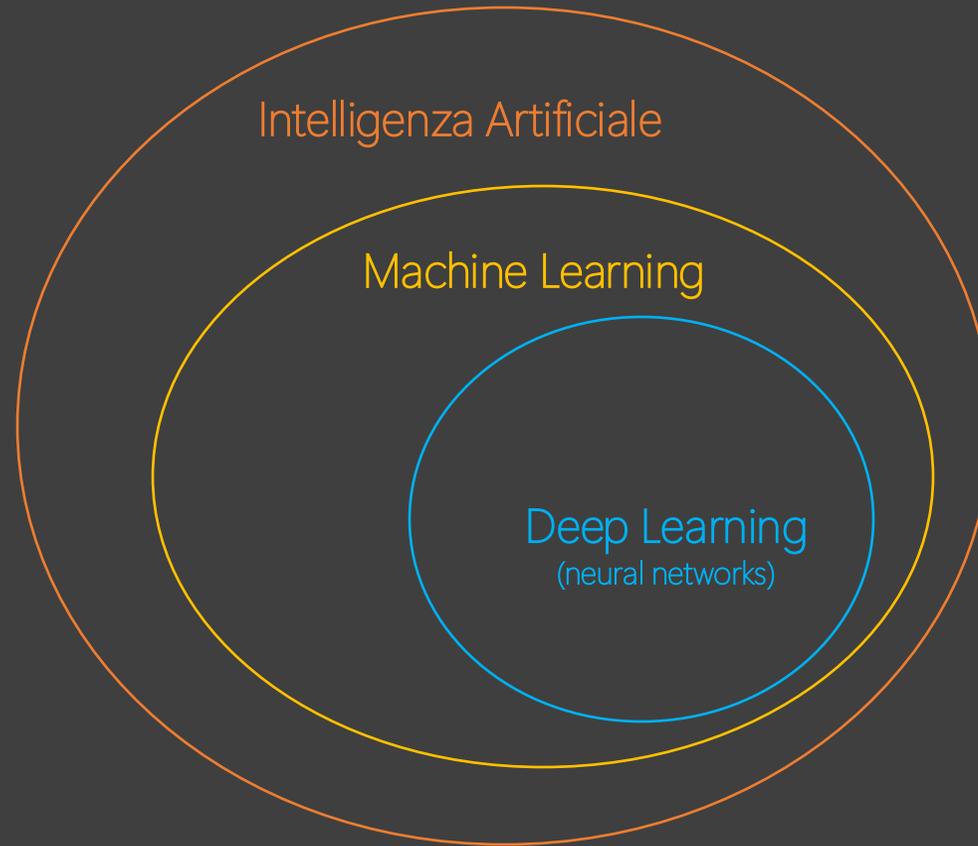


# Artificiale o Intelligente?

# Contesto



# What is Machine Learning?

ML is the study of **ALGORITHMIC** prediction

Machines that LEARN how to PREDICT

Algo1: LEARNER      Algo2: PREDICTOR

Algorithm: is a finite sequence of rigorous instructions, typically used to

SOLVE A CLASS of specific problems

# ML - definitions

We want the machine LEARN a PREDICTOR (the learner fixes parameters in the predictor)

Domain set:

$\mathcal{X}$  Usually, you have vectors  $\underline{x}$  of features

Label set:

$\mathcal{Y}$  for example  $[0,1]$

Training set:

$S = ((x_1, y_1), \dots, (x_N, y_N))$   
Sequence of pairs  $\mathcal{X} \times \mathcal{Y}$

|| What the learner can use to determine the values of the predictor's parameters

Learner output:

Is the predictor: a function  $h: \mathcal{X} \rightarrow \mathcal{Y}$

# Data-generation model

What are the data we use?

①

Instances from the training dataset are generated by some probability distribution

$\mathcal{D}$  is a probability distribution over the domain set  $\mathcal{X}$

The learner knows nothing about the distribution  $\mathcal{D}$

②

There is a "correct" labeling function  $f: \mathcal{X} \rightarrow \mathcal{Y}$

$$y_i = f(x_i) \forall i$$

This function is what the learner must learn.

The training dataset is the outcome of:

- Use  $\mathcal{D}$  to extract  $x_i$  from  $\mathcal{X}$
- Use  $f$  to compute  $y_i$  given  $x_i$

## Recap:

- We have a training set  $S$  sampled from an UNKNOWN distribution  $\mathcal{D}$  and labeled by some target function  $f$
- The learner output is a predictor

$$h_S: \mathcal{X} \rightarrow \mathcal{Y}$$

 The predictor depend on  $S$

GOAL: find  $h_S$  that minimizes the error with respect to the UNKNOWN  $\mathcal{D}$  and  $f$

# Empirical Risk Minimization & loss function

$L_{\mathcal{D},f}(h)$ , the true error, is not directly available because we do not know  $\mathcal{D}$  and  $f$

→ TRAINING ERROR: error of the predictor on the training samples

$$L_S(h) \stackrel{\text{def}}{=} \frac{|\{i \in [N] : h(x_i) \neq y_i\}|}{N} \quad [N] = \{1, \dots, N\} \quad \text{EMPIRICAL RISK/ERROR}$$

A loss function takes two inputs,  $\hat{y} = h(x)$  and the target values  $y$ , and returns a **real number**  $loss(\hat{y}, y)$  that we interpret as a quantified error for predicting  $\hat{y}$  when the target is  $y$ .

The risk associated with  $h$  is defined as

$$L(h) = \mathbb{E}(loss(h(\mathcal{X}), \mathcal{Y}))$$

$$L_S(h) = \frac{1}{N} \sum_{i=1}^N loss(h(x_i), y_i)$$

# Deep Learning

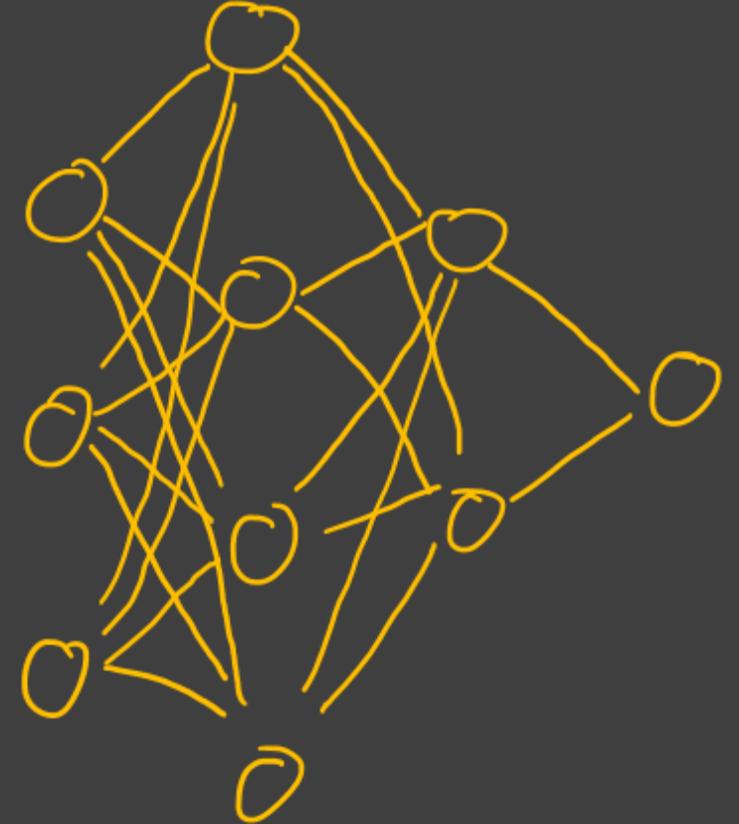
How do we choose  $h(x)$  ?

Many different algorithms available to approximate the target function

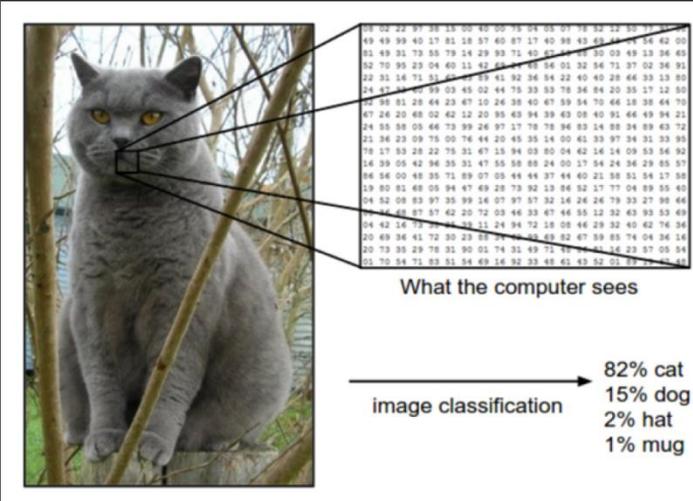
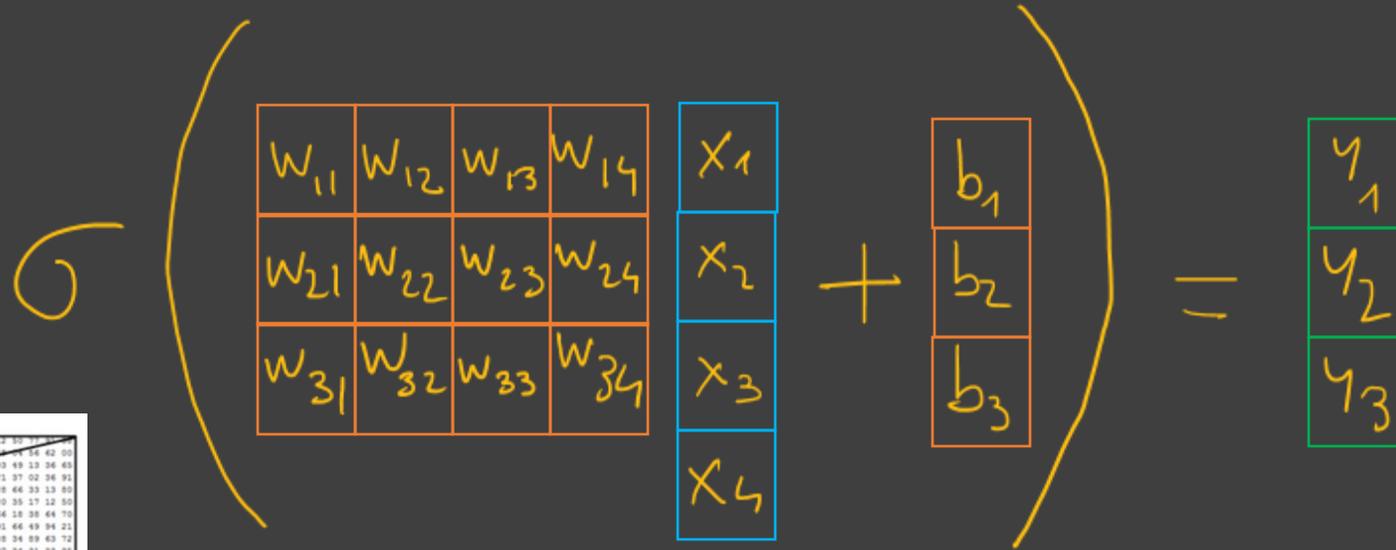
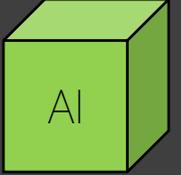
BUT

Today most of the people use:

## Deep Neural Networks



# The core of a neural network



$$Y_j = \sigma(X_i W_{ij} + b_j)$$

Kevin P. Murphy, Probabilistic Machine Learning, An Introduction

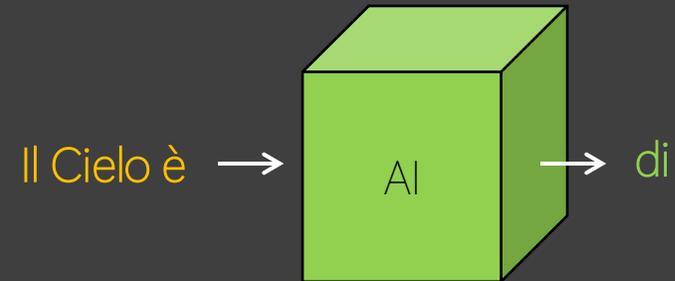
# Deep Learning: types of network

- Fully connected → Tabular data
  - Convolutional → Images
  - Recurrent → Sequence
  - Transformer → Sequence
  - [ Diffusion ] → Images
- Large Language models [ChatGPT]
- Images generative models [Dall-E]
-

# Large Language Models

# Intelligenza Artificiale: Teoria delle probabilità sotto steroidi

- 1 [Cielo](#) – Wikipedia
- 2 Rino Gaetano - Ma il [cielo](#) è sempre più [blu](#) (Official Video)
- 3 Ma il [cielo](#) è sempre più [blu](#) (Extended Version) – YouTube
- 4 Bungaro e Fiorella Mannoia - Il [Cielo](#) è di [Tutti](#) ... – YouTube
- 5 Perché il [cielo](#) è [blu](#)?
- 6 Perché il [cielo](#) è [blu](#)? La spiegazione semplice
- 7 Scienza Per [Tutti](#) - 0264. Perché il [cielo](#) è azzurro? – INFN
- 8 Perché il [cielo](#) è [blu](#) – base
- 9 Il [cielo](#) è di [tutti](#). Ediz. a colori: Il cielo e di [tutti](#) - Amazon.it
- 10 Il [cielo](#) è di [tutti](#)
- 11 Il [cielo](#) è di [tutti](#) - Gianni Rodari
- 12 Il [cielo](#) non è [blu](#), lo dice la scienza
- 13 Il [cielo](#) è dei leggeri - Matteo Munaretto
- 14 Il [cielo](#) è di [tutti](#), la terra è di [tutti](#) - Mirca Benetton | Ed. ETS
- 15 Il [cielo](#) è dipinto di stelle
- 16 Il [Cielo](#) è di [Tutti](#)! - Video – RaiPlay
- 17 Perché il [Cielo](#) è Azzurro? Come spiegarlo ai bambini
- 18 Il [cielo](#) è dei violenti di Flannery O'Connor
- 19 Il [cielo](#) è di [tutti](#) - Rodari/Costa | Emme Edizioni
- 20 Perché il [cielo](#) è [blu](#)?

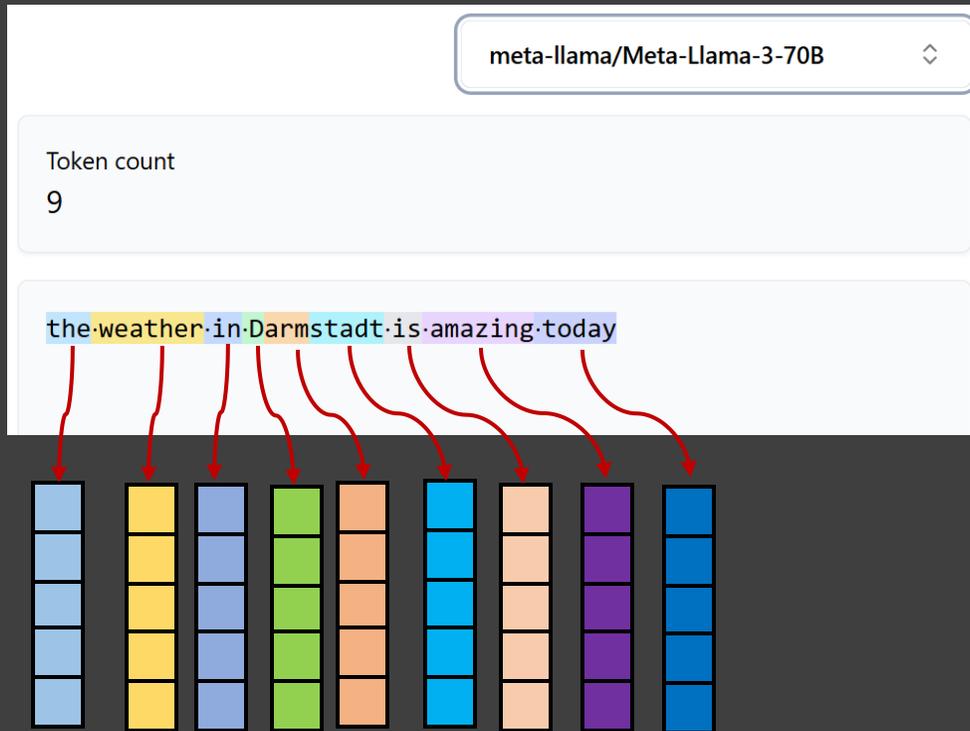


# Tokens and Embeddings

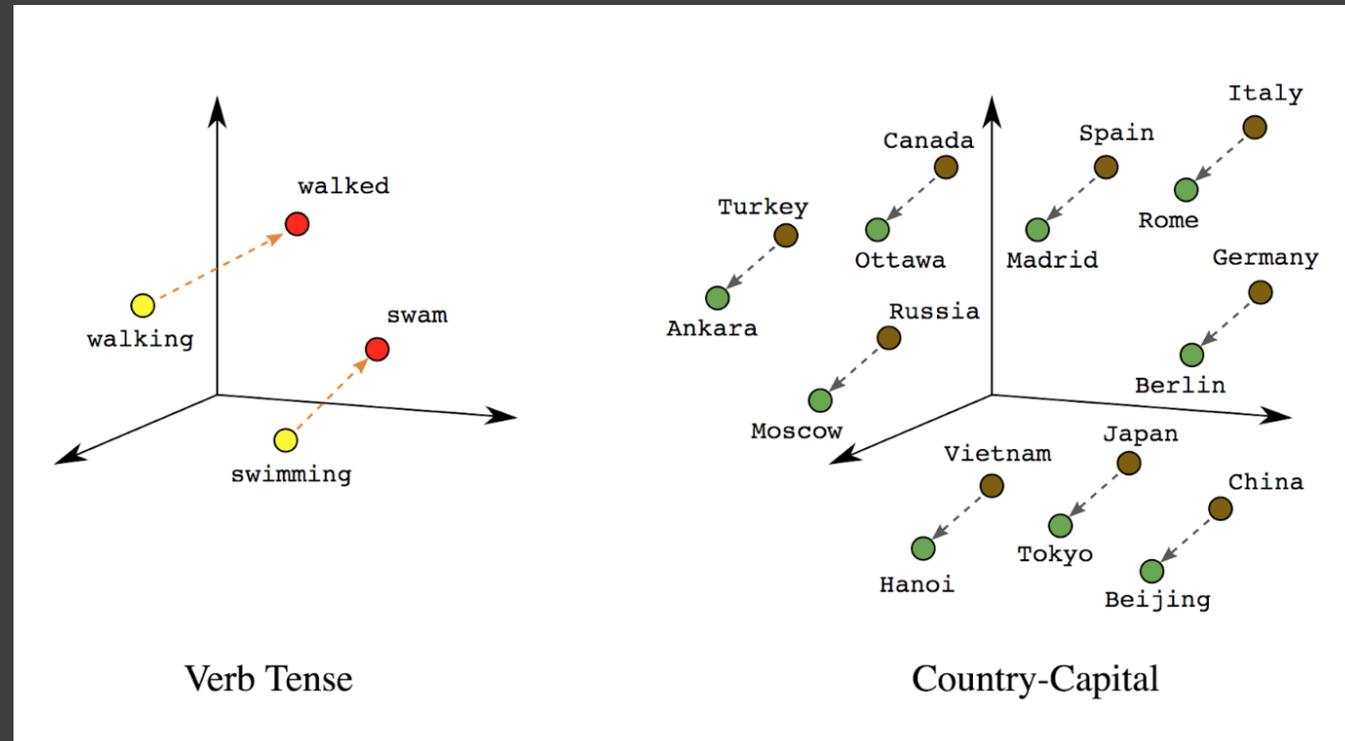
Computer cannot understand natural language, *words become vectors*.

Vector embedding

Can capture semantic meaning and relationships.

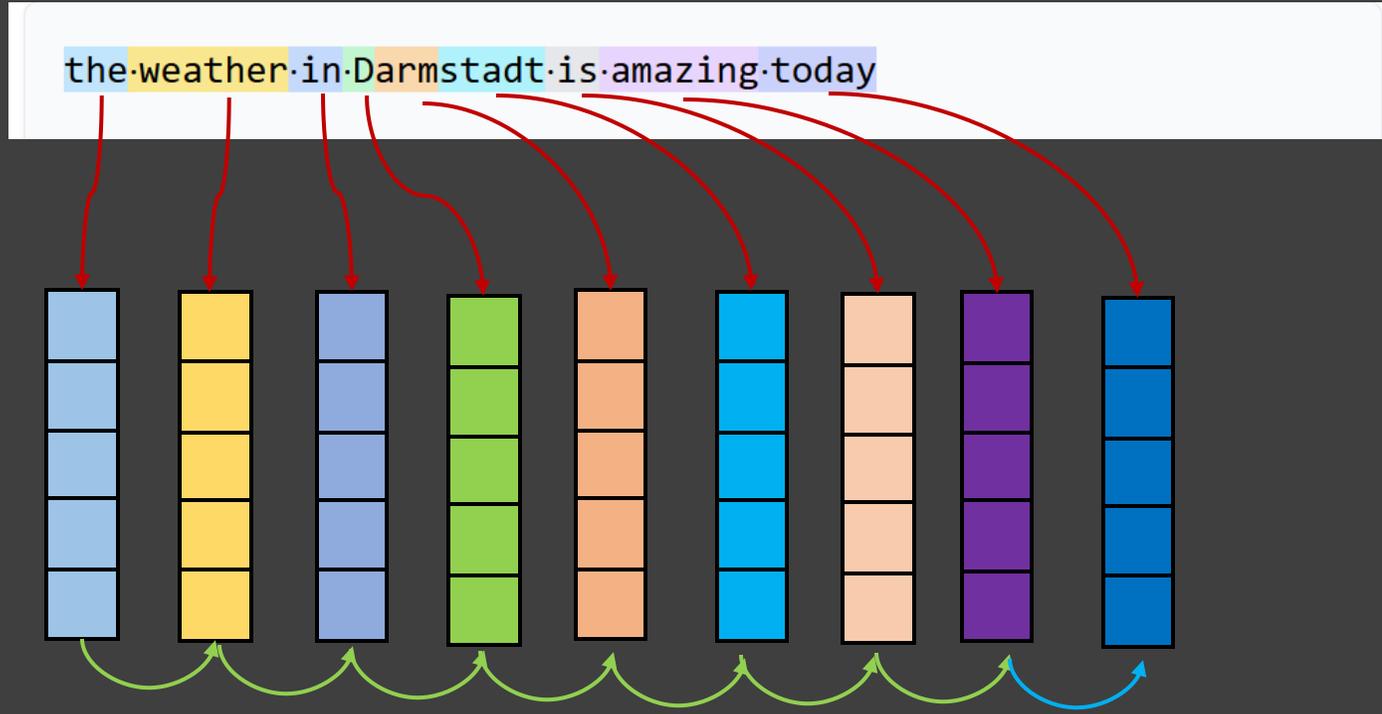


<https://tiktokenizer.vercel.app>



Source: [Google for Developers](https://developers.google.com/ai/vector-embeddings/)

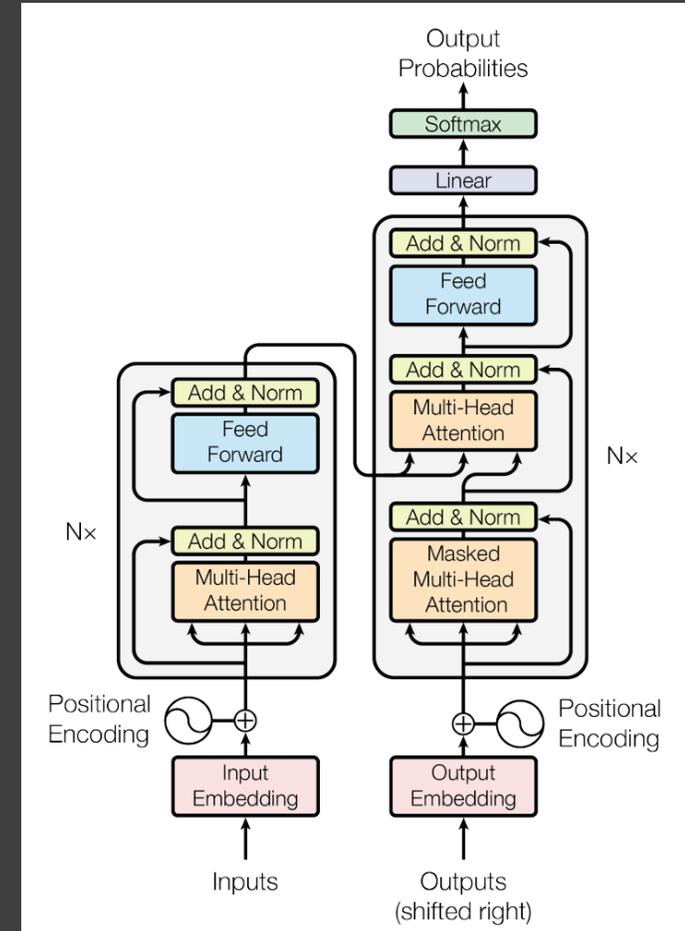
# Tokens and Embeddings



Tokens are processed autoregressively

# Attention mechanism

- Transformer is the go-to architecture for deep learning models, powering text-generative models like OpenAI's GPT, Meta's Llama (open), and Google's Gemini.
- Main innovation is the self-attention mechanism, which allows them to process entire sequences and capture long-range dependencies.
- Latest models can have more than 500 billion parameters, which are weights learned by the model during training.
- LLMs demand enormous data and compute power, balancing accuracy, cost, and sustainability.



Source: [Attention Is All You Need](#)

# What is an LLM

Auto-regressive system\* that, **given the user prompt**, will complete it with the **most likely words**.

The actress that played Rose in the 1997 film Titanic is named **Kate Winslet**. She was born on Octob

Kate = 65.16%
Gloria = 10.15%
Frances = 3.65%
Rose = 2.00%
Billy = 1.44%

Different from humans (typically) thinking first about what they want to say and then how.

\*the output variable depends on its previous values.

# Due questioni:

- ① Quali sono i rischi legati allo sviluppo dei sistemi basati su Intelligenza Artificiale?
- ② L'Intelligenza Artificiale è intelligente?

## 2018 Turing Award winners



Yoshua Bengio

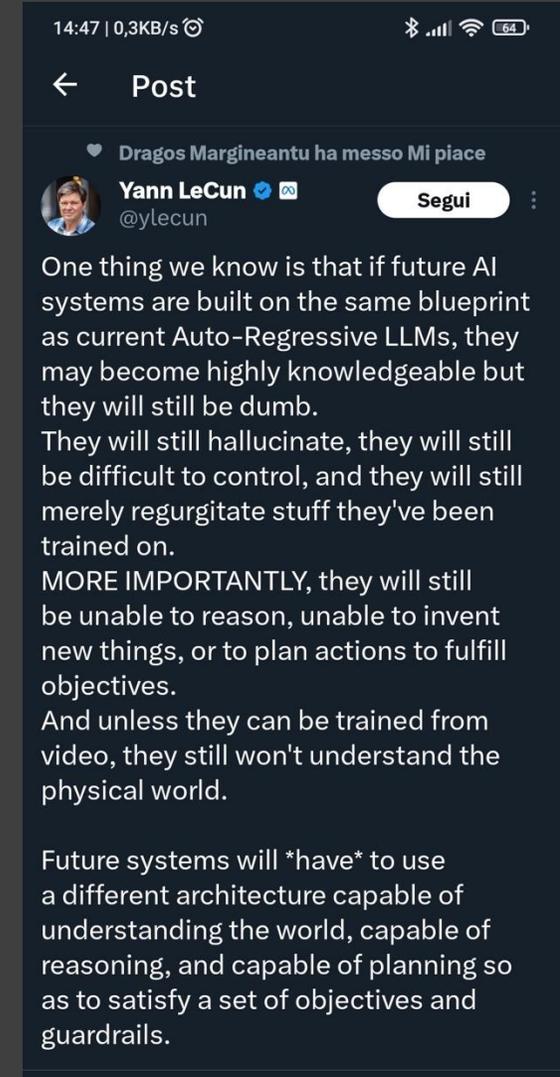
Geoffrey Hinton

Yann LeCun

### Managing AI Risks in an Era of Rapid Progress

Yoshua Bengio      Mila - Quebec AI Institute, Université de Montréal, Canada CIFAR AI Chair  
 Geoffrey Hinton      University of Toronto, Vector Institute  
 Andrew Yao      Tsinghua University

Yoshua Bengio, Geoffrey Hinton, Andrew Yao, Dawn Song, Pieter Abbeel, Yuval Noah Harari, Ya-Qin Zhang, Lan Xue, Shai Shalev-Shwartz, Gillian Hadfield, Jeff Clune, Tegan Maharaj, Frank Hutter, Atılım Güneş Baydin, Sheila McIlraith, Qi Qi Gao, Ashwin Acharya, David Krueger, Anca Dragan, Philip Torr, Stuart Russell, Daniel Kahneman, Jan Brauner, Sören Mindermann, arXiv:2310.17688, 26 Oct 2023



Twitter, 28/10/2023

# Quali sono i rischi legati allo sviluppo dei sistemi basati su Intelligenza Artificiale?

## Societal-scale risks

AI systems could rapidly come to outperform humans in an increasing number of tasks. If such systems are not carefully designed and deployed, they pose a range of societal-scale risks. They threaten to amplify social injustice, erode social stability, and weaken our shared understanding of reality that is foundational to society. They could also enable large-scale criminal or terrorist activities. Especially in the hands of a few powerful actors, AI could cement or exacerbate global inequities, or facilitate automated war-

fare, customized mass manipulation, and pervasive surveillance<sup>12,13</sup>.

Many of these risks could soon be amplified, and new risks created, as companies are developing *autonomous AI*: systems that can plan, act in the world, and pursue goals. While current AI systems have limited autonomy, work is underway to change this<sup>14</sup>. For example, the non-autonomous GPT-4 model was quickly adapted to browse the web<sup>15</sup>, design and execute chemistry experiments<sup>16</sup>, and utilize software tools<sup>17</sup> including other AI models<sup>18</sup>.

Quali sono i rischi legati allo sviluppo dei sistemi basati su Intelligenza Artificiale?

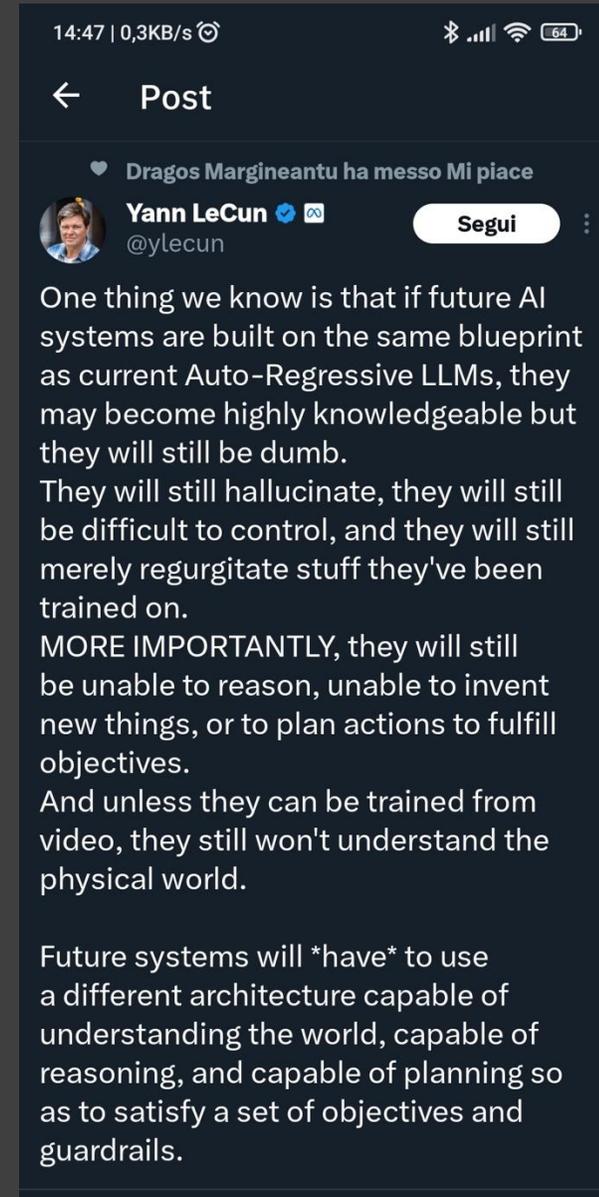
## LLAMA 2: Open Foundation and Fine-Tuned Chat Models

### Contents

<b>1</b>	<b>Introduction</b>	<b>3</b>
<b>2</b>	<b>Pretraining</b>	<b>5</b>
2.1	Pretraining Data	5
2.2	Training Details	5
2.3	LLAMA 2 Pretrained Model Evaluation	7
<b>3</b>	<b>Fine-tuning</b>	<b>8</b>
3.1	Supervised Fine-Tuning (SFT)	9
3.2	Reinforcement Learning with Human Feedback (RLHF)	9
3.3	System Message for Multi-Turn Consistency	16
3.4	RLHF Results	17
<b>4</b>	<b>Safety</b>	<b>20</b>
4.1	Safety in Pretraining	20
4.2	Safety Fine-Tuning	23
4.3	Red Teaming	28
4.4	Safety Evaluation of LLAMA 2-CHAT	29
<b>5</b>	<b>Discussion</b>	<b>32</b>
5.1	Learnings and Observations	32
5.2	Limitations and Ethical Considerations	34
5.3	Responsible Release Strategy	35
<b>6</b>	<b>Related Work</b>	<b>35</b>
<b>7</b>	<b>Conclusion</b>	<b>36</b>
<b>A</b>	<b>Appendix</b>	<b>46</b>
A.1	Contributions	46
A.2	Additional Details for Pretraining	47
A.3	Additional Details for Fine-tuning	51
A.4	Additional Details for Safety	58
A.5	Data Annotation	72
A.6	Dataset Contamination	75
A.7	Model Card	77



# l'Intelligenza Artificiale è intelligente?



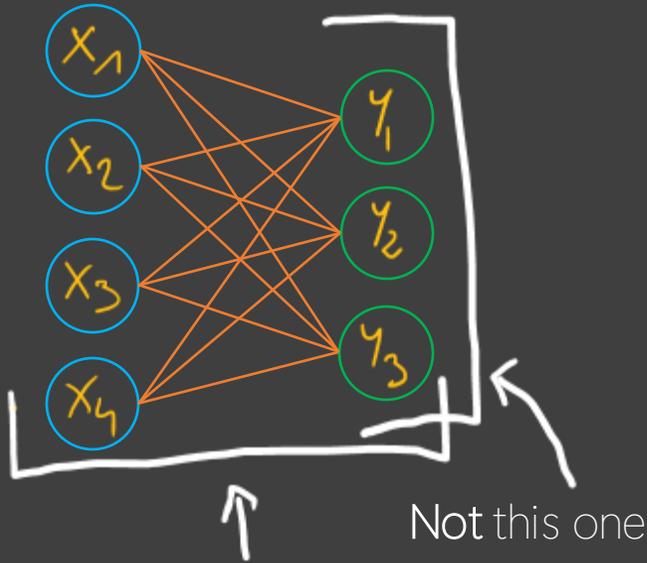
Twitter, 28/10/2023

l'Intelligenza Artificiale è intelligente?

Un inizio di risposta è «la domanda»  
l'Intelligenza Artificiale non si pone domande

# Backup

# Neural Network: the layer



The layer for me is this one

Not this one

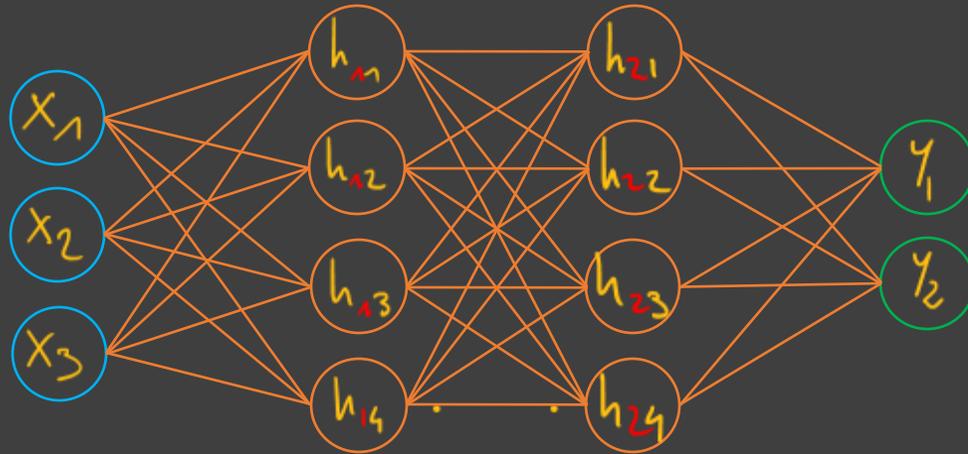
Activation function

$$\sigma \left( \begin{pmatrix} w_{11} & w_{12} & w_{13} & w_{14} \\ w_{21} & w_{22} & w_{23} & w_{24} \\ w_{31} & w_{32} & w_{33} & w_{34} \end{pmatrix} \begin{pmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{pmatrix} + \begin{pmatrix} b_1 \\ b_2 \\ b_3 \end{pmatrix} \right) = \begin{pmatrix} y_1 \\ y_2 \\ y_3 \end{pmatrix}$$

$$y_j = \sigma(x_i w_{ij} + b_j)$$

Linear combination +  
Non linear activation function

# Neural Network with hidden layers



$$h_{1i} = \sigma(x_i w_{ij}^1 + b_j^1)$$

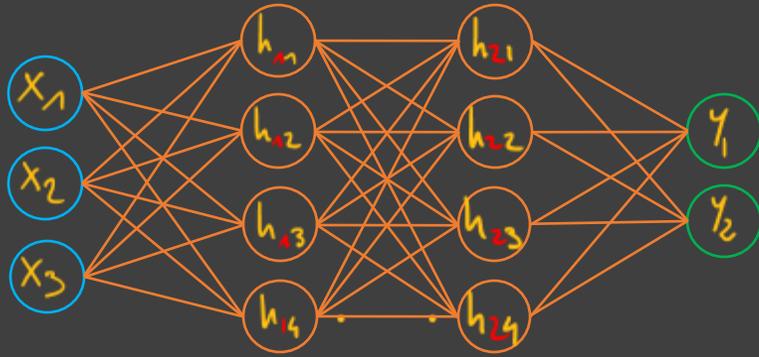
$$h_{2i} = \sigma(h_{1i} w_{ij}^2 + b_j^2)$$

$$y_j = \sigma(h_{2i} w_{ij}^3 + b_j^3)$$

E.g.: Sigmoid activation function

$$\sigma = \frac{1}{1 + e^{-x}}$$

# Empirical risk minimization



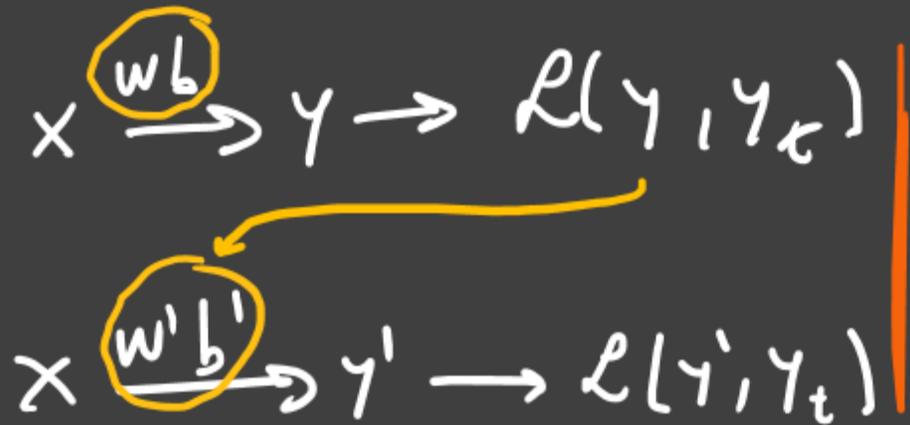
$$\begin{aligned}
 h_{1i} &= \sigma(x_i w_{ij}^1 + b_j^1) \\
 h_{2i} &= \sigma(h_{1i} w_{ij}^2 + b_j^2) \\
 y_j &= \sigma(h_{2i} w_{ij}^3 + b_j^3)
 \end{aligned}$$

$$\mathcal{L}_{\mathcal{D}_P}(h)$$

$$\begin{aligned}
 x \xrightarrow{w, b} y &\rightarrow \mathcal{L}(y, y_t) \\
 x \xrightarrow{w', b'} y' &\rightarrow \mathcal{L}(y', y_t)
 \end{aligned}$$

A yellow arrow points from the  $w, b$  term in the first equation to the  $w', b'$  term in the second equation, illustrating the process of finding a better model.

# Empirical risk minimization



## Gradient Descent

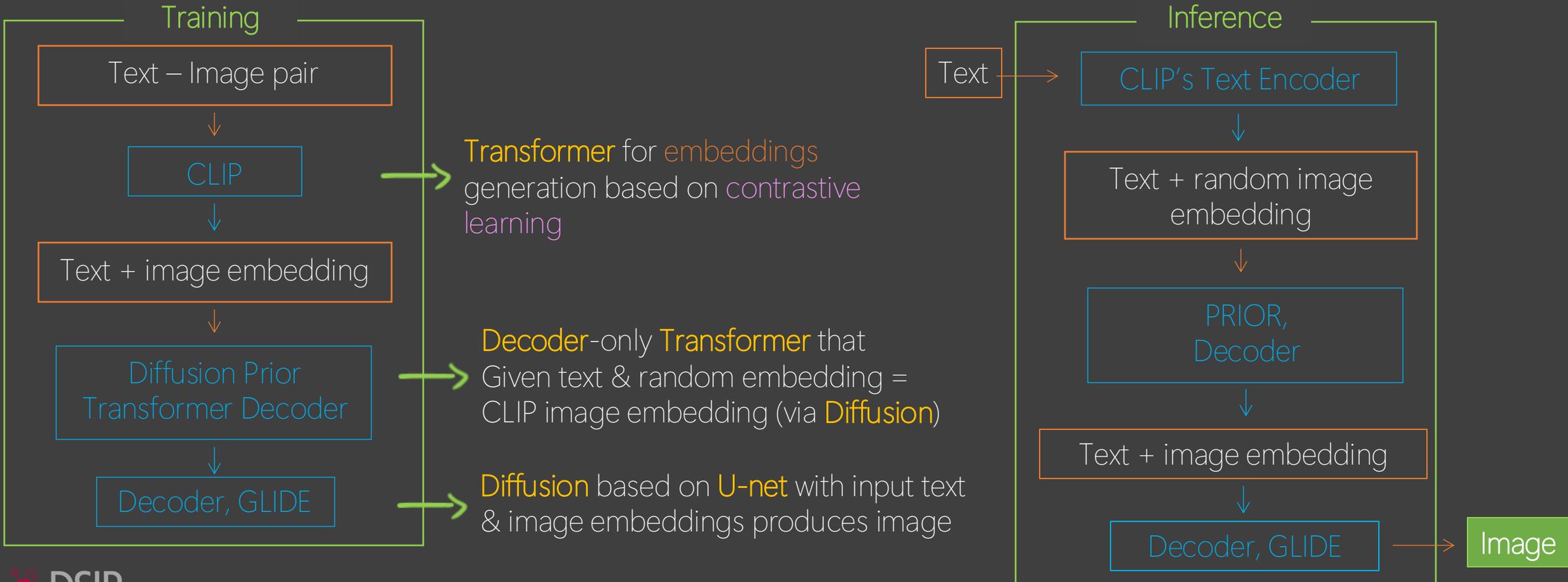
$$w' = w - \eta \nabla_w \mathcal{L}(w)$$

Learning rate

Need to compute the gradient wrt each weight

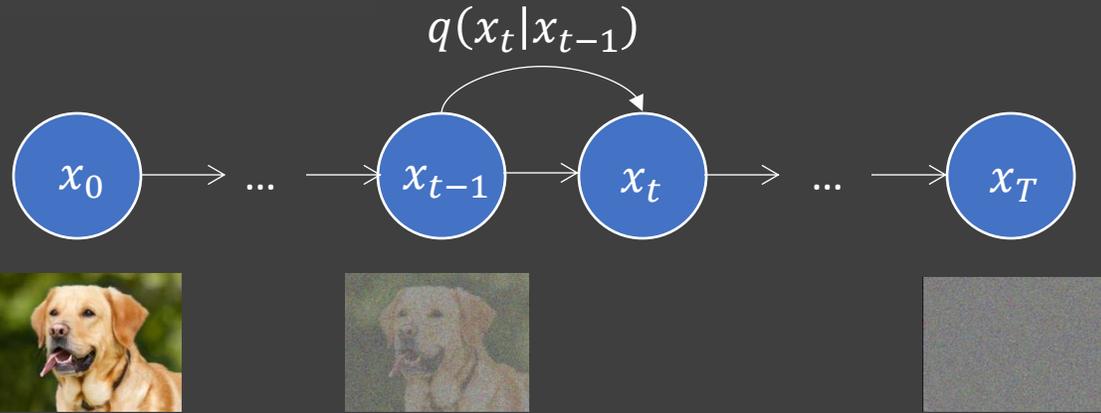
# Deep Learning example: DALL-E

Text-to-image generator → Multimodal



# Diffusion model

## Forward diffusion



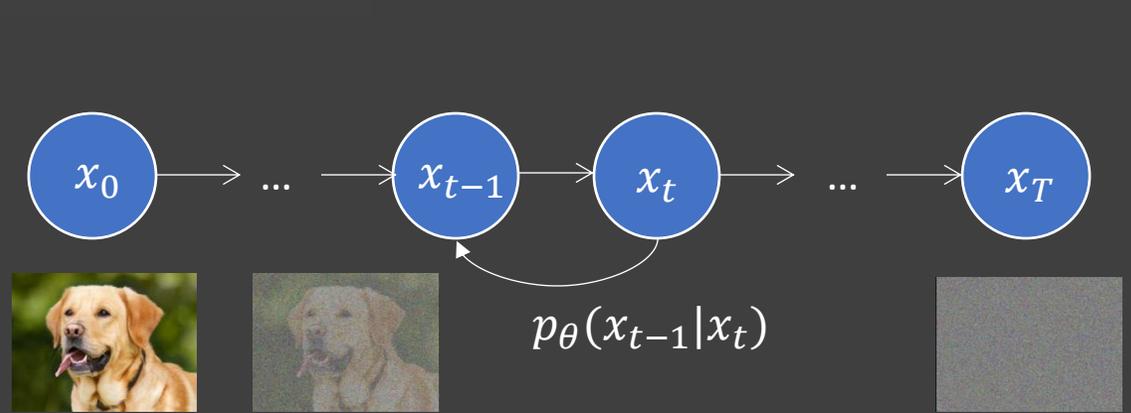
$$q(x_t|x_{t-1}) = \mathcal{N}(x_t; \mu_t = \sqrt{1 - \beta_t}x_{t-1}, \Sigma_t = \beta_t)$$

Reparameterization trick

$$x_t = \sqrt{\bar{\alpha}_t}x_0 + \sqrt{1 - \bar{\alpha}_t}\epsilon_t \quad \alpha_t = 1 - \beta_t \quad \bar{\alpha}_t = \prod_{s=0}^t \alpha_s$$

$$\epsilon_t \sim \mathcal{N}(0,1)$$

## Reverse diffusion



$$p_{\theta}(x_{t-1}|x_t) = \mathcal{N}(x_{t-1}; \mu_{\theta}(x_t, t), \Sigma_{\theta}(x_t, t))$$

Neural Network

This works because the process is Gaussian