# Agenda

- Brief intro to Artificial Intelligence
- Machine Learning and Neural Networks
- Large Language Model (LLM)
- Transformers architecture
- Top-k and temperature
- Retrieval Augmented Generation (RAG)
- Embedding and Vector Search
- Lab: LangChain, Llama 3.2, Elasticsearch



Image generated using dall-e-3

elastic

# Artificial intelligence

- Many definitions have been proposed:

  - The ability of a digital computer to perform tasks commonly associated with intelligent beings

  - The theory and development of computer systems able to perform tasks normally requiring human intelligence, such as visual perception, speech recognition, decision-making, and translation between languages

  - An umbrella term for a range of algorithm-based technologies that solve complex tasks by carrying out functions that previously required human thinking
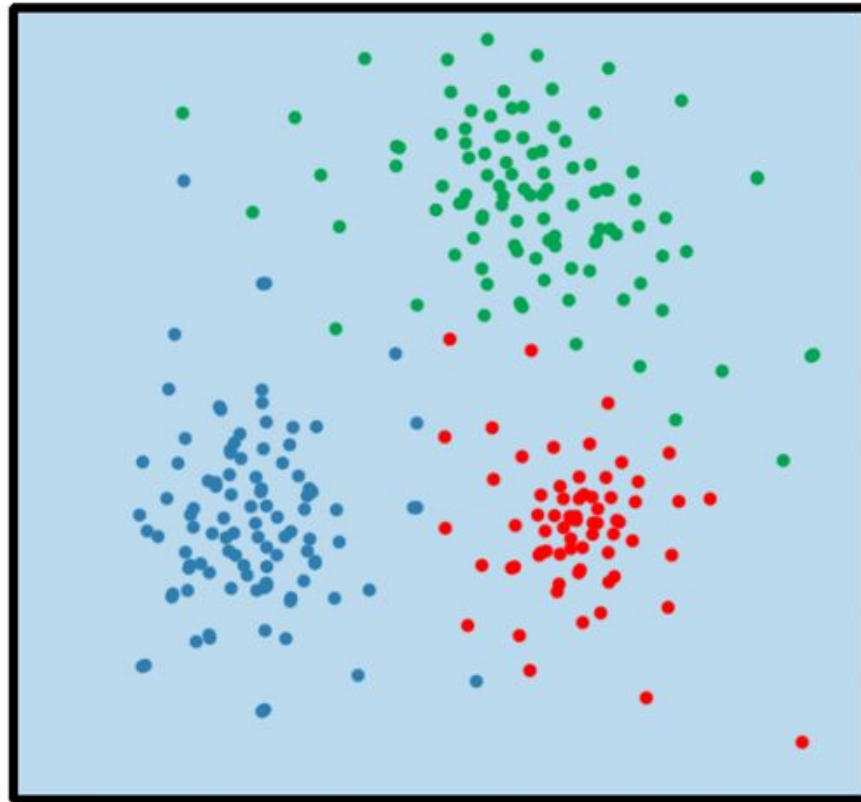
# AI examples

- February 10, 1996, [Deep Blue](#) beat Garry Kasparov in the first game of a six-game match—the first time a computer had ever beat a human in a formal chess game

- March 15, 2016, [AlphaGo](#) beat Lee Sedol 4-1 in a formal Go game

- June 2020, [Google's DeepMind A.I.](#) beats doctors in breast cancer screening trial

- June 2022, the Google [LaMDA](#) (Language Model for Dialog Applications) chatbot apparently passed the **Turing Test**. Many experts in the field, pointing out that a language model appearing to mimic human conversation does not indicate that any intelligence is present
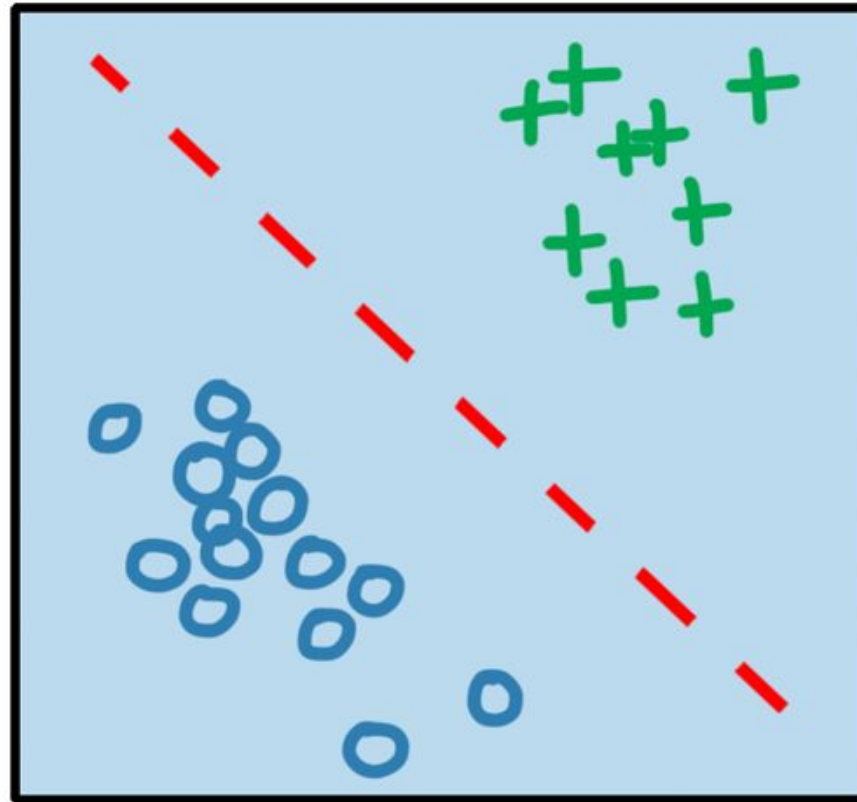
# Machine learning

- **Machine Learning (ML)** is the use and development of computer systems that are able to **learn** and **adapt** <u>without following explicit instructions</u>, by using algorithms and **statistical models** to analyse and draw inferences from patterns in data

- We can have 3 types of ML:

  - **Supervised learning**: use of labeled datasets to train algorithms that to classify data or predict outcomes (eg. image and speech recognition, recommendation systems, fraud detection)

  - **Unsupervised learning**: algorithms learn patterns exclusively from unlabeled data (eg. clustering, anomaly detection, preparing data for supervised learning)

  - **Reinforcement learning:** training method based on rewarding desired behaviors and punishing undesired ones
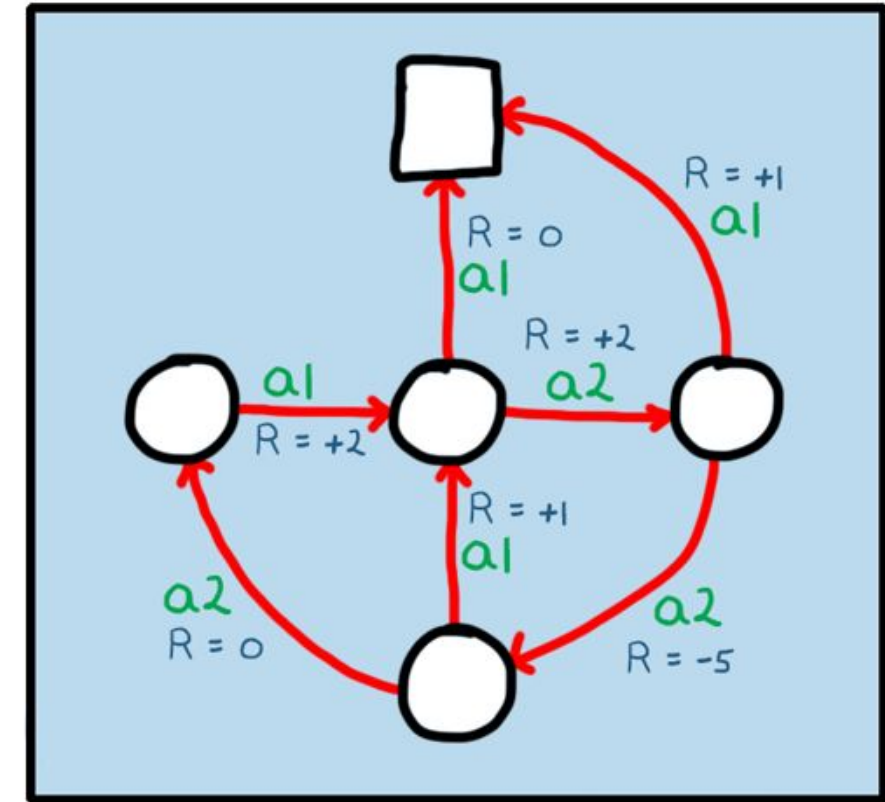
# machine learning

## unsupervised learning
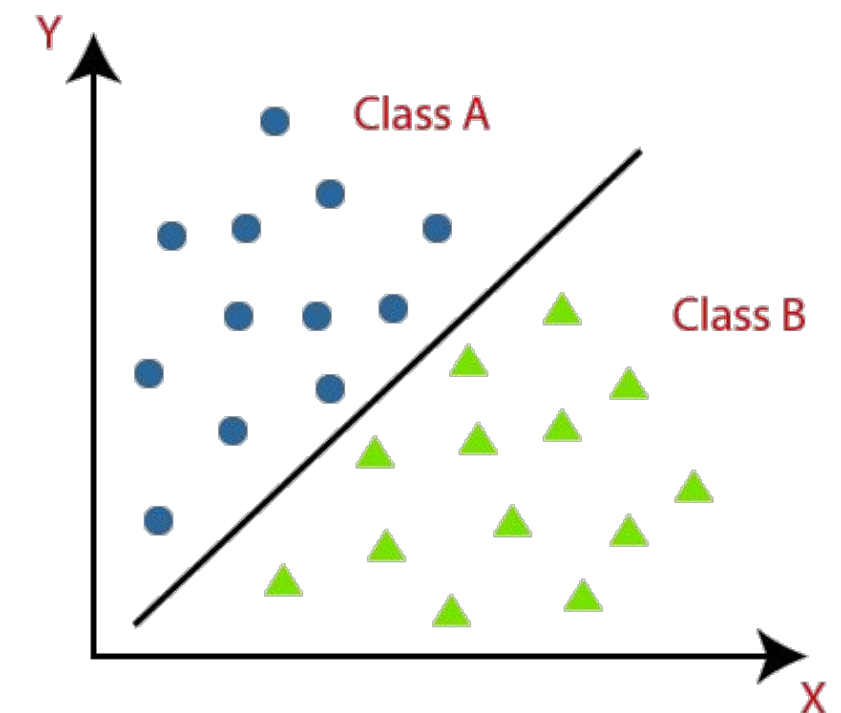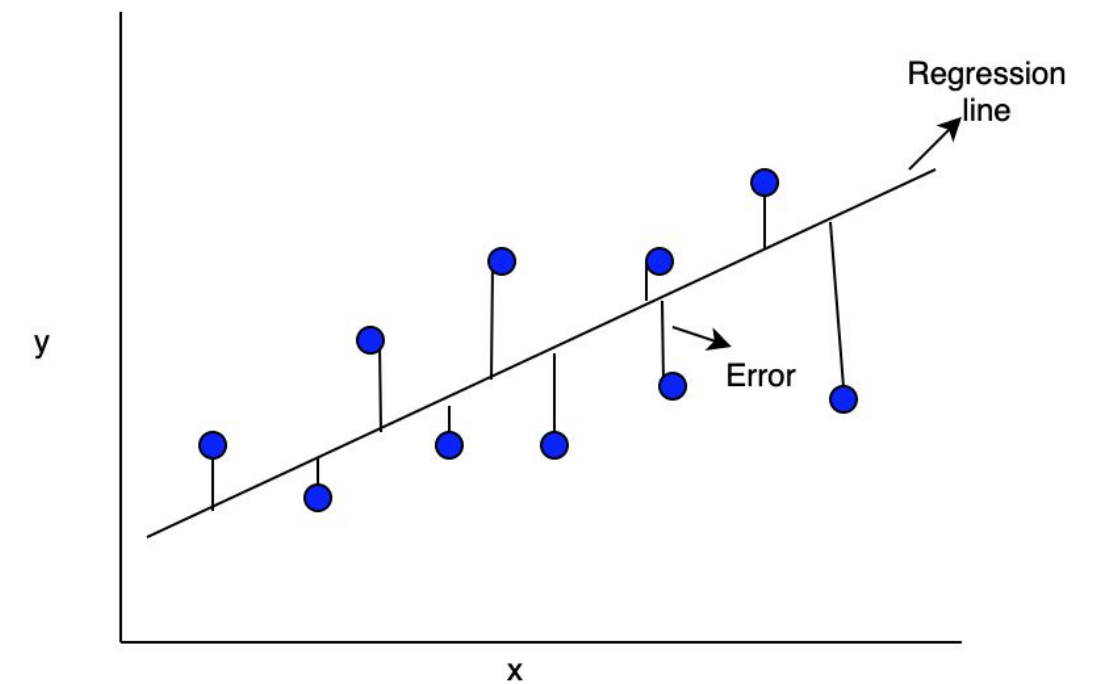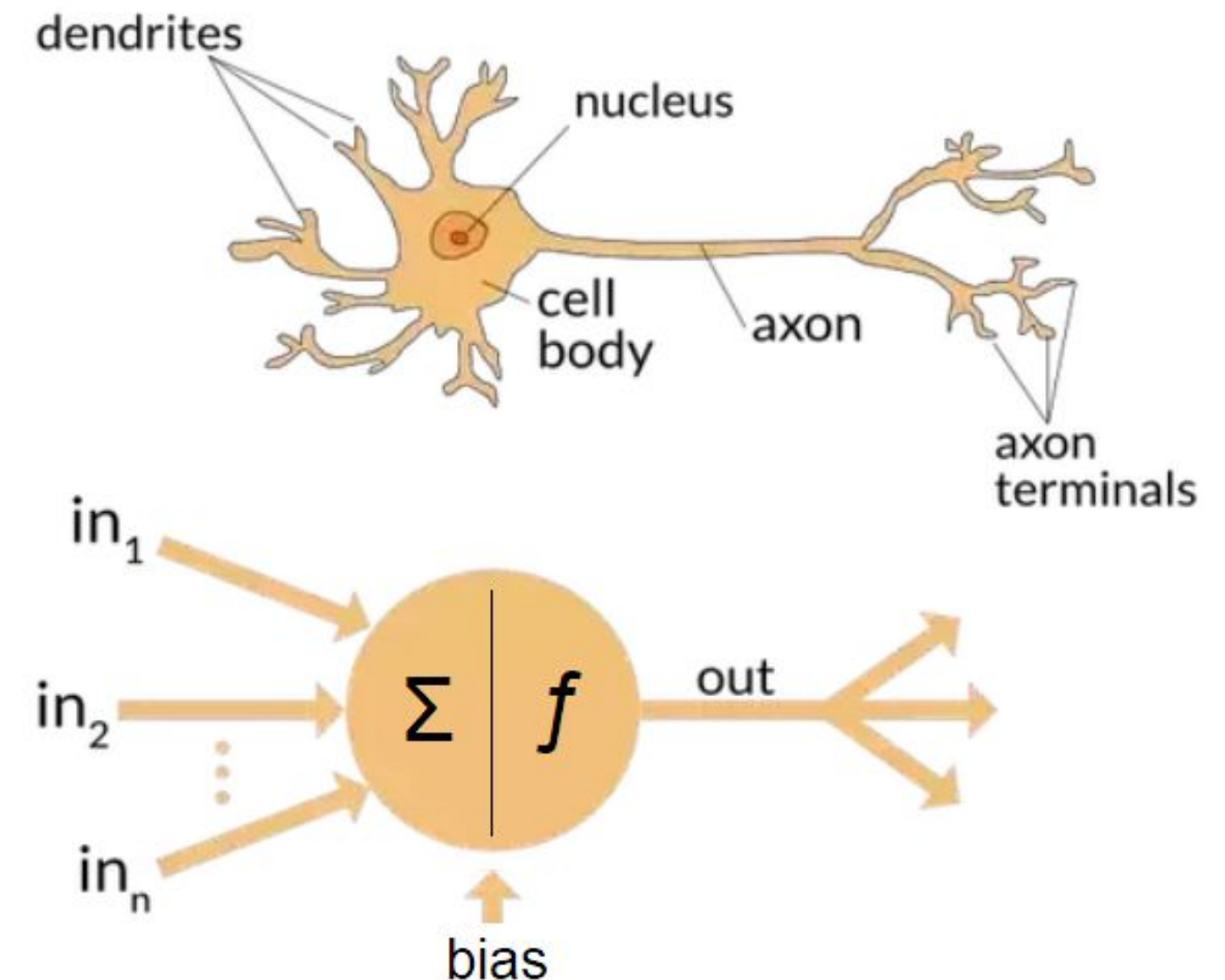
## supervised learning

## reinforcement learning

# Classification and regression

- In the supervised learning we have two types of classifications:
  - **Regression**, when the target variable is continuous. Example: predict the salary of a person based on education degree, work experience and geo location

  - **Classification**, when the target variable is discrete. Example: sentiment analysis of a piece of text (e.g. comments of a product)
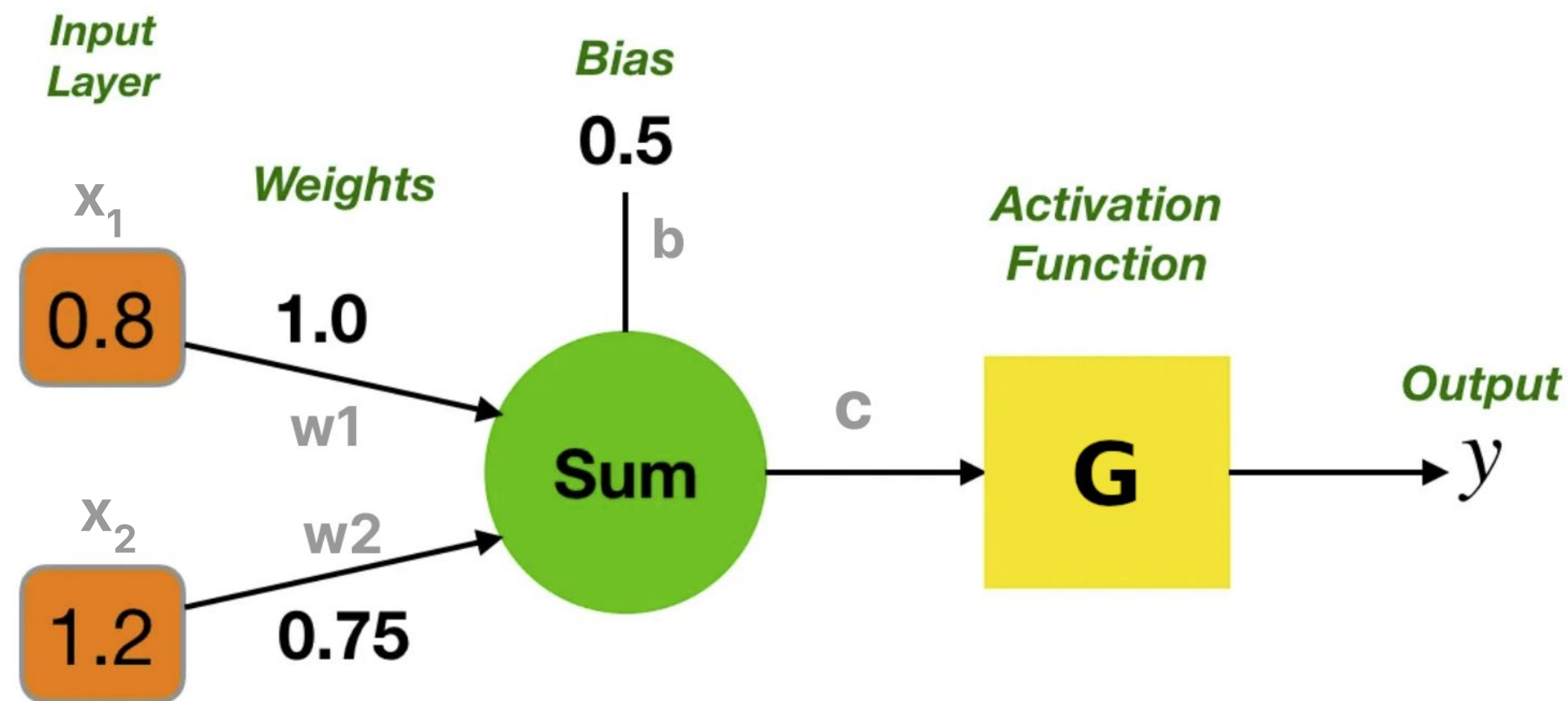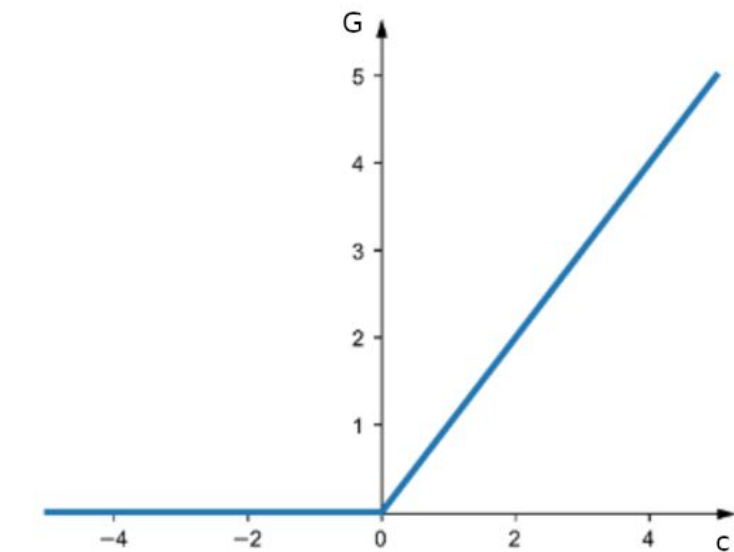
# Neural Network

- A **neural network** is a method in artificial intelligence that teaches computers to process data in a way that is inspired by the human brain

- Collection of **nodes** (artificial neurons) with inputs and outputs. A neuron computes some non-linear function of the sum of its inputs

- The nodes are collected in **layers**

- If the number of layers is greater than 3 we say **deep learning network**
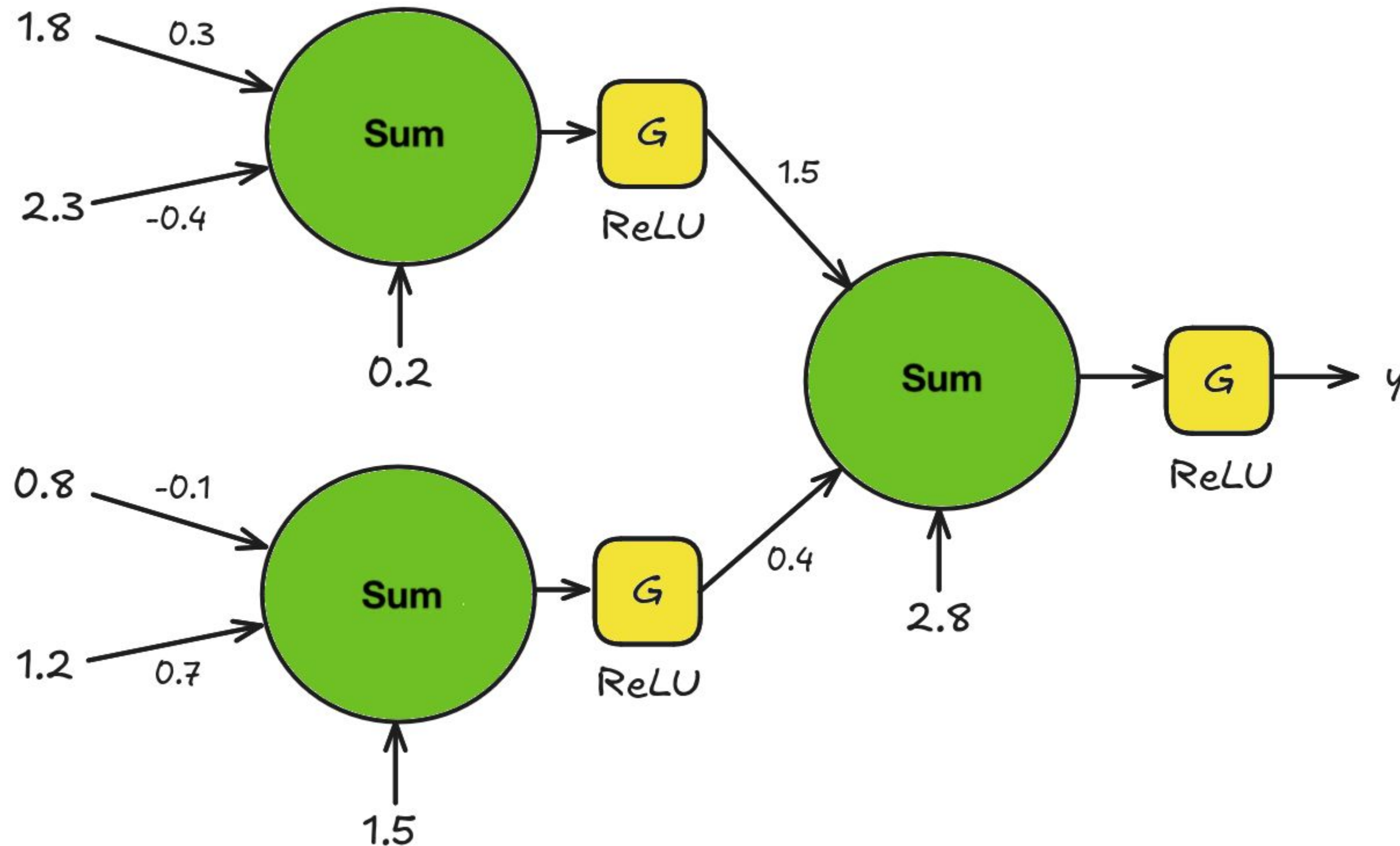
# Single layer neuron
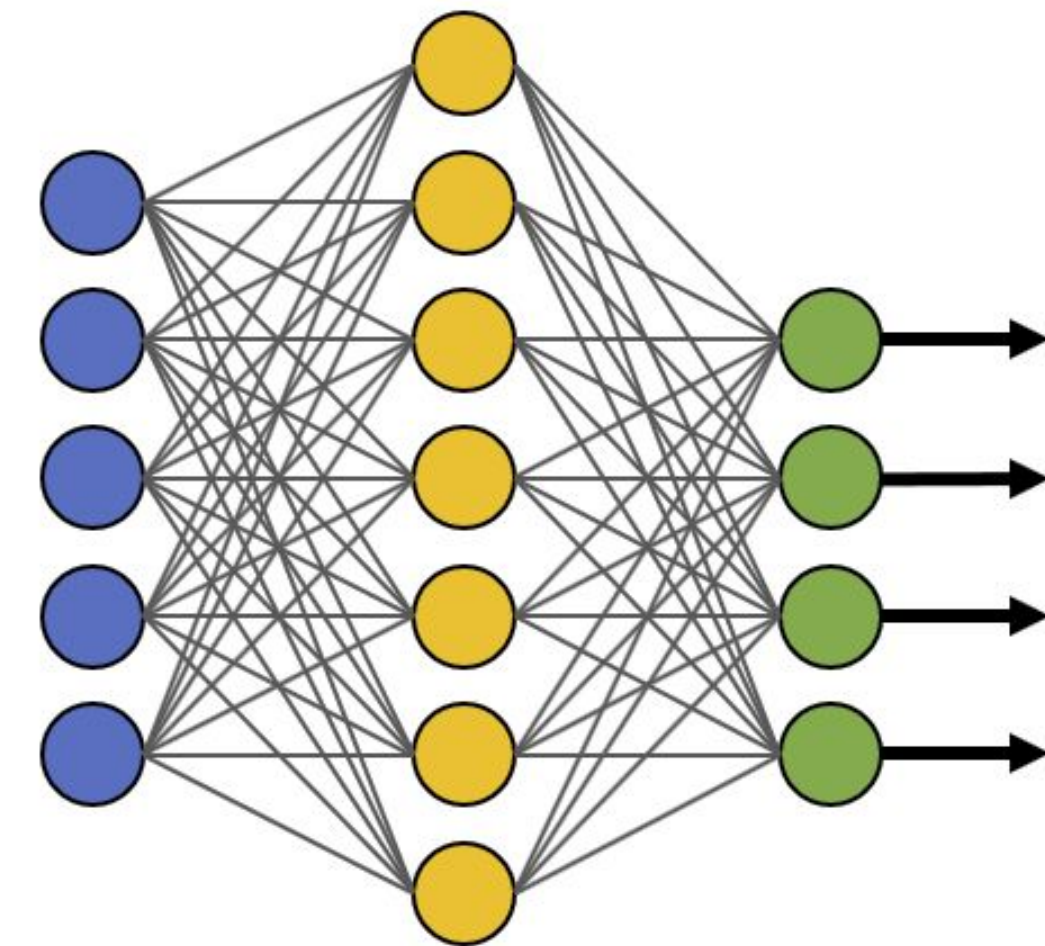
**Example: ReLU activation**

$$G(c) = c \geq 0 \ ? \ c \ : \ 0$$



$$c = x_1 * w_1 + x_2 * w_2 + b$$
$$= 0.8 * 1.0 + 1.2 * 0.75 + 0.5$$
$$= 2.2$$

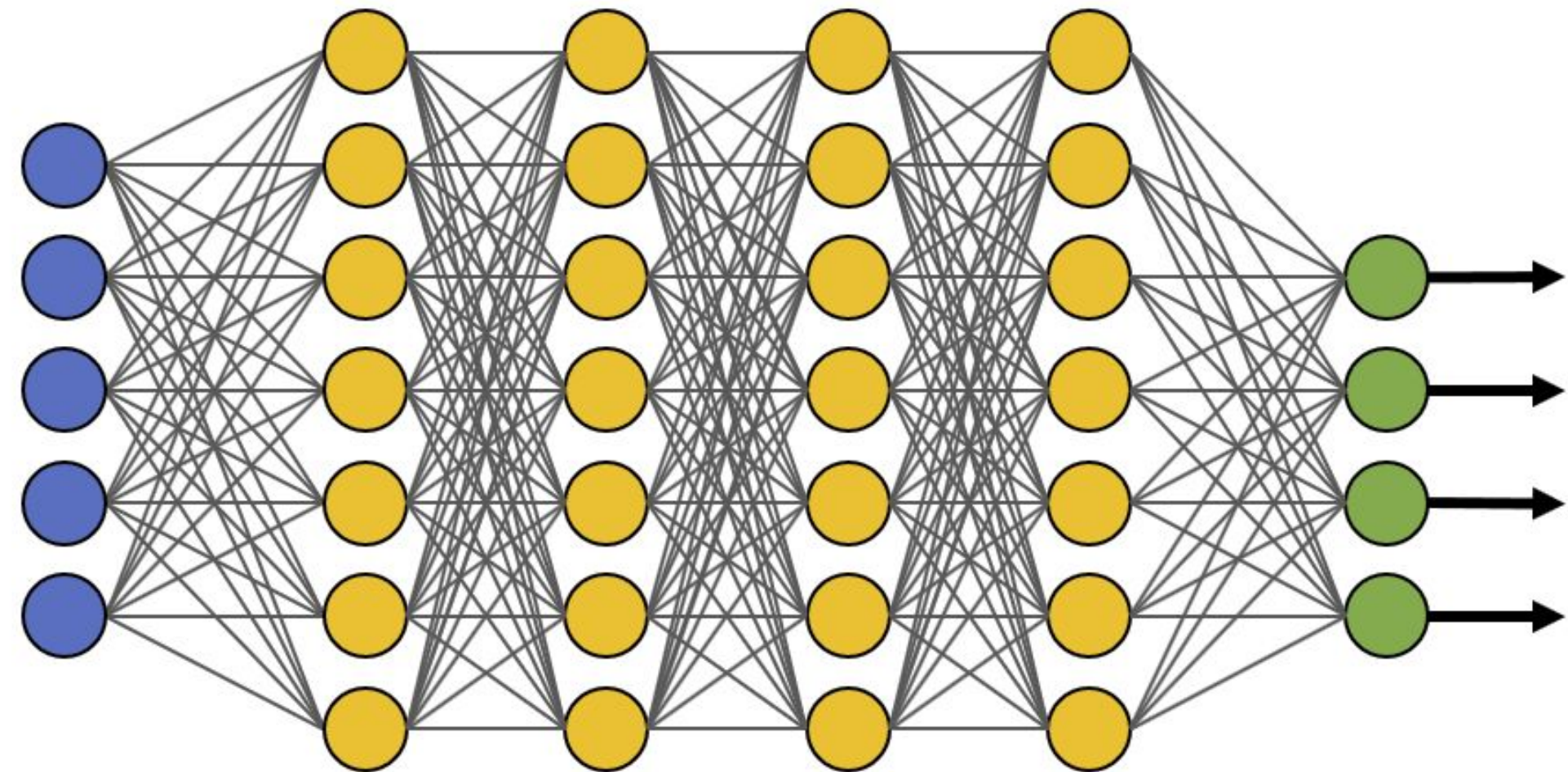$$y = G(c) = G(2.2) = 2.2$$

# Exercise: what is the value of y?

Simple Neural Network

Deep Learning Neural Network

Input Layer Hidden Layer Output Layer

# Generative AI

- **Generative Artificial Intelligence** (GenAI) is artificial intelligence capable of generating text, images, or other media, using generative models

- GenAI models **learn the patterns and structure** of their input training data and then generate new data that has **similar characteristics**

- It's used in many industries, including art, writing, script writing, software development, product design, healthcare, finance, gaming, marketing, and fashion

- The [GenAI market size](#) has been valued at $36 billion in 2024, and is projected to reach $191.8 billion by 2032
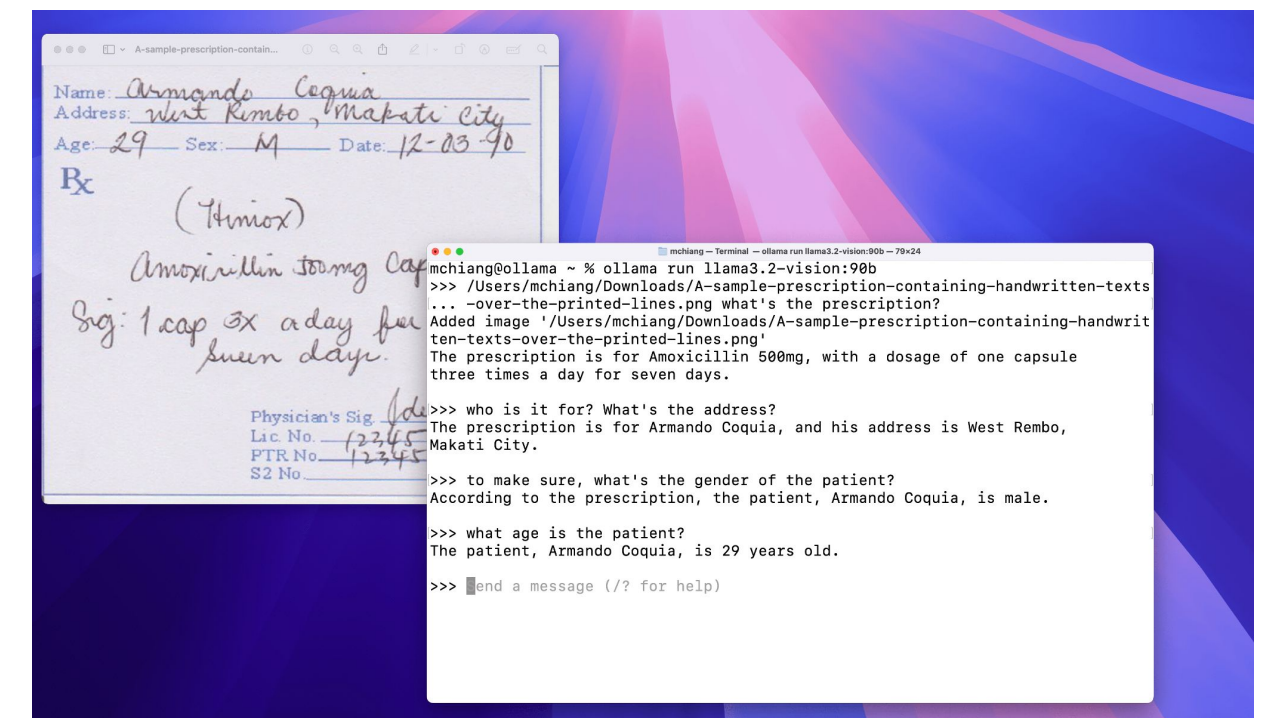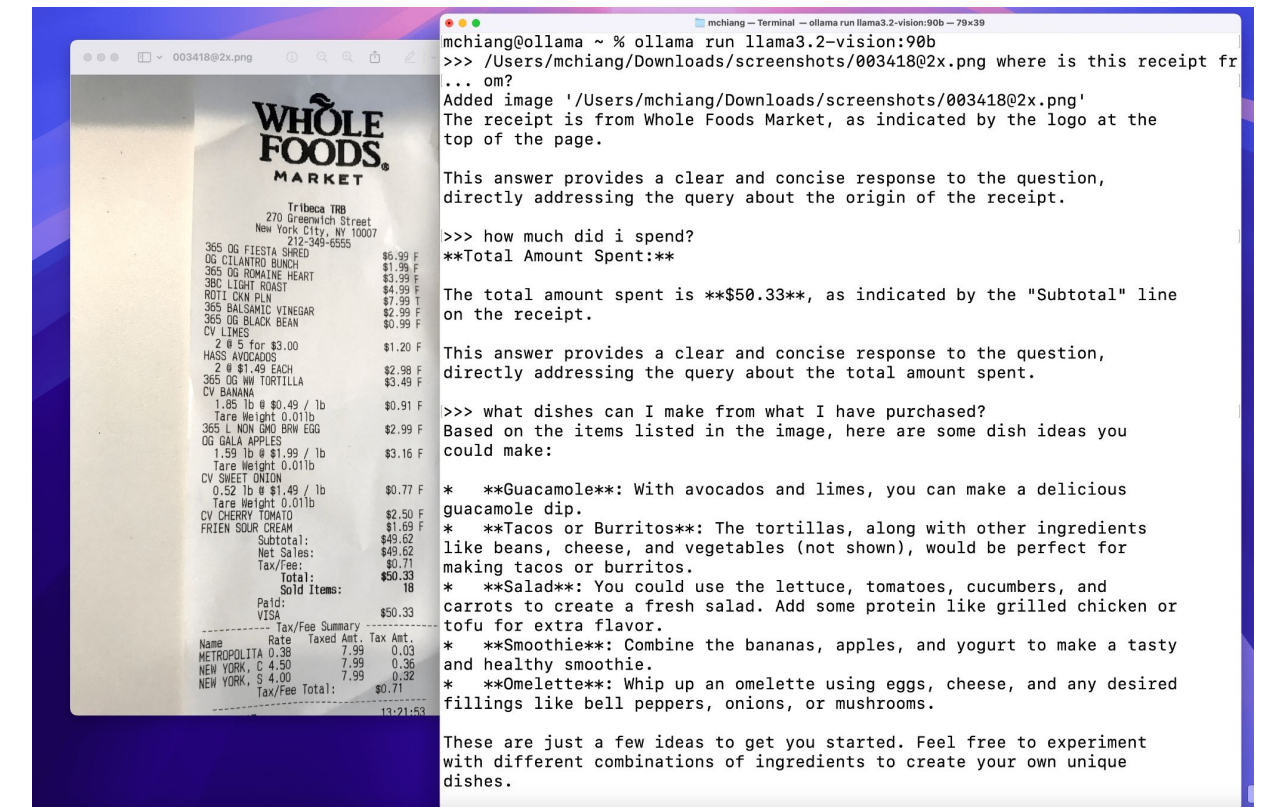
# Examples



**Prompt for GPT-4o**
Create a picture of an astronaut programming a computer on the Moon

Audio file generated using tts-1

# Llama3.2-vision

- A multimodal small (11B) and medium-sized (90B) vision LLM by Meta

- Support image reasoning, such as document-level understanding including charts and graphs, captioning of images, and more

- Released Sept. 25, 2024

- A quick [DEMO](#) using ollama

# AI ⊃ ML ⊃ DL ⊃ GenAI

## Artificial Intelligence

The ability of a machine to show human ability like reasoning, learning, such as creativity.

### Machine Learning

The set of algorithms that make intelligent machines capable of improving with time and experience.

#### Deep Learning

A type of ML based on *deep* neural networks made of multiple layers of processing.

##### Generative AI

# State of AI 2024

- [stateof.ai](stateof.ai) is a very good report about AI in Research, Industry, Politics, Safety and Predictions

# Large Language Model (LLM)

# RNN, before LLM

- **Recurrent Neural Networks** (RNN)

- Prediction of the next words based on the previous words

- RNN does not scale

- To complete a sentence the model needs to understand the structure of the entire sentence

- Eg. "The teacher taught the students with the book"

  - Did the teacher teach using the book?

  - Did the student have the book?

  - Or was it both?

# RNN backpropagation

# LLM

- **Large Language Model** (LLM) are probabilistic models that produce sentence in natural language

- These models work by completing sentences

AI is transforming   →   LLM   →   AI is transforming the way we work

(PROMPT)

elastic

# Transformer architecture

- Introduced in [Attention is All You Need](#) paper in 2017

- Basement of all LLMs

- The sentences are analyzed using a **self-attention** mechanism: each part of a sentence is evaluated in relation to every other part to understand contextual relationships and assign appropriate weights

# RNN vs Transformers

**RNN**



The teacher taught the student with the book.

**Transformers**



The teacher taught the student with the book.

# Attention map



eg. **book** is strongly connected with **teacher** and **student**

**self-attention**

# LLM

- **Large Language Model** (LLM) consisting of a neural network with many parameters (typically billions of weights or more), trained on large quantities of unlabelled text using self-supervised learning

- A message is splitted in **tokens**

- Each token is translated in a number using an operation called **embeddings**

- LLM **repeatedly predicting** the next token

elastic

# Size of GPT-4

- Around **1.76 trillion** parameters

- Neural network with **120** layers

- Process up to **25,000** words at once

- Estimated training cost is $200M using 10,000 [Nvidia A100 GPU](#) for 11 months

GPT-3                    GPT-4

175.000.000.000          1.000.000.000.000.00

# Transformers architecture

# Tokenizer

- We need to convert a sentence in numbers using a **tokenizer**

# Embedding



Output

e.g. 512

| 342 | 879 | 432 | 342 |

Embedding    Embedding    Embedding

Inputs

**Example with 3 dimensions**

fire

fox

student

book

internet

computer

Angle measures distance between words

# Positional encoding

# Self-attention

# Feed forward network

It helps refine the token representations
learned from self-attention

# Softmax



Normalized to a probability score

$$\sigma(\mathbf{z})_i = \frac{e^{z_i}}{\sum_{j=1}^{K} e^{z_j}} \quad \text{for } i = 1, \ldots, K \text{ and } \mathbf{z} = (z_1, \ldots, z_K) \in \mathbb{R}^K.$$

# Transformer Explainer



https://poloclub.github.io/transformer-explainer/

# Predict the next word

Softmax output

| prob | word |
|------|------|
| 0.20 | cake |
| 0.10 | donut |
| 0.02 | banana |
| 0.01 | apple |
| … | … |

Choose the one with greatest probability (greedy algorithm)

elastic

# Top-k

| prob | word |
|------|------|
| 0.20 | cake |
| 0.10 | donut |
| 0.02 | banana |
| 0.01 | apple |
| ... | ... |

Softmax output

*k=3*

**top-k**: select an output from the top-k results after applying random-weighted strategy using the probabilities

elastic

# Temperature



Temperature setting

## Cooler temperature (e.g <1)

| prob | word |
|------|------|
| 0.001 | apple |
| 0.002 | banana |
| 0.400 | cake |
| 0.012 | donut |
| … | … |

Strongly peaked probability distribution

## Higher temperature (>1)

| prob | word |
|------|------|
| 0.040 | apple |
| 0.080 | banana |
| 0.150 | cake |
| 0.120 | donut |
| … | … |

Broader, flatter probability distribution

Softmax output

elastic

# LLM visualization



https://bbycroft.net/llm

# Prompt engineering

- You can encounter situations where the model doesn't produce the outcome that you want on the first try

- You may have to revisit the language several times to get a good answer

- The development and improvement of the prompt is known as **prompt engineering**

- One powerful strategy is to include examples of the task that you want the model to carry out inside the prompt

- This is called **In-Context Learning (ICL)**

elastic

# ICL - zero shot inference

**Prompt**
Classify this review:
I loved this movie!
Sentiment:

LLM

**Completion**
Classify this review:
I loved this movie!
Sentiment:
Positive

elastic

# ICL - one shot inference

**Prompt**
Classify this review:
I loved this movie!
Sentiment:
Positive
Classify this review:
I don't like this chair.
Sentiment:

LLM

**Completion**
Classify this review:
I loved this movie!
Sentiment:
Positive
Classify this review:
I don't like this chair.
Sentiment:
Negative

elastic

# ICL - few shot inference

**Prompt**
Classify this review:
I loved this movie!
Sentiment:
Positive
Classify this review:
I don't like this chair.
Sentiment:
Negative
Classify this review:
This is not great.
Sentiment:

LLM

**Completion**
Classify this review:
I loved this movie!
Sentiment:
Positive
Classify this review:
I don't like this chair.
Sentiment:
Negative
Classify this review:
This is not great.
Sentiment:
Negative

elastic

# Ollama

- [Ollama](#) is a software for downloading and running LLMs locally
- Llama 3, Phi 3, Mistral, Gemma, and [other models](#)
- Simple command line tool:
  - ollama pull llama3.2:3b
  - ollama run llama3.2:3b

elastic

# Laboratory:
## LLM examples with Ollama and OpenAI



[Sorgenti Python](#)

elastic

# Retrieval-Augmented Generation (RAG)

# Retrieval-Augmented Generation (RAG)

- **RAG** is a technique in natural language processing that combines information retrieval systems with **Large Language Models** (LLM) to generate more informed and accurate responses
- It is composed by the following parts:
  - **Retrieval-Augmented**
  - **Generation**

elastic

# Generation

- LLMs are very powerful but have some limitations:
  - **No source** (potential hallucinations)
    - How can I verify the information coming from an LLM?
    - What sources has been used to generate the answer?
  - **Out of date**
    - An LLM is trained in a period of time
    - For update we need to retraining the model (very expensive)

elastic

# Retrieval-Augmented

- We collect sets of private or public document
- We build a **retrieval system** (e.g. a database) to extract a subset of documents using a **question**
- Then we pass the **question + documents found** to an LLM as prompt with a context
- The LLM can give an answer using the updated documents

elastic

# RAG architecture

# Retrieve documents from a question

- How we can retrieve documents in a database using a question?
- We need to use **semantic search**
- One solution is to use a **vector database**
- A vector database is a system that uses **vectors** (set of numbers) to retrieve information

elastic

# What is a vector?

- A vector is a set of numbers
- Example: a vector of 3 elements [2, 5, -10]
- A vector can be represented in a multi-dimensional space (eg. Llama3.2 uses 3072 dimensions)



elastic

# Similarity between two vectors

- Two vectors are (semantically) similar if they are close to each other
- We need to define a way to measure the similarity

**Cosine Distance**

$$1 - \frac{A \cdot B}{||A|| \quad ||B||}$$

**Squared Euclidean (L2 Squared)**

$$\sum_{i=1}^{n} (x_i - y_i)^2$$

**Dot Product**

$$A \cdot B = \sum_{i=1}^{n} A_i B_i$$

**Manhattan (L1)**

$$\sum_{i=1}^{n} |x_i - y_i|$$

elastic

# Embedding

- Embedding is the translation of an input (document, image, sound, movie, etc) to a vector
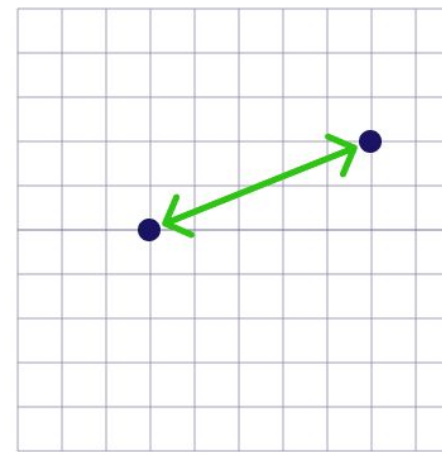- There are many techniques, using an LLM typically this is done by a neural network
- The goal is to group information that are semantically related to each other
- https://projector.tensorflow.org/



### Words As Vectors

# Vector database + LLM

- The search query (**question**) is in natural language
- We use semantic search to retrieve top-n relevant documents (**context**)
- We send the following prompt to the LLM (example):
  - *Given the following **{context}** answer to the following **{question}***

elastic

# Split the documents in chunk

- We need to store data in the vector database using chunk of information
- We cannot use big documents since we need to pass it in the context part of the prompt for an LLM that typically has a token limit (e.g. Llama3.2 up to 128K)
- We need to split the documents in **chunk** (part of words)

elastic

# Elasticsearch (vector database)

- [Elasticsearch](#) is Free and Open Source ([AGPL](#)), Distributed, RESTful Search Engine
- Distributed search and analytics engine, scalable data store and **vector database** optimized for speed and relevance on production-scale workloads.
- You can run it locally with a single command:
    - **curl -fsSL https://elastic.co/start-local | sh**

elastic

# LangChain

- LangChain is an open source composable framework to build with LLMs
- Supports all the LLMs (see here)
- Integrations with many vector databases (e.g. Chroma, Elasticsearch, Milvus, Qdrant, Redis)
- Available for Python (98K ⭐) and Javascript (13K ⭐)
- MIT license
- Other interesting projects: LangGraph (MIT license) and LangSmith (commercial)



elastic

# Laboratory:
# RAG with LangChain + Llama 3.2 + Elasticsearch

Google Colab                    Sorgenti Python

# References

- [What is retrieval-augmented generation?](#) IBM research
- Ashish Vaswan et al., [Attention Is All You Need](#), Proceedings of 31st Conference on Neural Information Processing Systems (NIPS 2017)
- Albert Ziegler, John Berryman, [A developer's guide to prompt engineering and LLMs](#), Github blog post
- Sebastian Raschka, [Build a Large Language Model (From Scratch)](#), Manning, 2024
- [Elasticsearch as vector database](#), Elastic Search Labs
- [Elasticsearch search relevance](#), Elastic Search Labs
- E.Zimuel, [Retrieval-Augmented Generation for talking with your private data using LLM](#), AI Heroes 2023 conference, Turin (Italy)
- L. Gianfagna, E. Zimuel, [Explainable AI (XAI) and Large Language Models (LLM): an impossible pairing?](#), AI Heroes 2024 conference, Turin (Italy)

elastic

# Thanks!

More information: www.elastic.co

Contact information: enrico.zimuel (at) elastic.co

elastic