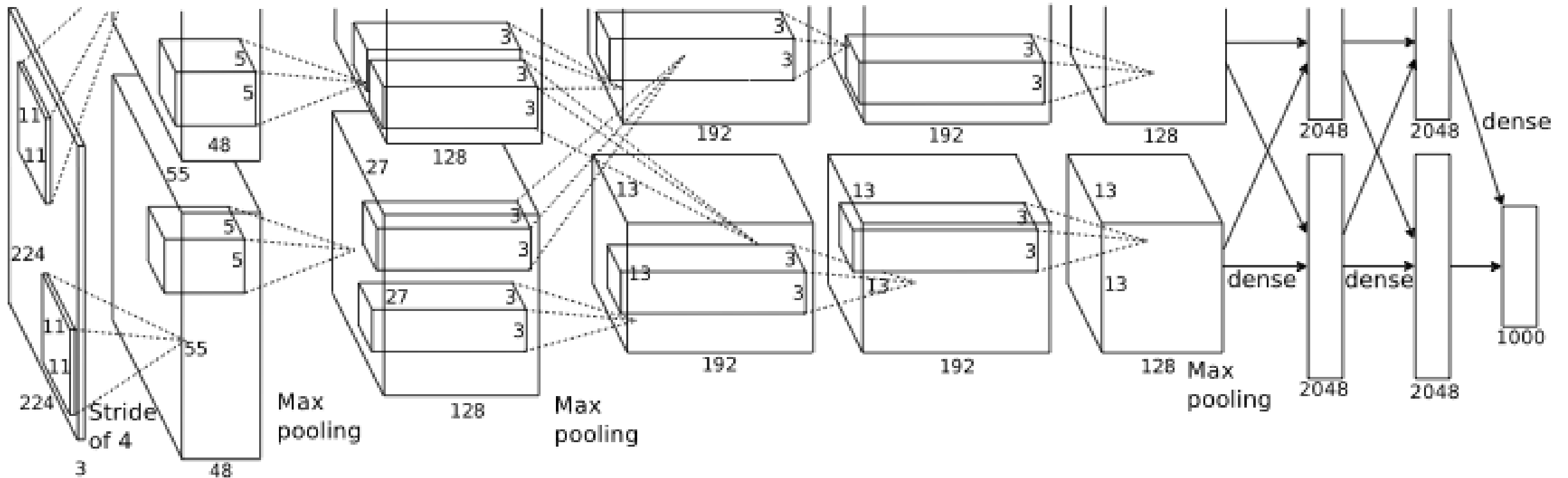


Winners' Error Rates on ImageNet Challenge

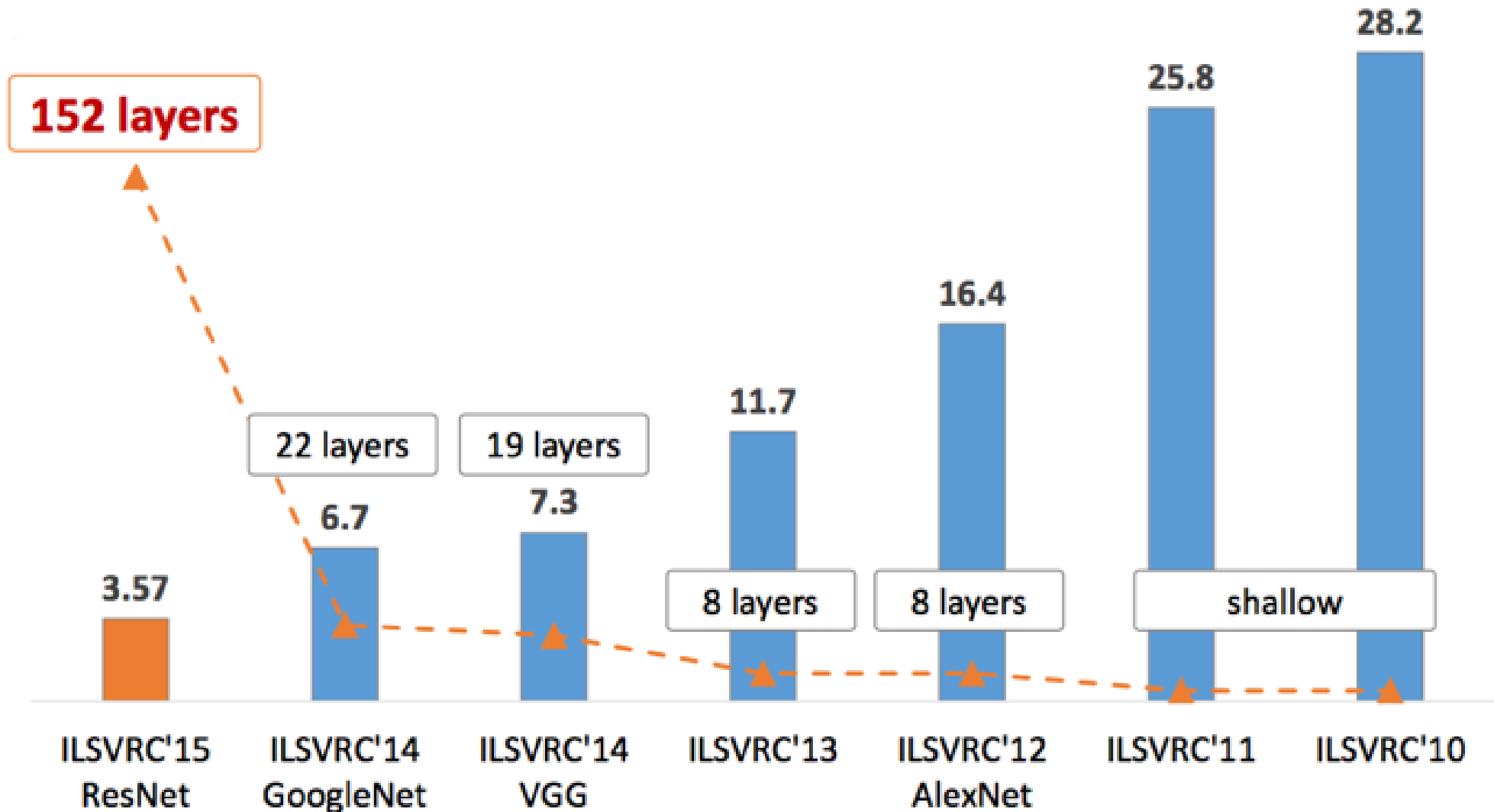


Convolutional Neural Network (AlexNet)

- Alex Krizhevsky, Geoffery Hinton et al., 2012



Winners' Error Rates on ImageNet Challenge



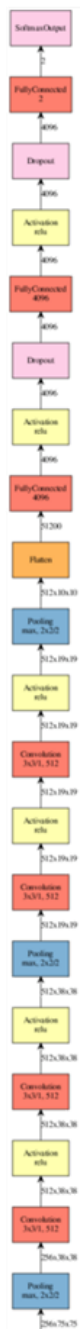
LeNet



AlexNet



VGG



GoogLeNet



Inception V3

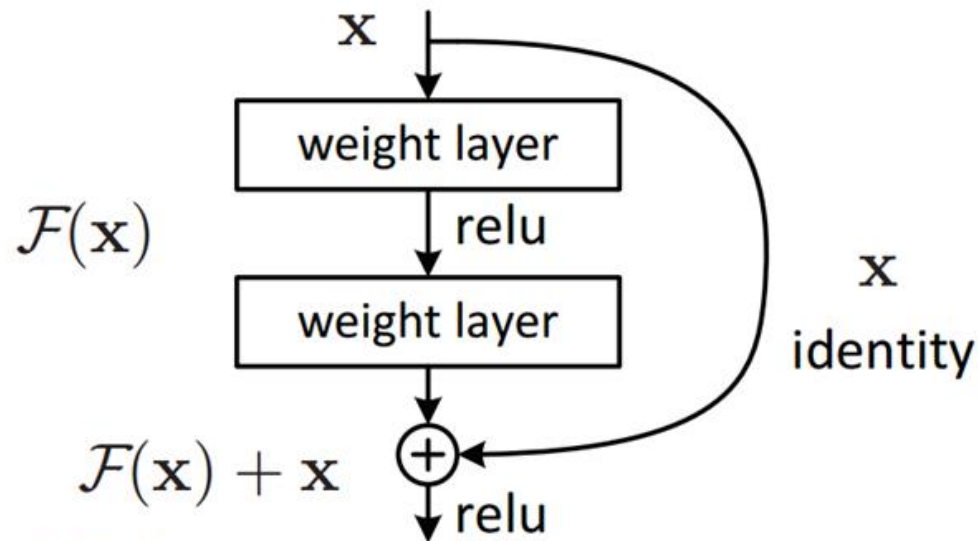


Inception BN

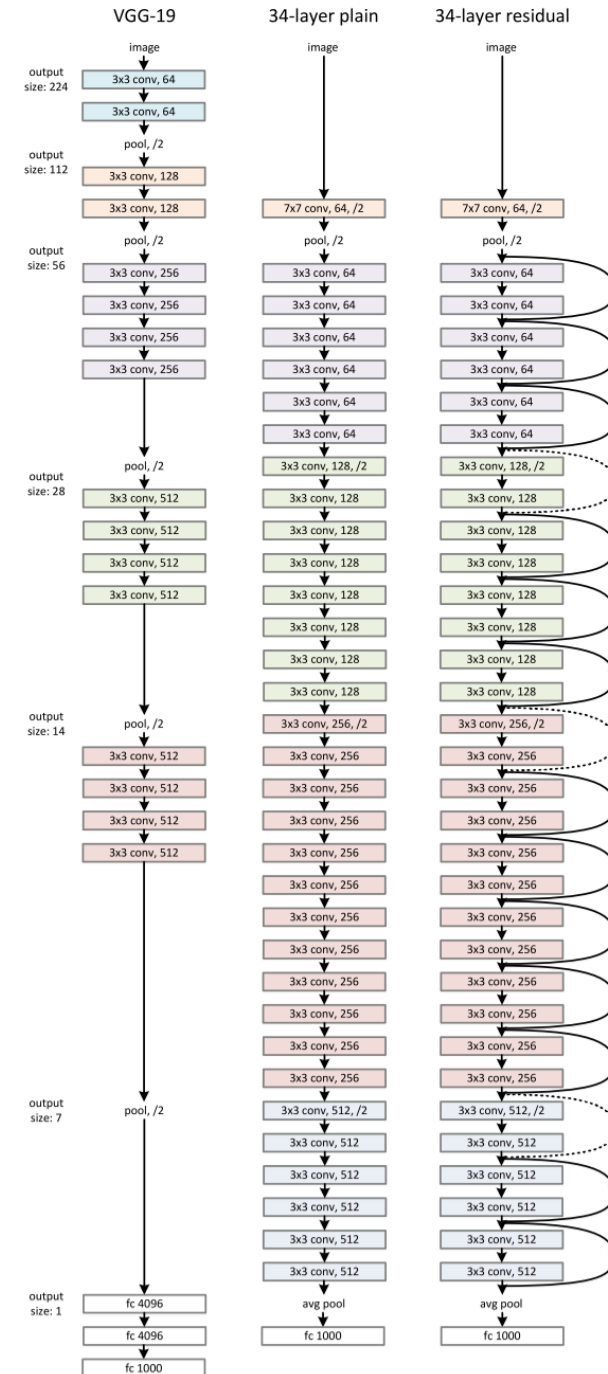


ResNet (2015)

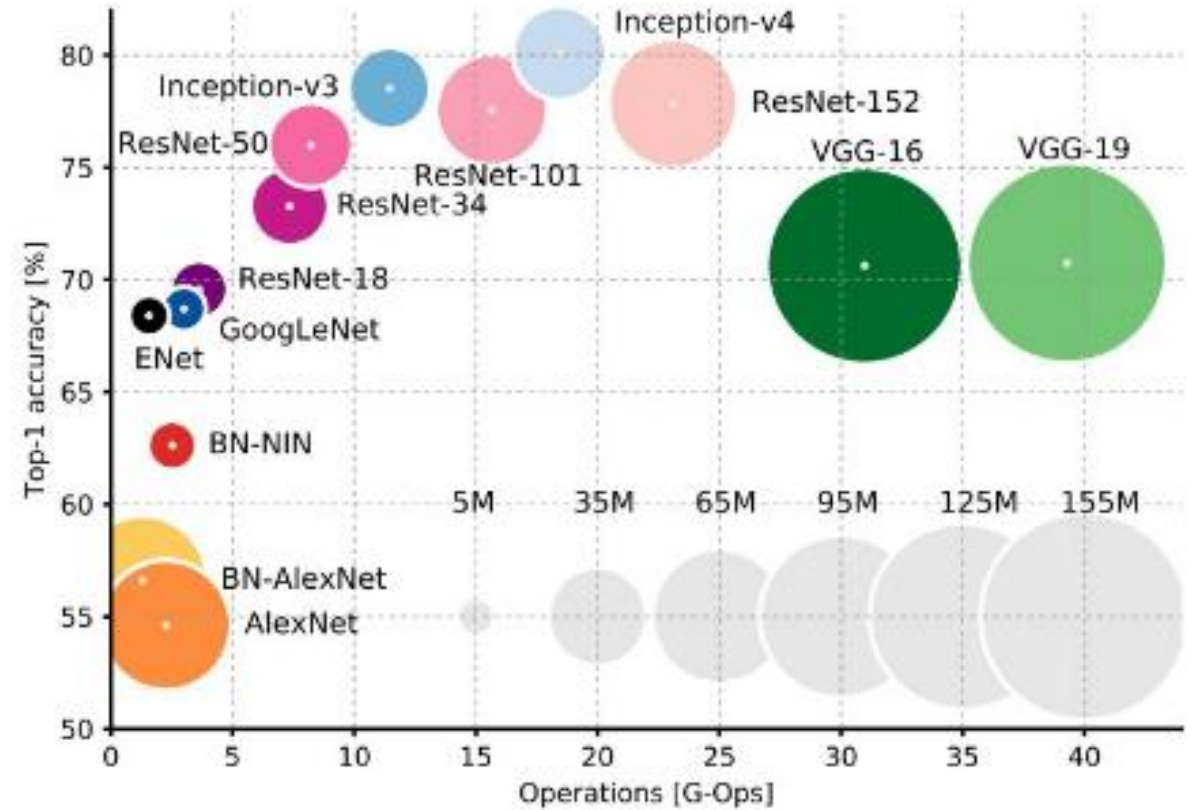
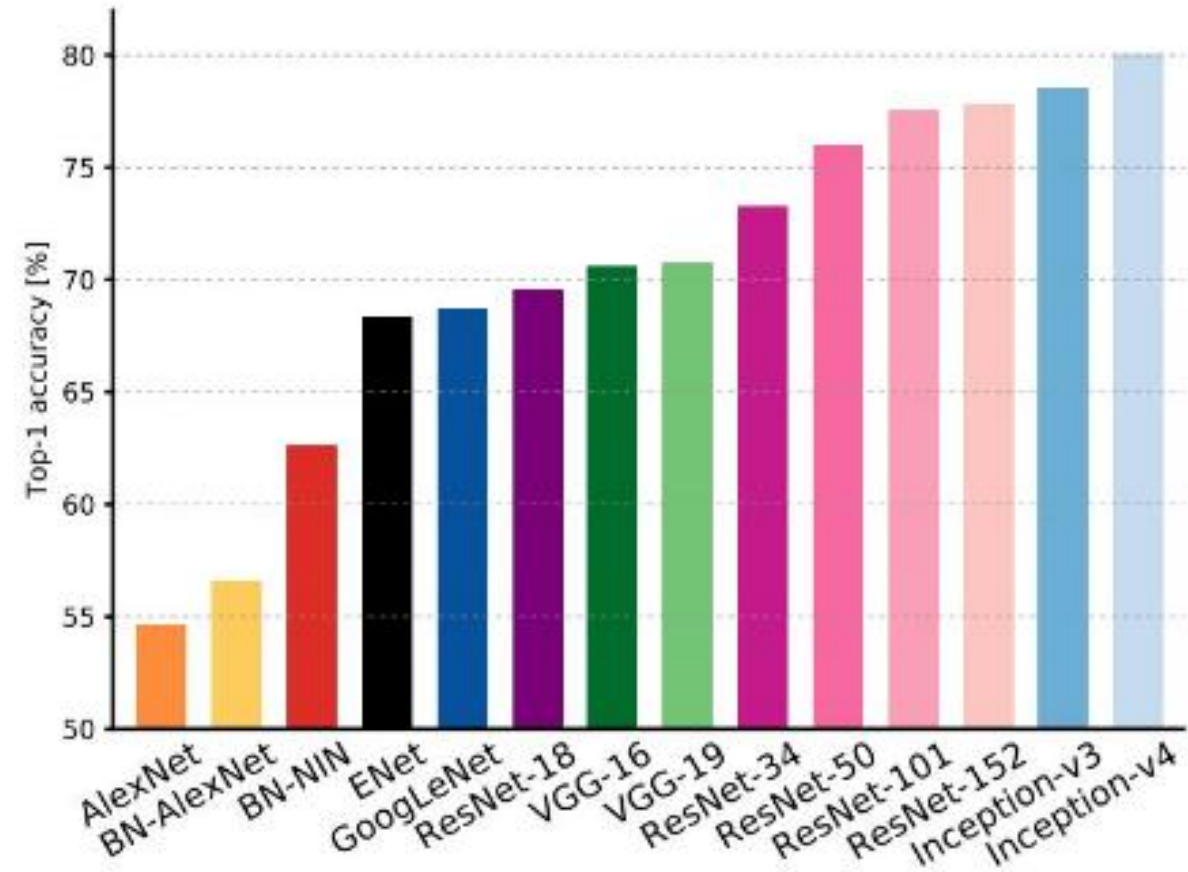
- Residual Neural Network
- Proposed “skip connection”
- 152-layer with 3.57% error rate



A residual block



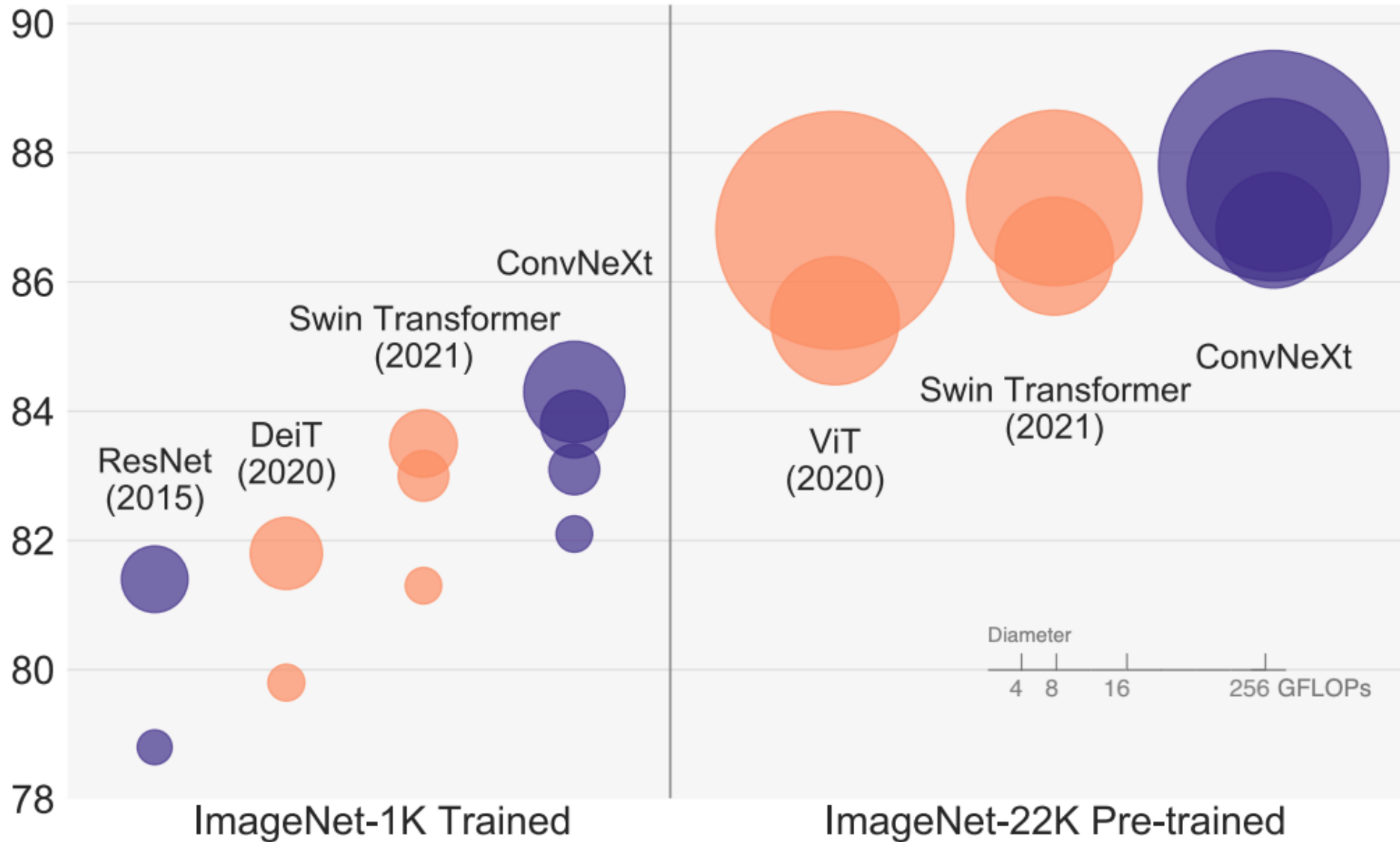
CNN Comparison



ConvNeXt (2022)

[\[2201.03545\] A ConvNet for the 2020s](#)

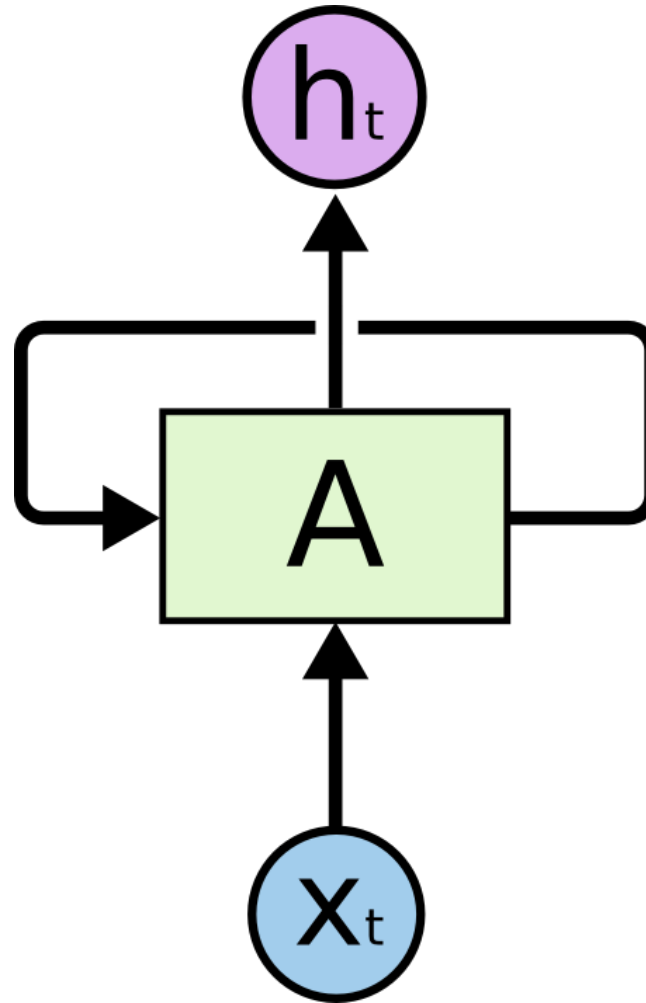
ImageNet-1K Acc.



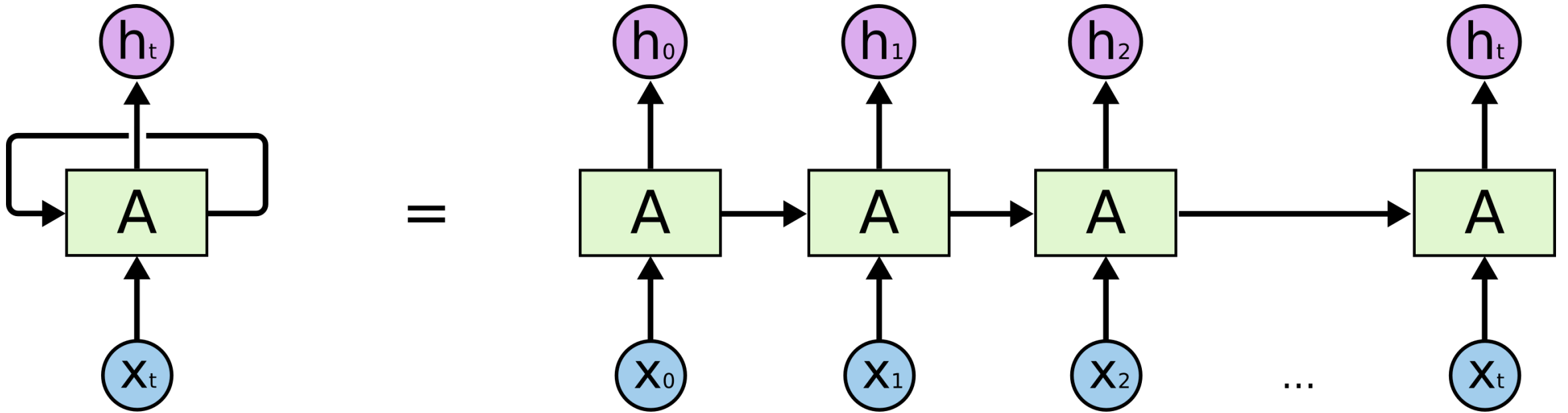
Recurrent Neural Networks (RNNs) and LSTM



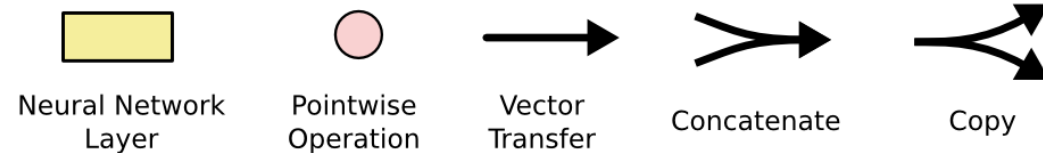
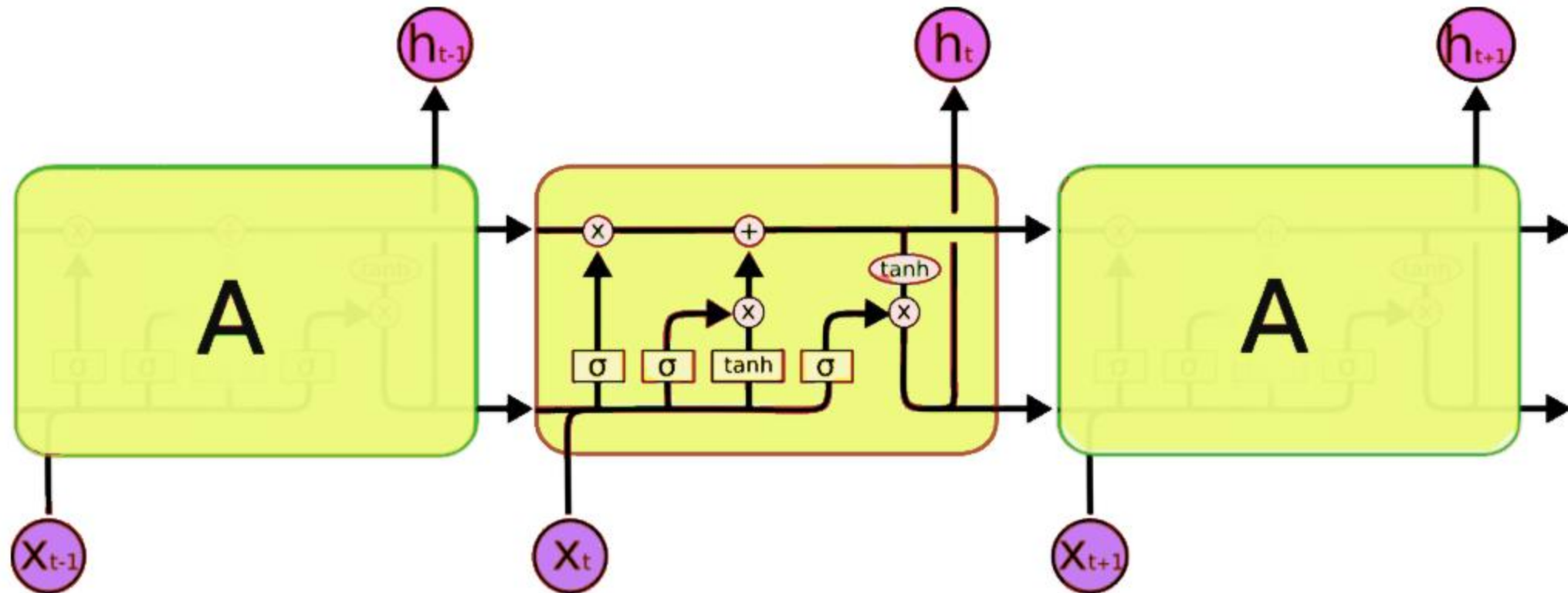
Recurrent Neural Networks (RNNs)



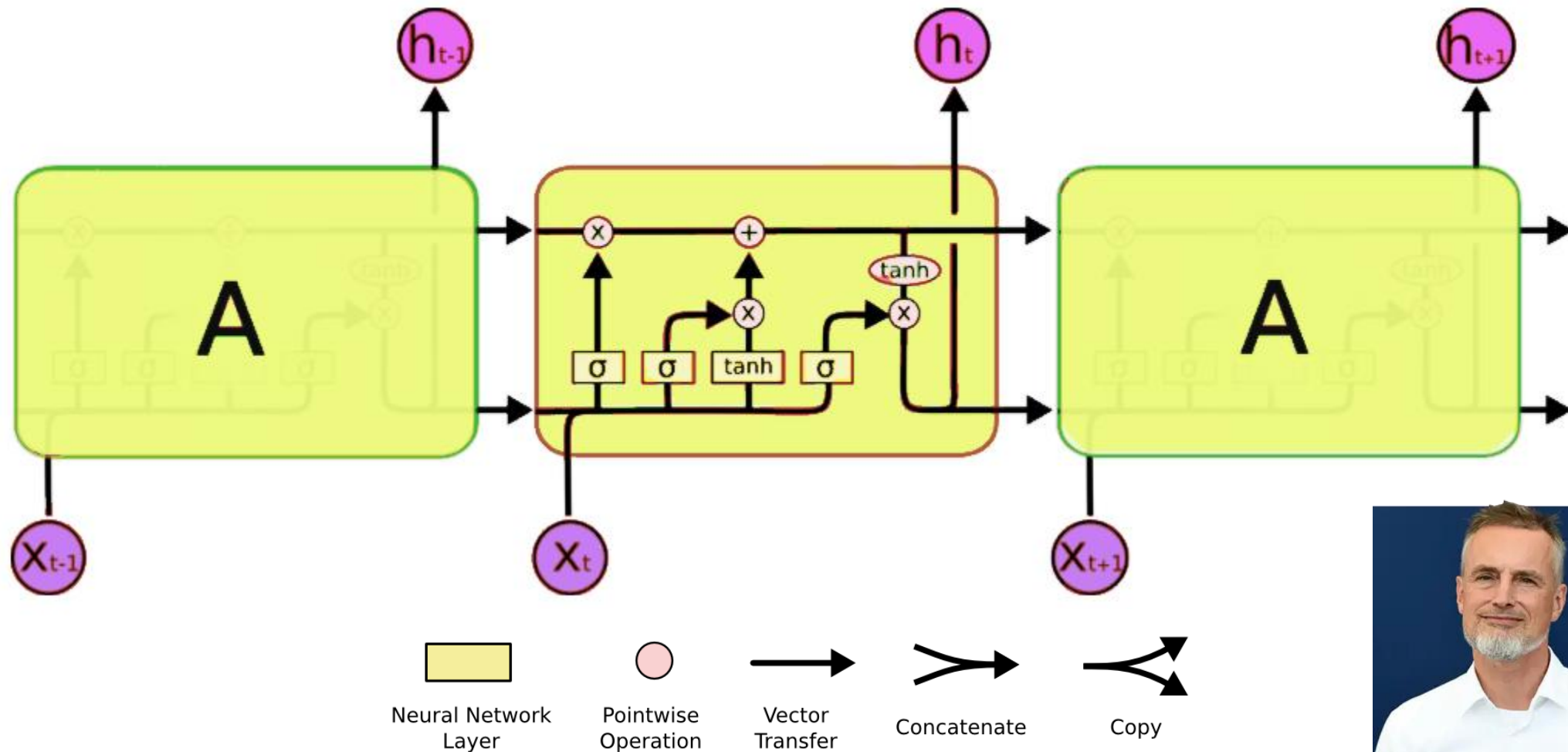
Unroll the RNN



Long Short-term Memory (LSTM)



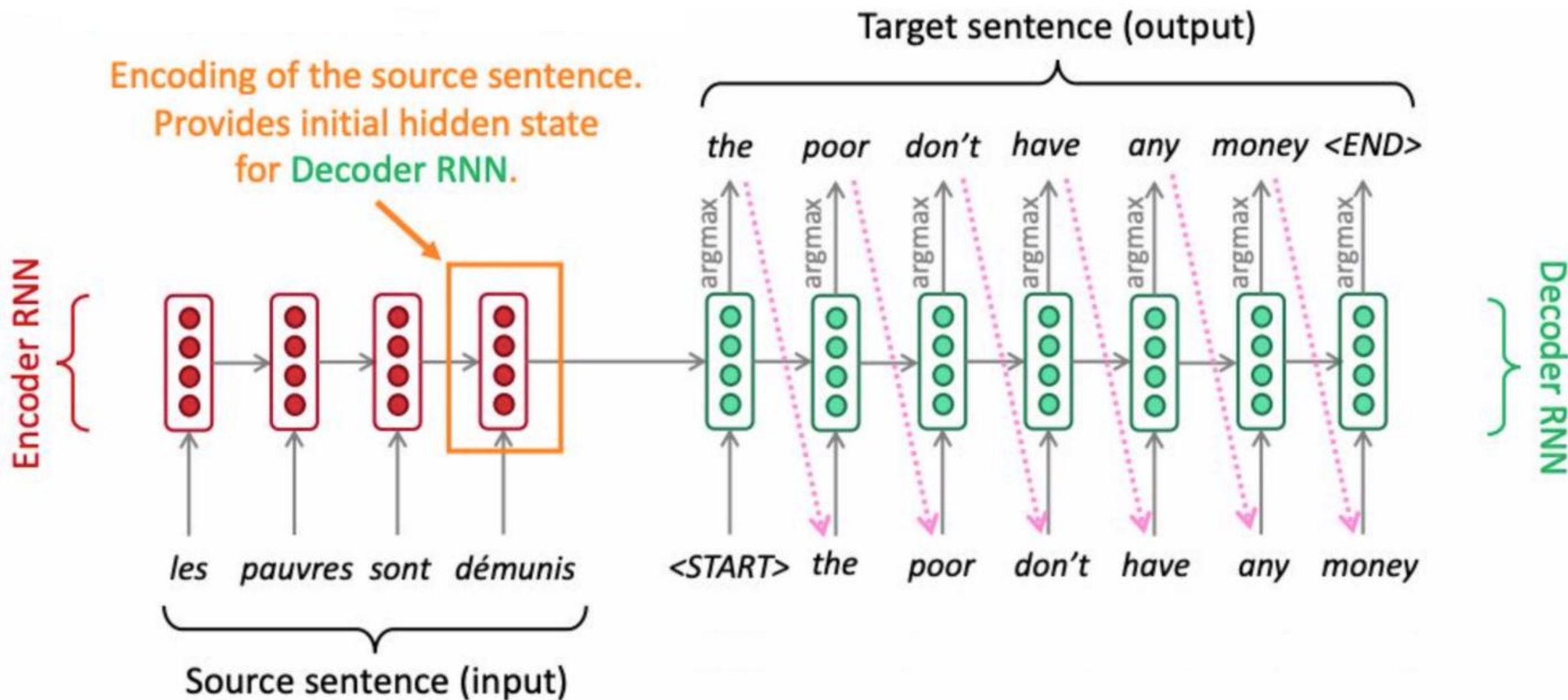
Long Short-term Memory (LSTM)



Jürgen Schmidhuber

<https://colah.github.io/posts/2015-08-Understanding-LSTMs/>

Sequence-2-Sequence model (Language Translation)





Attention is All You Need!

Ashish Vaswani*
Google Brain
avaswani@google.com

Noam Shazeer*
Google Brain
noam@google.com

Niki Parmar*
Google Research
nikip@google.com

Jakob Uszkoreit*
Google Research
usz@google.com

Llion Jones*
Google Research
llion@google.com

Aidan N. Gomez* †
University of Toronto
aidan@cs.toronto.edu

Lukasz Kaiser*
Google Brain
lukaszkaizer@google.com

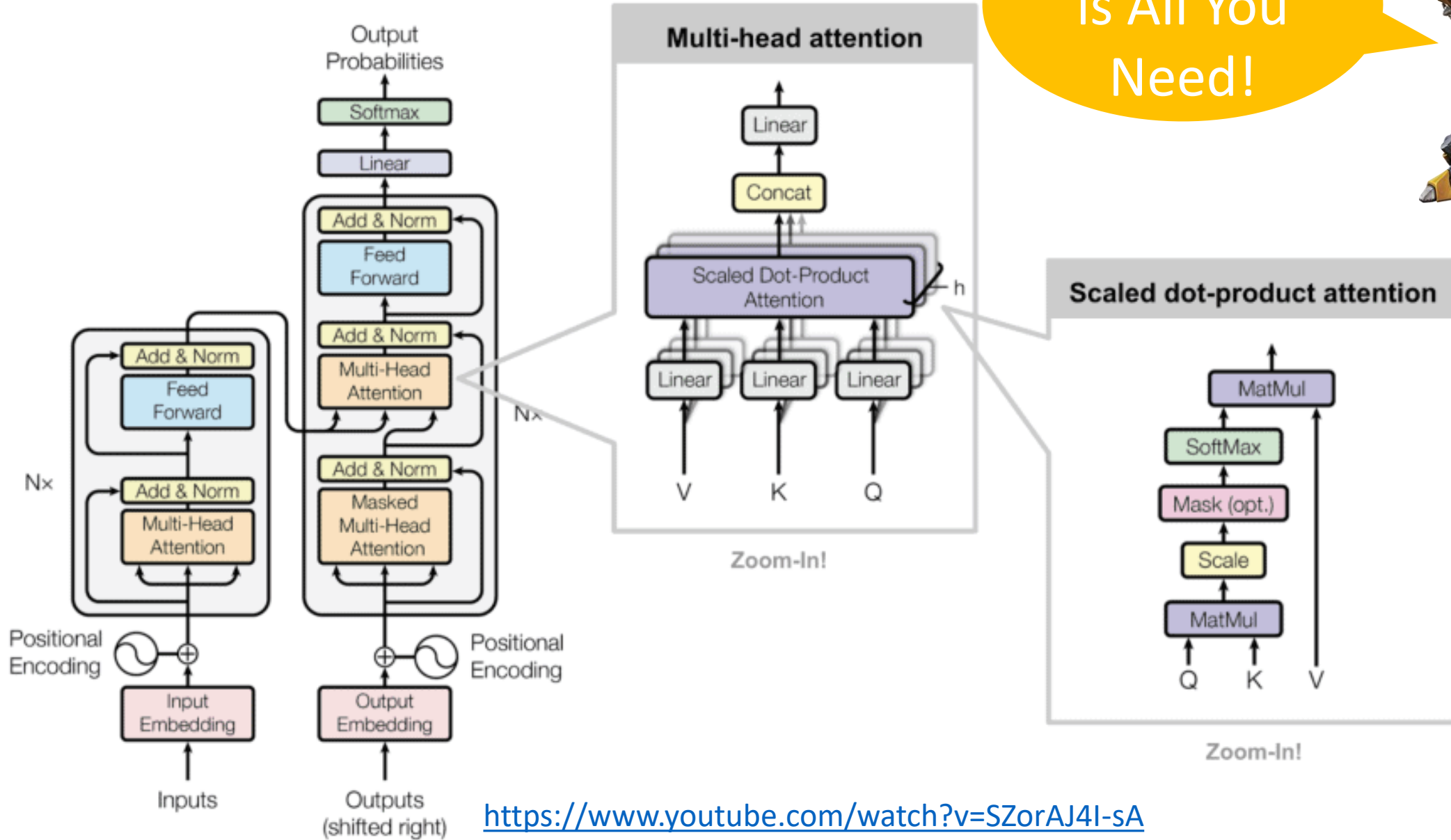
Illia Polosukhin* ‡
illia.polosukhin@gmail.com



Google Brain & University of Toronto, *NIPS*, 2017

Transformer

Attention is All You Need!

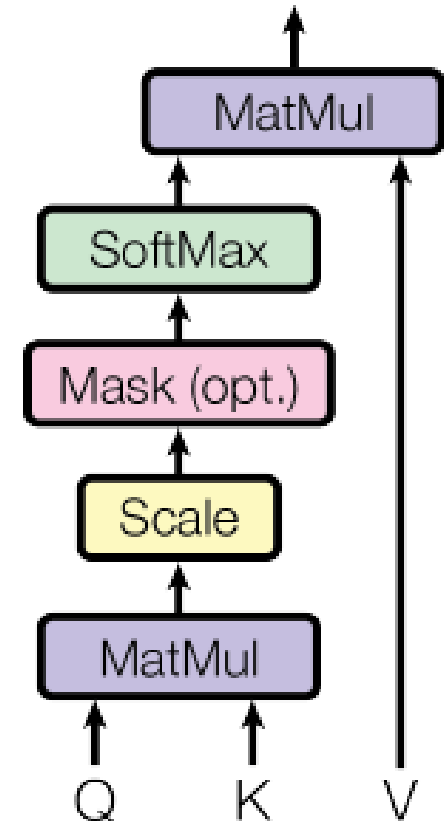


<https://www.youtube.com/watch?v=SZorAJ4I-sA>

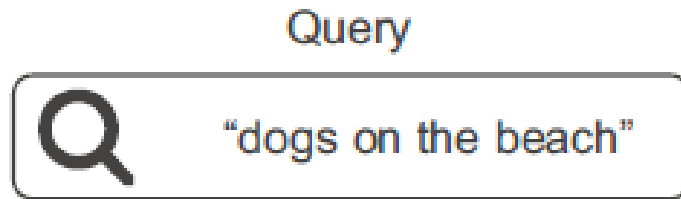
Attention Module in Transformer

- Query (Q), Key (K), Value (V) attention

$$\text{Attention}(Q, K, V) = \text{softmax}\left(\frac{QK^T}{\sqrt{d_k}}\right)V$$



Query, Keys, Values



Retrieving
images from
a database

Keys

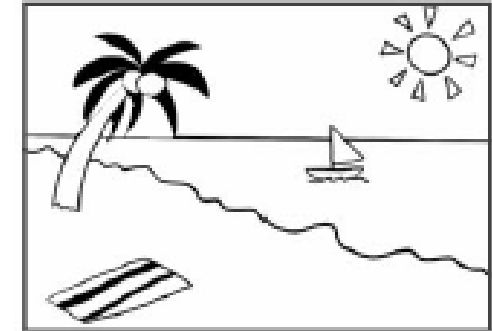
Values

match: 0.5

Beach

Tree

Boat

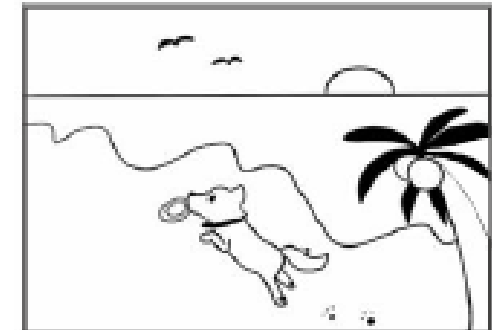


match: 1.0

Beach

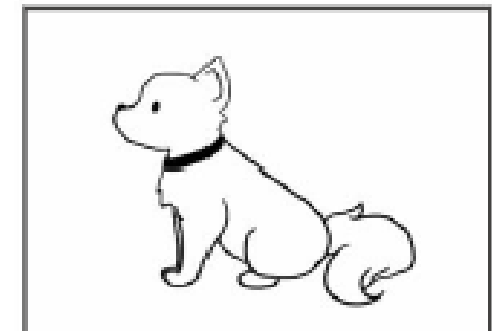
Dog

Tree



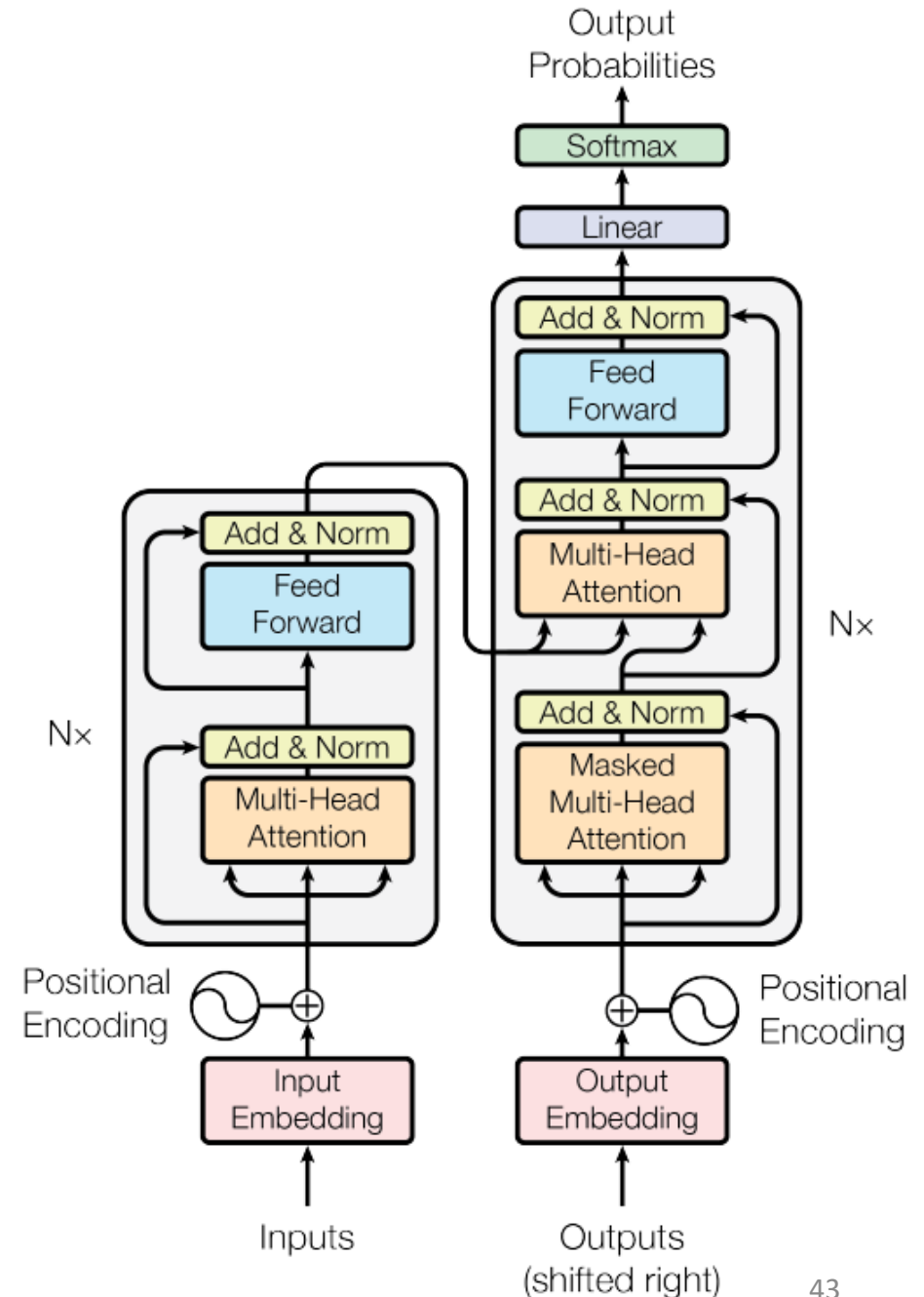
match: 0.5

Dog



The Transformer Model

- Encoder-decoder architecture
- Multi-head attention
 - Self-attention in encoders
 - Masked Self-attention in decoders
 - Encoder-decoder attention
- Positional encoding



Visualizing Attention

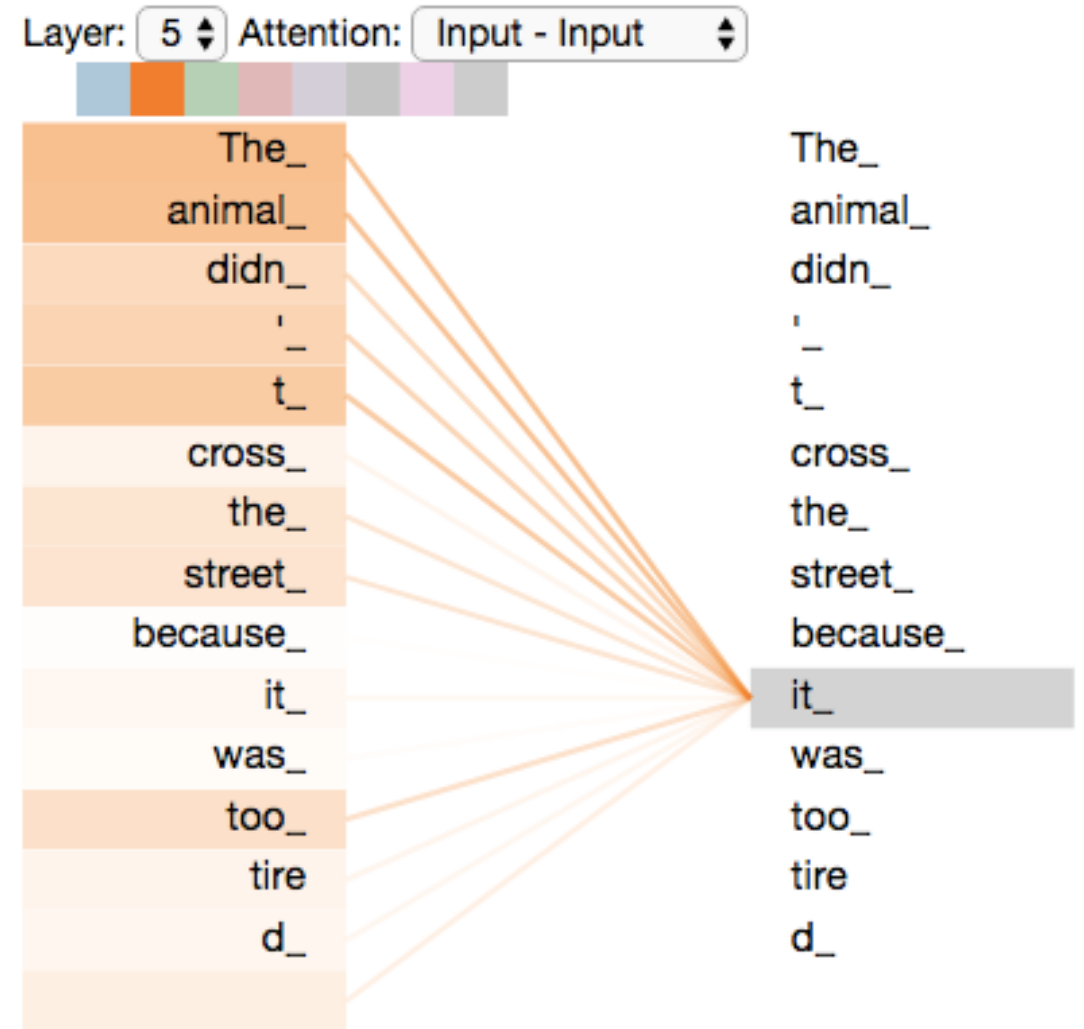
- Tensor2Tensor Notebook

https://colab.research.google.com/github/tensorflow/tensor2tensor/blob/master/tensor2tensor/notebooks/hello_t2t.ipynb

Inputs: The animal didn't cross the street because it was too tired



Outputs: Das Tier überquerte die Straße nicht, weil es zu müde war, weil es zu müde war.



Deep Reinforcement Learning



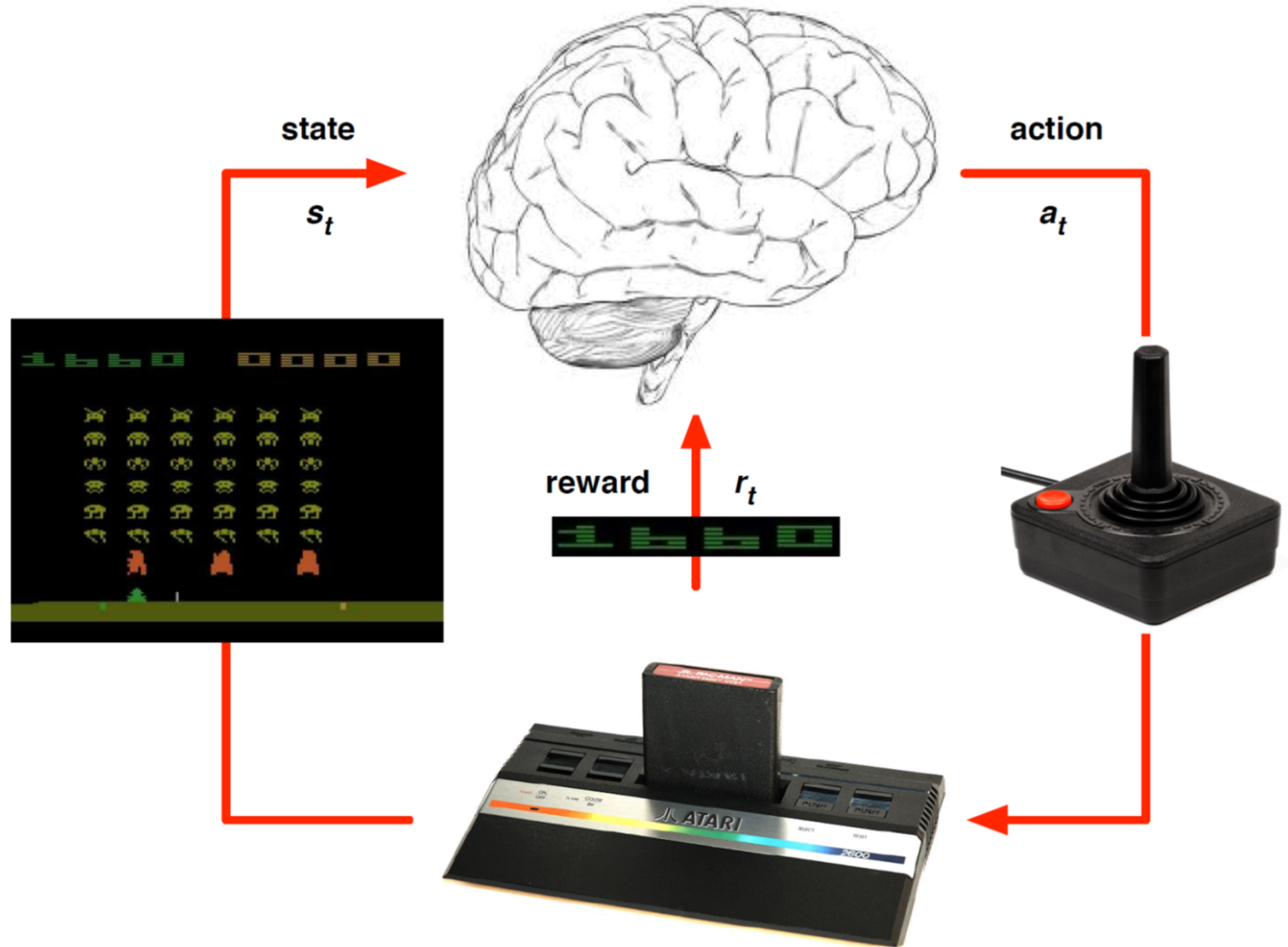
DeepMind: DRL in Atari



Demis Hassabis



Mustafa Suleyman

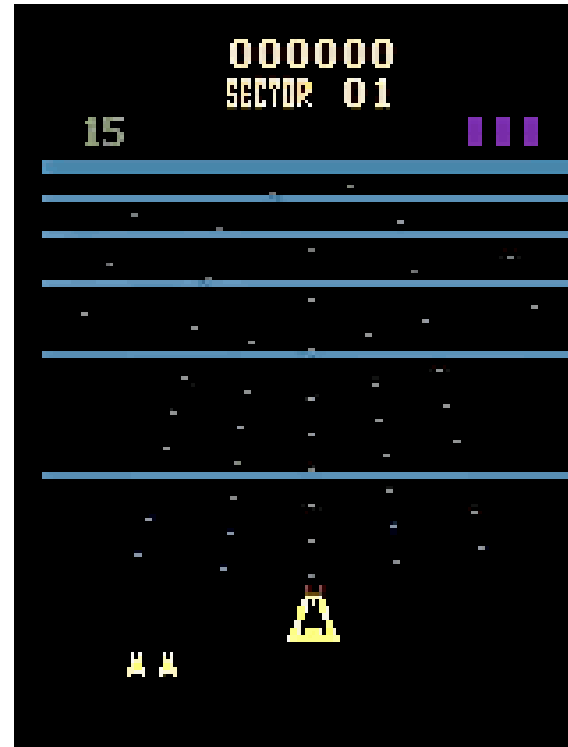
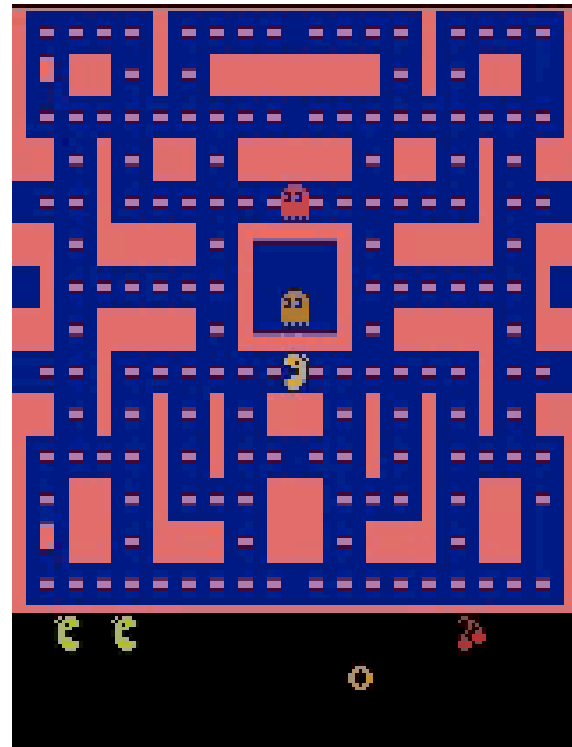
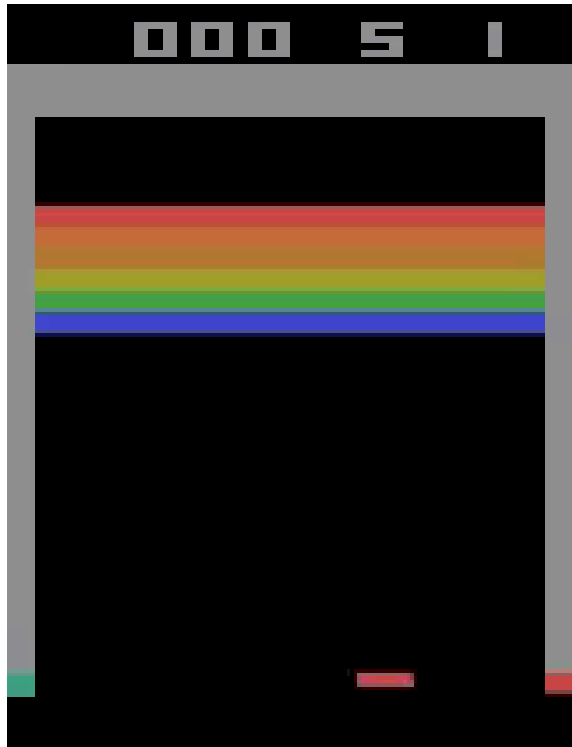


Mnih et al., "Human Level Control through Deep Reinforcement Learning," *Nature*, 2015

Learning to Play Atari Games



David
Silver



AlphaGo vs. Lee Sedol



Stable- Baselines 3

- [Stable Baselines3 \(SB3\)](#) is a set of reliable implementations of reinforcement learning algorithms in PyTorch

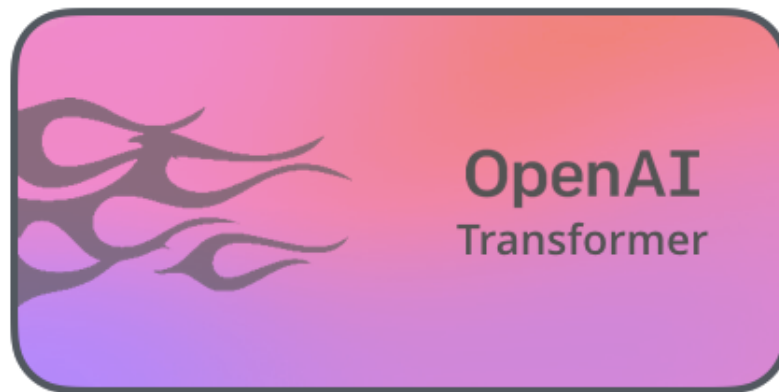


Generative AI

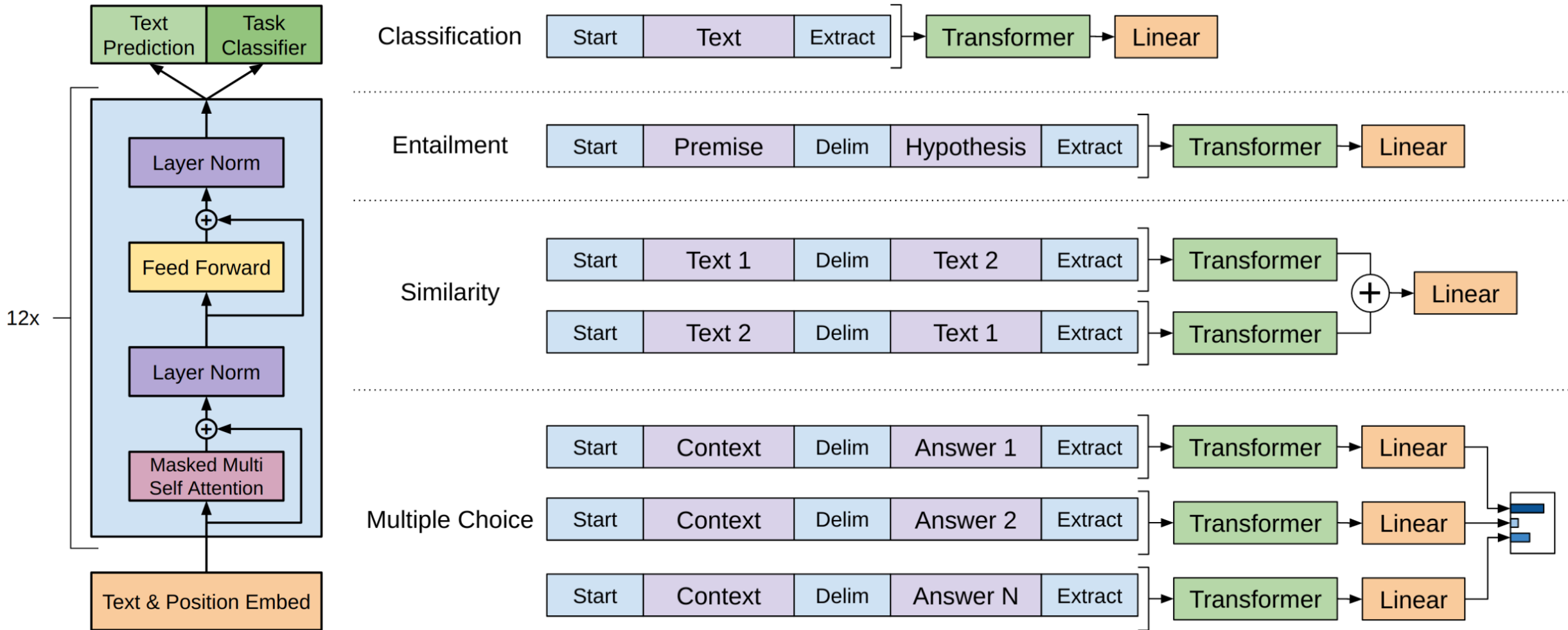
The image features a stylized illustration of a person's profile on the left, with a glowing blue face and a body composed of translucent, overlapping shapes. The background is a vibrant, abstract composition of musical staves with notes, rendered in a warm, glowing palette of orange, yellow, and purple. The overall aesthetic is futuristic and artistic, suggesting the intersection of technology and music.

OpenAI GPT: Pre-training Transformer Decoders

- Unsupervised pre-train transform decoders for predicting the next word (GPT: Generative Pre-Training)
- Use 12 Transformer decoders in GPT-1
 - GPT-1: [Improving Language Understanding with Unsupervised Learning \(2018\)](#)
 - GPT-2: [Better Language Models and Their Implications \(2019\)](#)
 - GPT-3: [Language Models are Few-Shot Learners \(2020\)](#)



OpenAI GPT for Different Tasks



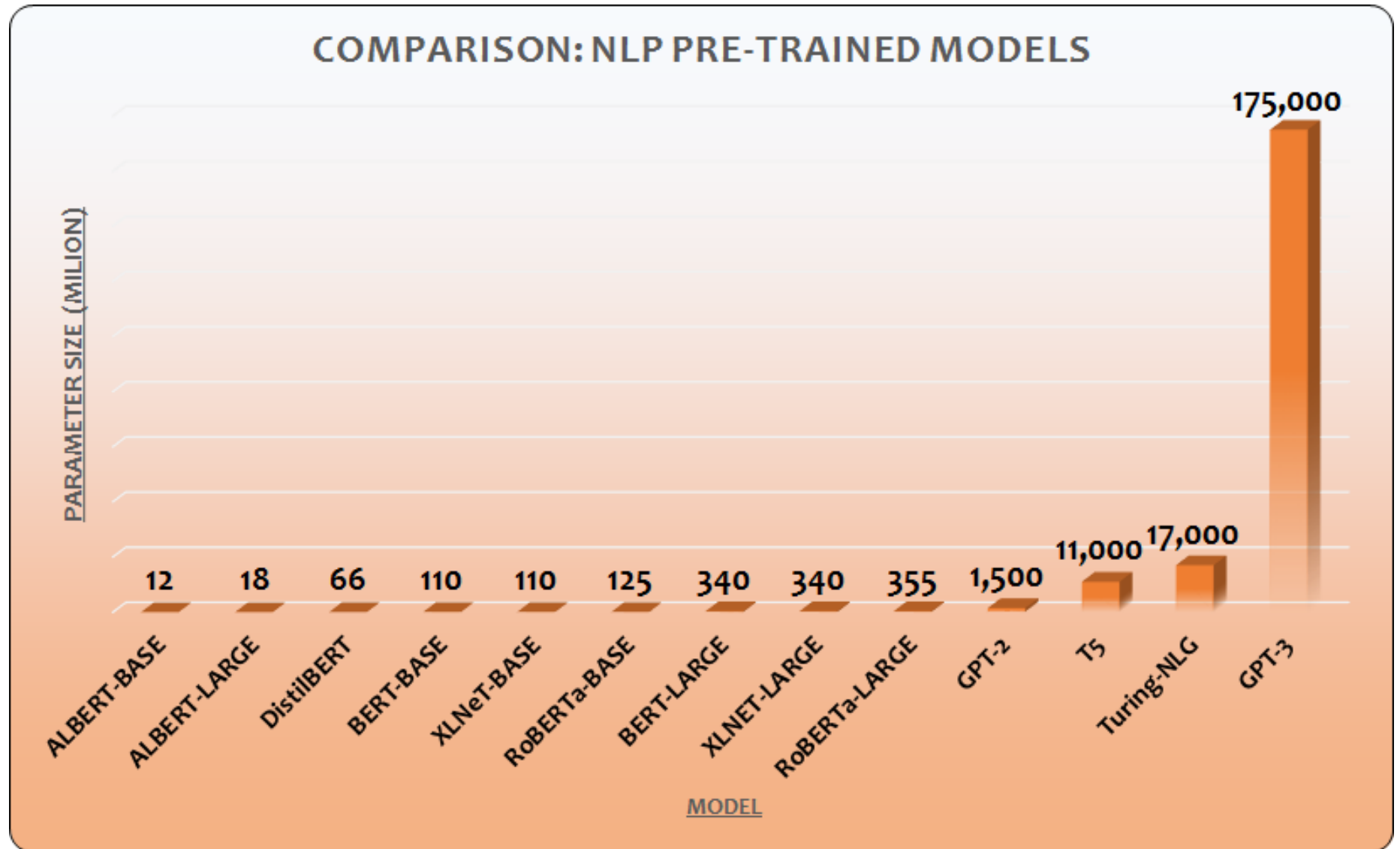
OpenAI GPT-2

- Pre-trained using 40GB of Internet text
- Scale-up of GPT with 10X parameters trained with 10X data
- Other tricks
 - Layer normalization was moved to the input of each sub-block
 - An additional layer normalization was added after the final self-attention block

Parameters	Layers	d_{model}
117M	12	768
345M	24	1024
762M	36	1280
1542M	48	1600

Size does Matter! GPT-3

- **175 Billion Parameters!**
- $175 \times 4 = 700\text{GB}$
- 55 years and \$4,600,000 to train - even with the lowest priced GPU cloud on the market.

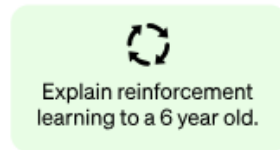


OpenAI ChatGPT

Step 1

Collect demonstration data and train a supervised policy.

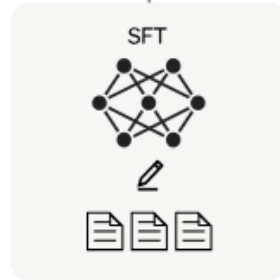
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



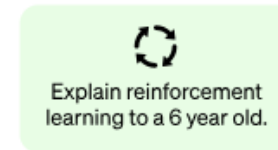
This data is used to fine-tune GPT-3.5 with supervised learning.



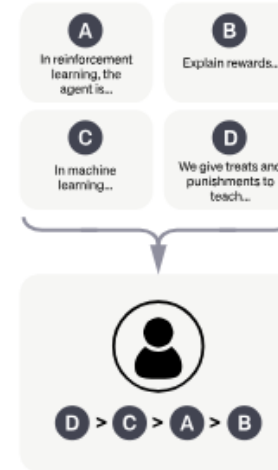
Step 2

Collect comparison data and train a reward model.

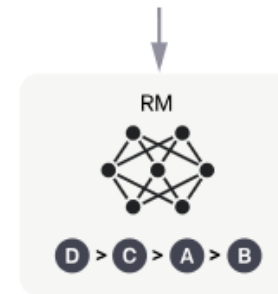
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using the PPO reinforcement learning algorithm.

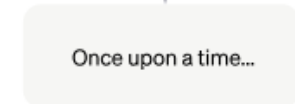
A new prompt is sampled from the dataset.



The PPO model is initialized from the supervised policy.



The policy generates an output.



The reward model calculates a reward for the output.

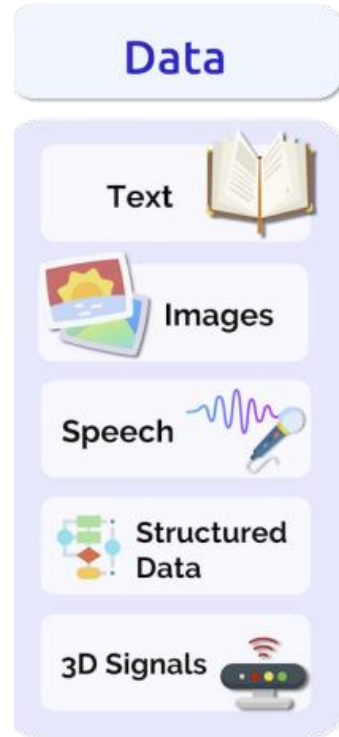


The reward is used to update the policy using PPO.

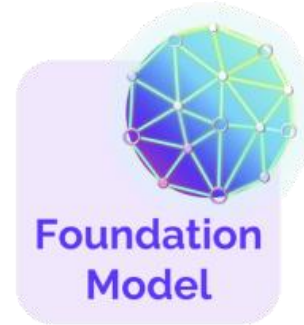


Foundation Models (基石模型)

- One model for All (2021)



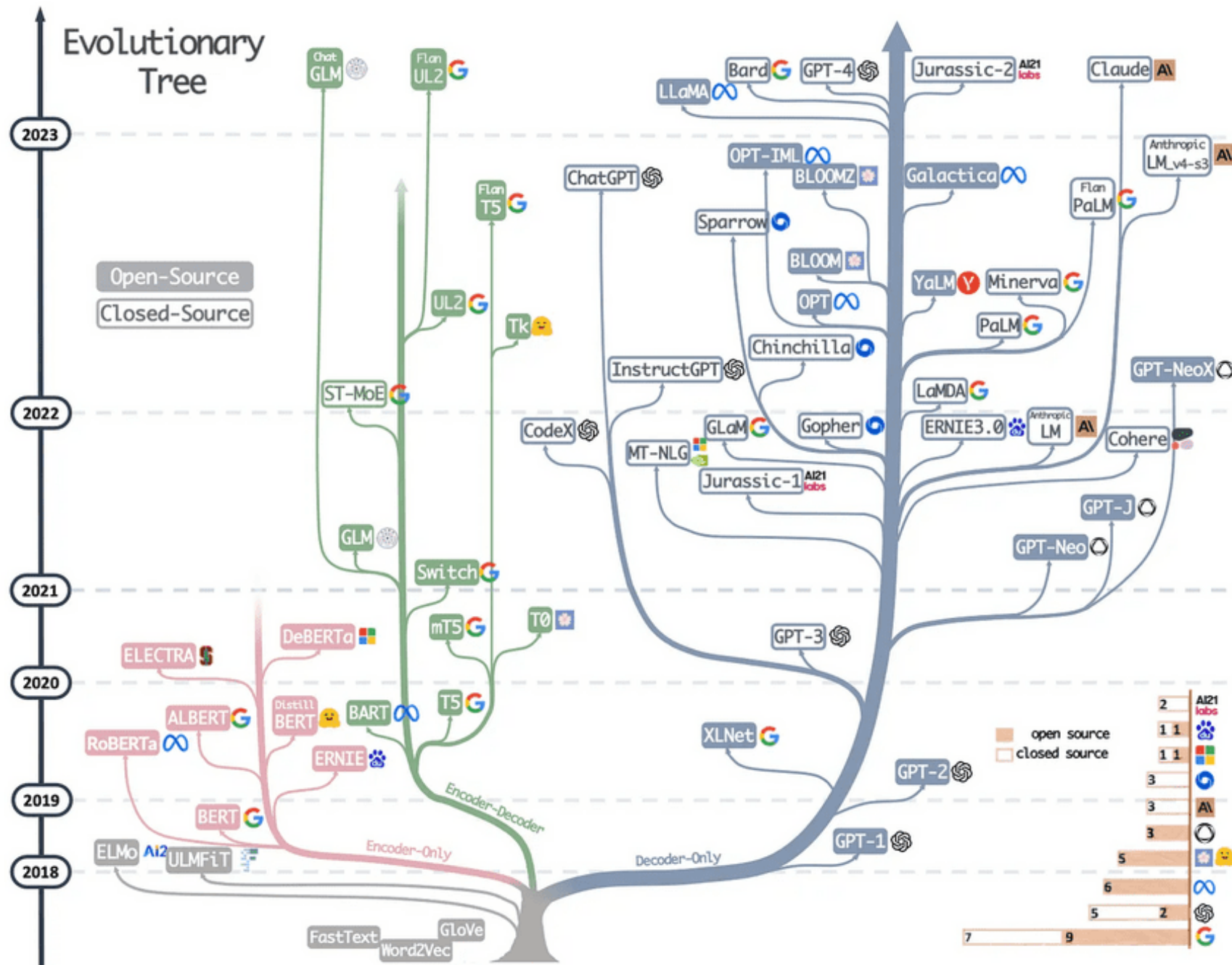
Training



Adaptation



Large Language Model (LLM) Practical Guide





Hugging Face: Free LLM models

 **Hugging Face**

Search models, datasets, users...

[Models](#) [Datasets](#) [Spaces](#) [Docs](#) [Solutions](#) [Pricing](#)

Tasks [Libraries](#) [Datasets](#) [Languages](#) [Licenses](#) [Other](#)

Filter Tasks by name

Multimodal

- Feature Extraction
- Text-to-Image
- Image-to-Text
- Text-to-Video
- Visual Question Answering
- Document Question Answering
- Graph Machine Learning

Computer Vision

- Depth Estimation
- Image Classification
- Object Detection
- Image Segmentation
- Image-to-Image
- Unconditional Image Generation
- Video Classification
- Zero-Shot Image Classification

Natural Language Processing


- Text Classification
- Token Classification
- Table Question Answering
- Question Answering
- Zero-Shot Classification
- Translation
- Summarization
- Conversational
- Text Generation
- Text2Text Generation

Models 235,314

Filter by name

new Full-text search


Sort: Most Downloads

 jonatasgrosman/wav2vec2-large-xlsr-53-english

 Updated Mar 25 • 71.9M • 182

xlm-roberta-large

 Updated Apr 7 • 42.6M • 160

 openai/clip-vit-large-patch14

 Updated Oct 4, 2022 • 16.8M • 460

roberta-base


 Updated Mar 6 • 12.2M • 176

distilbert-base-multilingual-cased

 Updated Apr 6 • 11.6M • 60

xlm-roberta-base

 Updated Apr 7 • 9.14M • 325

 microsoft/deberta-base

 Updated Sep 26, 2022 • 6.41M • 43

bert-large-uncased

 Updated Nov 15, 2022 • 5.18M • 33


bert-base-uncased


 Updated 26 days ago • 50.5M • 923

gpt2

 Updated Dec 16, 2022 • 17.3M • 1.18k

 sociocom/MedNER-CR-JA


 Updated Apr 5 • 15.7M • 5

 laion/CLIP-ViT-B-16-laion2B-s34B-b88K

 Updated Apr 20 • 11.7M • 6

distilbert-base-uncased


 Updated Nov 16, 2022 • 10.9M • 216

 microsoft/layoutlmv3-base

Updated Apr 12 • 8.19M • 168

bert-base-cased

 Updated Nov 16, 2022 • 6.38M • 114

 deepset/sentence_bert

Updated May 19, 2021 • 4.92M • 15

<https://huggingface.co/learn/nlp-course/chapter1/1>

BigScience Large Open-science Open-access Multilingual Language Model (BLOOM)

- With its 176 billion parameters, BLOOM is able to generate text in 46 natural languages and 13 programming languages.



a BigScience initiative

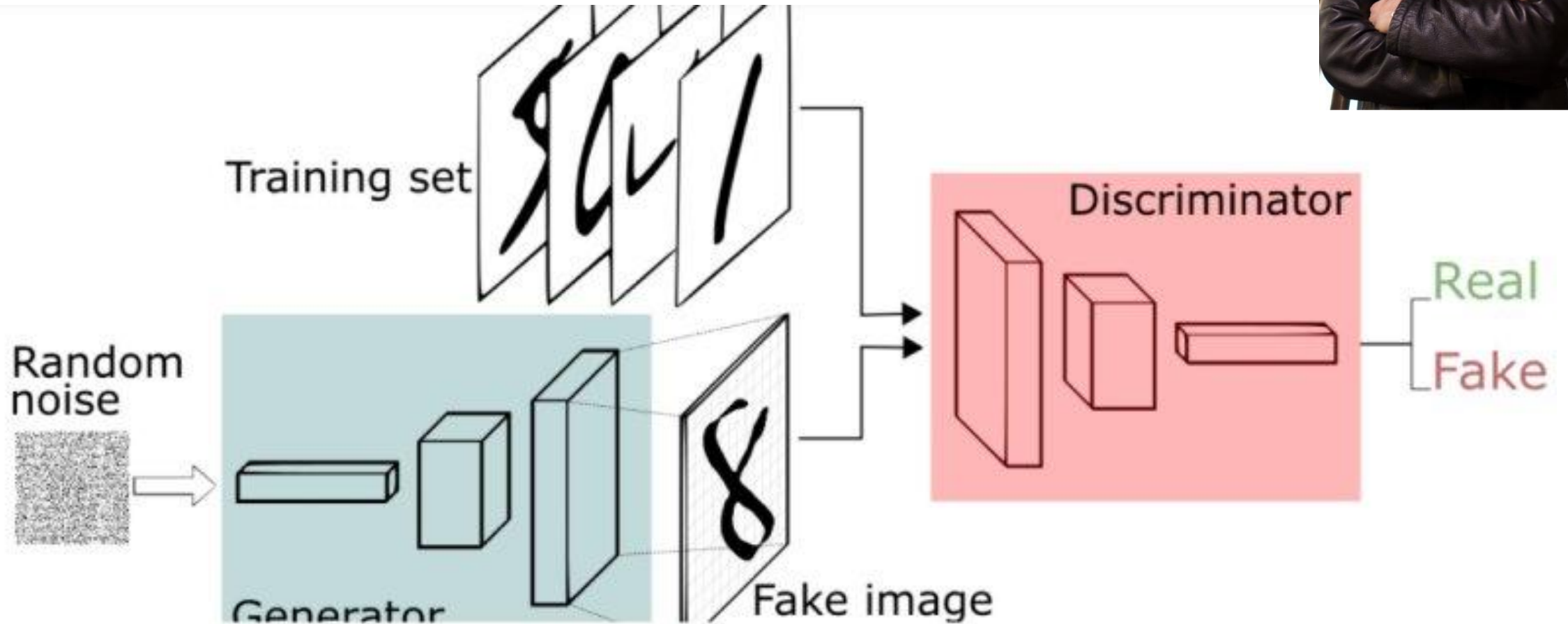
BL   **M**

176B params · 59 languages · Open-access

<https://huggingface.co/bigscience/bloom>

Generative Adversarial Networks (GAN)

- Ian Goodfellow



Super Resolution

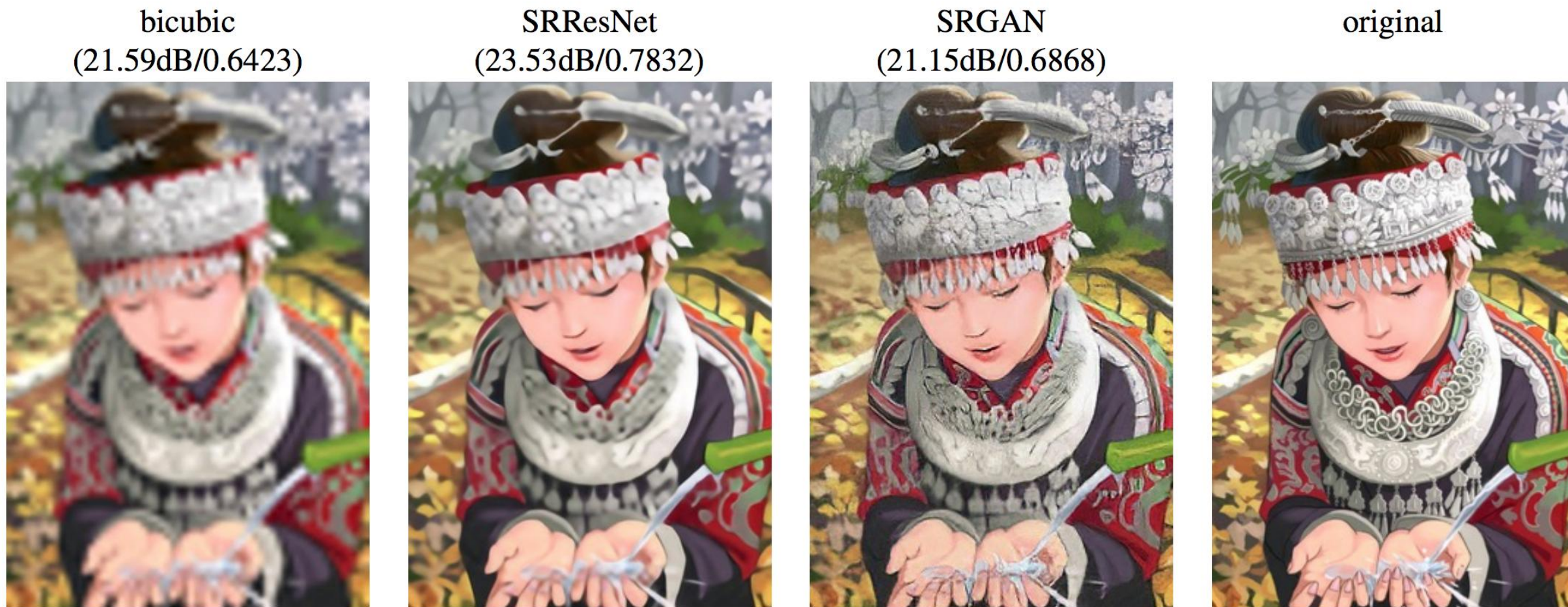
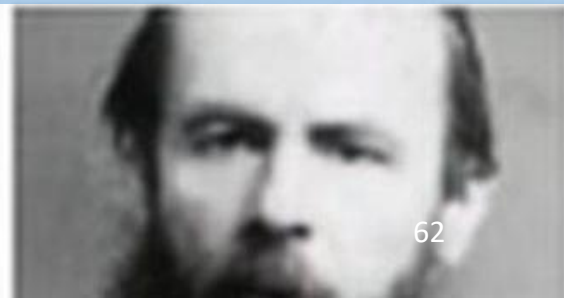


Figure 2: From left to right: bicubic interpolation, deep residual network optimized for MSE, deep residual generative adversarial network optimized for a loss more sensitive to human perception, original HR image. Corresponding PSNR and SSIM are shown in brackets. [4× upscaling]

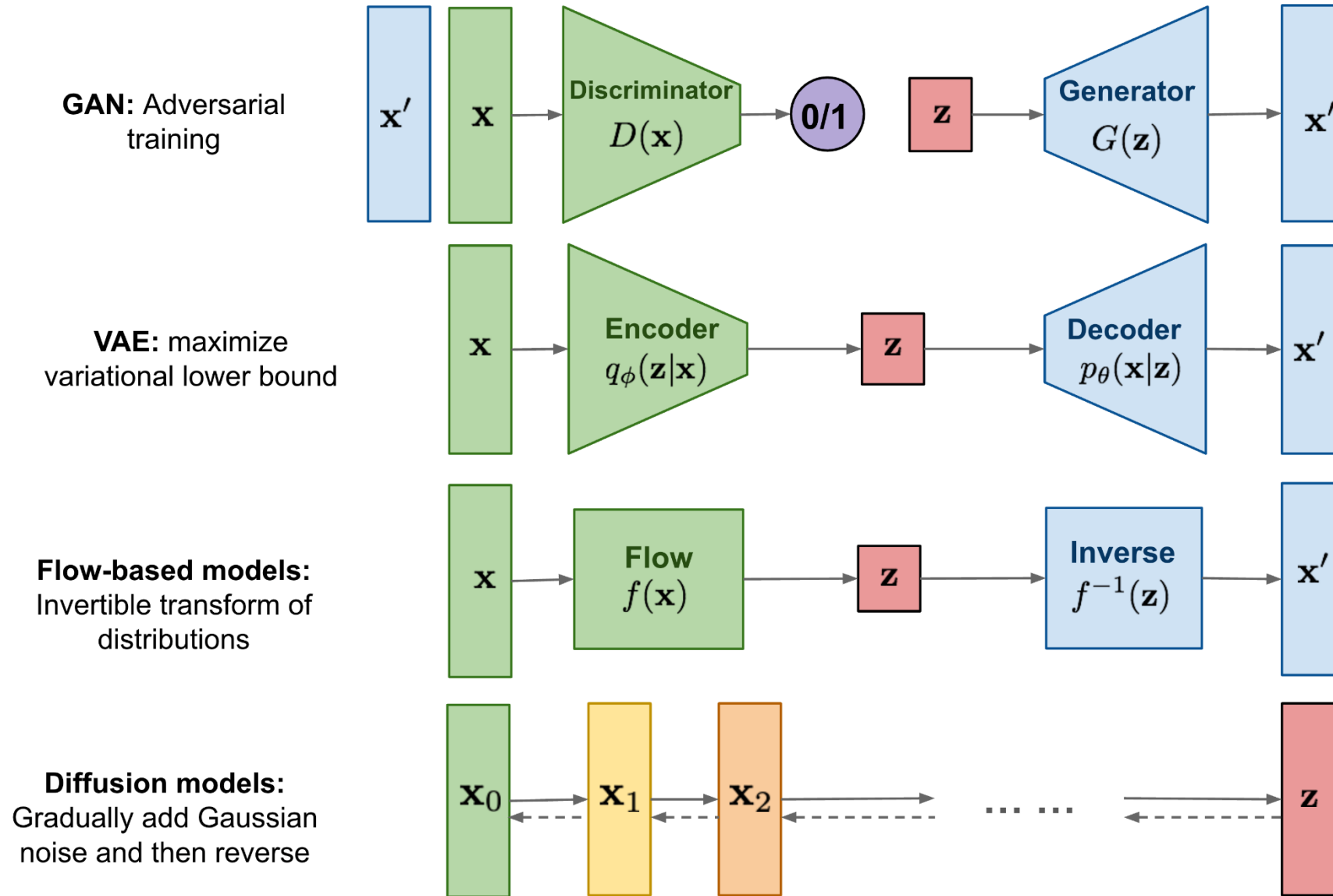




DeepFake: Is this you?



Overview of Different Generative Models



Diffusion is All You Need!

- Reverse diffusion process
- Flexible and tracible

Deep Unsupervised Learning using Nonequilibrium Thermodynamics

Jascha Sohl-Dickstein
Stanford University

JASCHA@STANFORD.EDU

Eric A. Weiss
University of California, Berkeley

EAWISS@BERKELEY.EDU

Niru Maheswaranathan
Stanford University

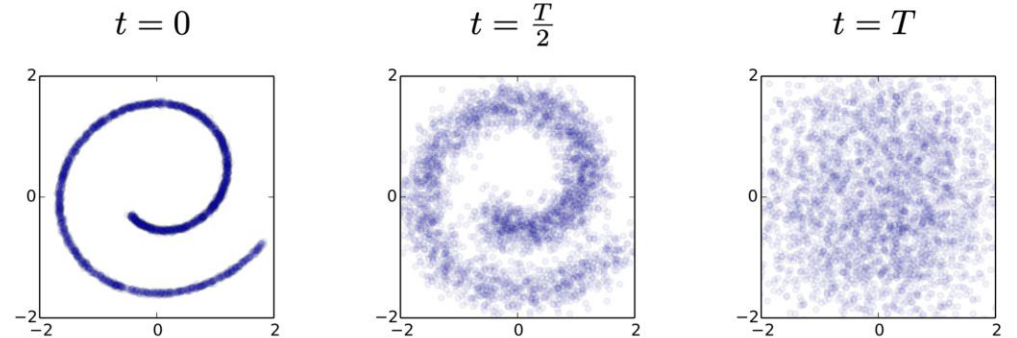
NIRUM@STANFORD.EDU

Surya Ganguli
Stanford University

SGANGULI@STANFORD.EDU

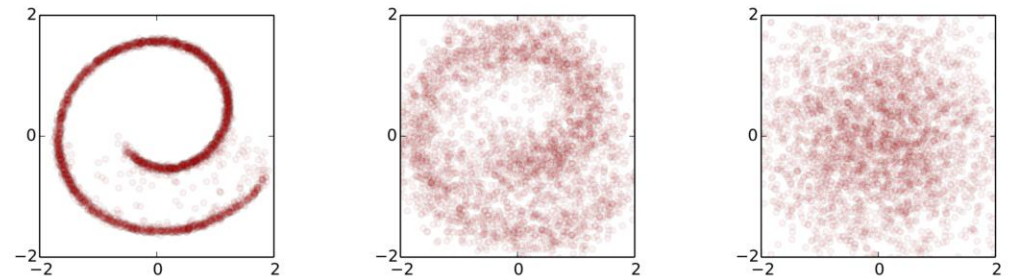
The forward trajectory

$$q(\mathbf{x}_{0:T})$$



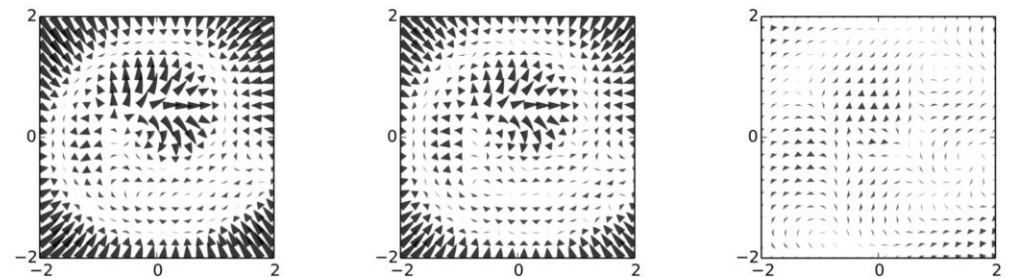
The reverse trajectory

$$p_{\theta}(\mathbf{x}_{0:T})$$



The drifting term

$$\mu_{\theta}(\mathbf{x}_t, t) - \mathbf{x}_t$$



<https://arxiv.org/pdf/1503.03585.pdf>

Key Takeaways

1. Deep learning is a branch of Machine Learning, which is a sub-field of Artificial Intelligence.
2. There are two stages in machine learning: training (learning) and testing (inference).
3. Gradient Descent is used to train NN models by updating weights to minimize the prediction errors.
4. Convolutional Neural Networks (CNN) are used to recognize images.
5. RNN and LSTM are used to recognize sequential data such as text or speech.
6. Transformer told us that attention is all you need!
7. Generative Adversarial Networks (GANs) can be used to generate fake data, but now maybe diffusion is all you need.