

# ECON441B : Intro Machine Learning Lab

Week 7, Lecture 7 | Large Language Models

Sam Borghese

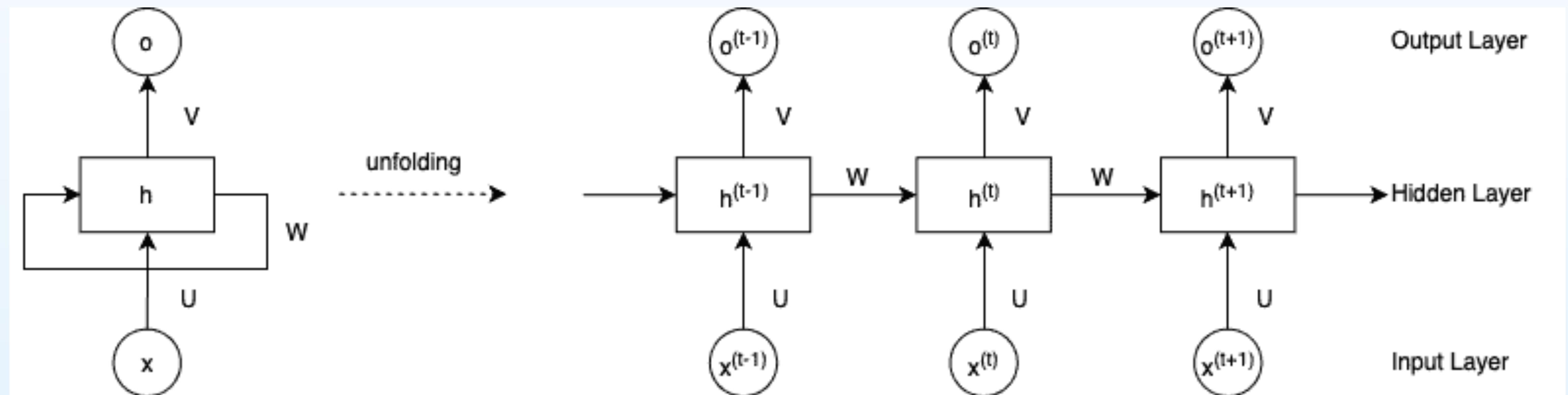
Wednesday, February 14th, 2024

1. How do Transformers work?
2. History of LLMs
3. Custom LLMs
4. OpenAI API
5. In-Class Example

# How do Transformers work?

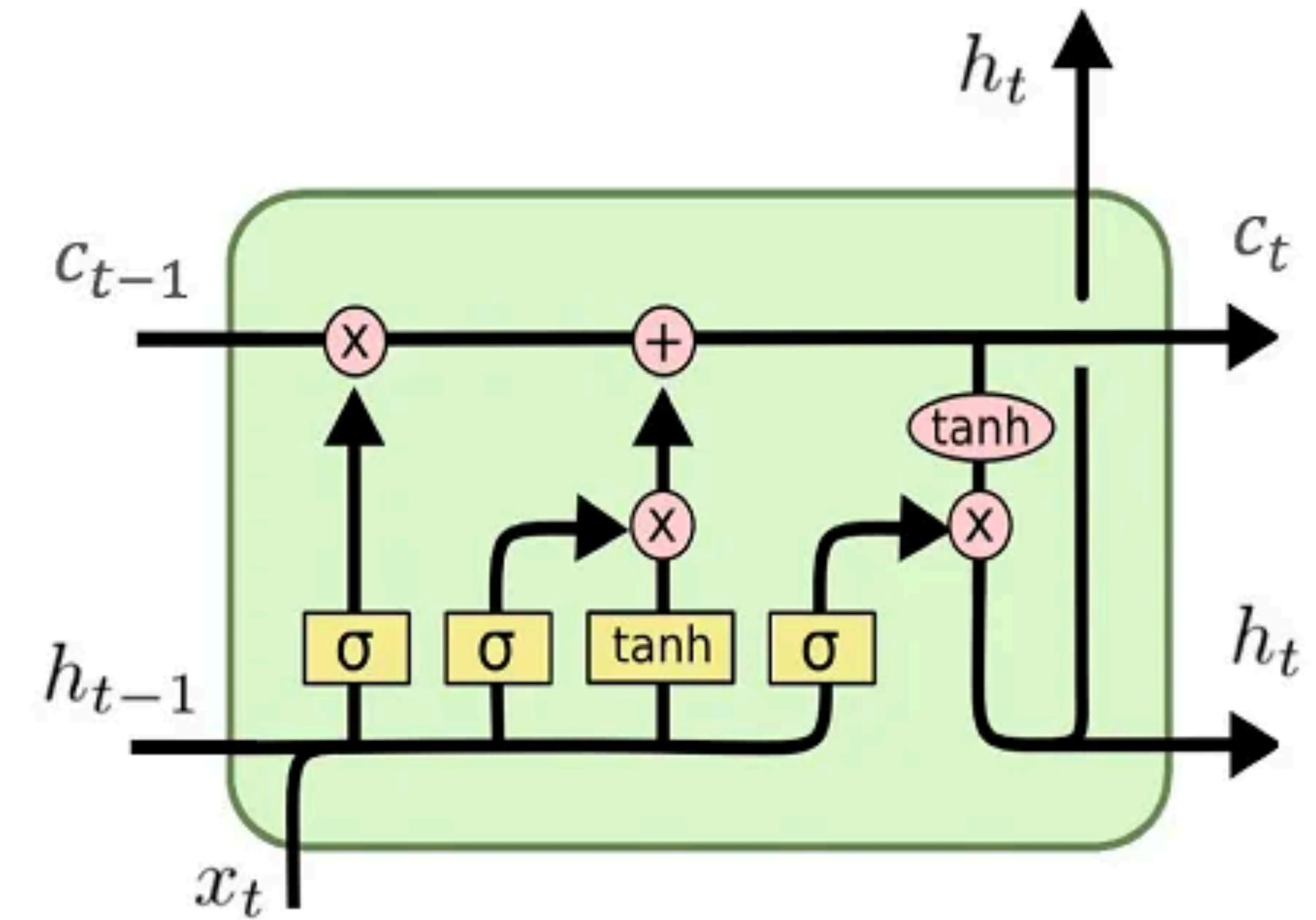
# RNN(1986)->LSTM->Transformers

RNN architecture is used as a building block for Transformers  
These cells do not necessarily have to be LSTM cells



# LSTM Cells (1997)

These now allow for information in sequential data to be stored and passed on to the next cells



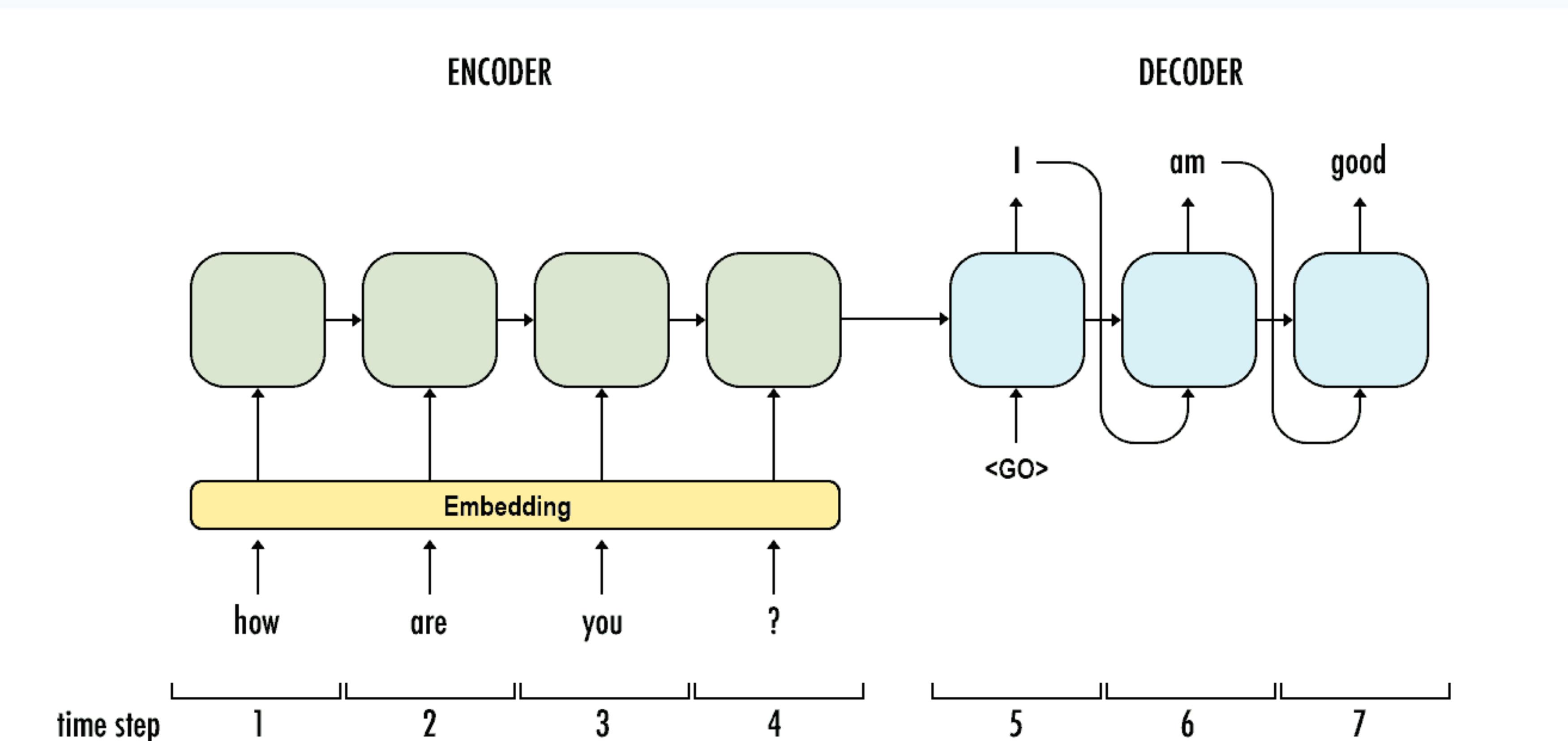
LSTM  
(Long-Short Term Memory)

# Encoder-Decoder Structure

The small difference with this structure is every output in the sequence knows what the last output was

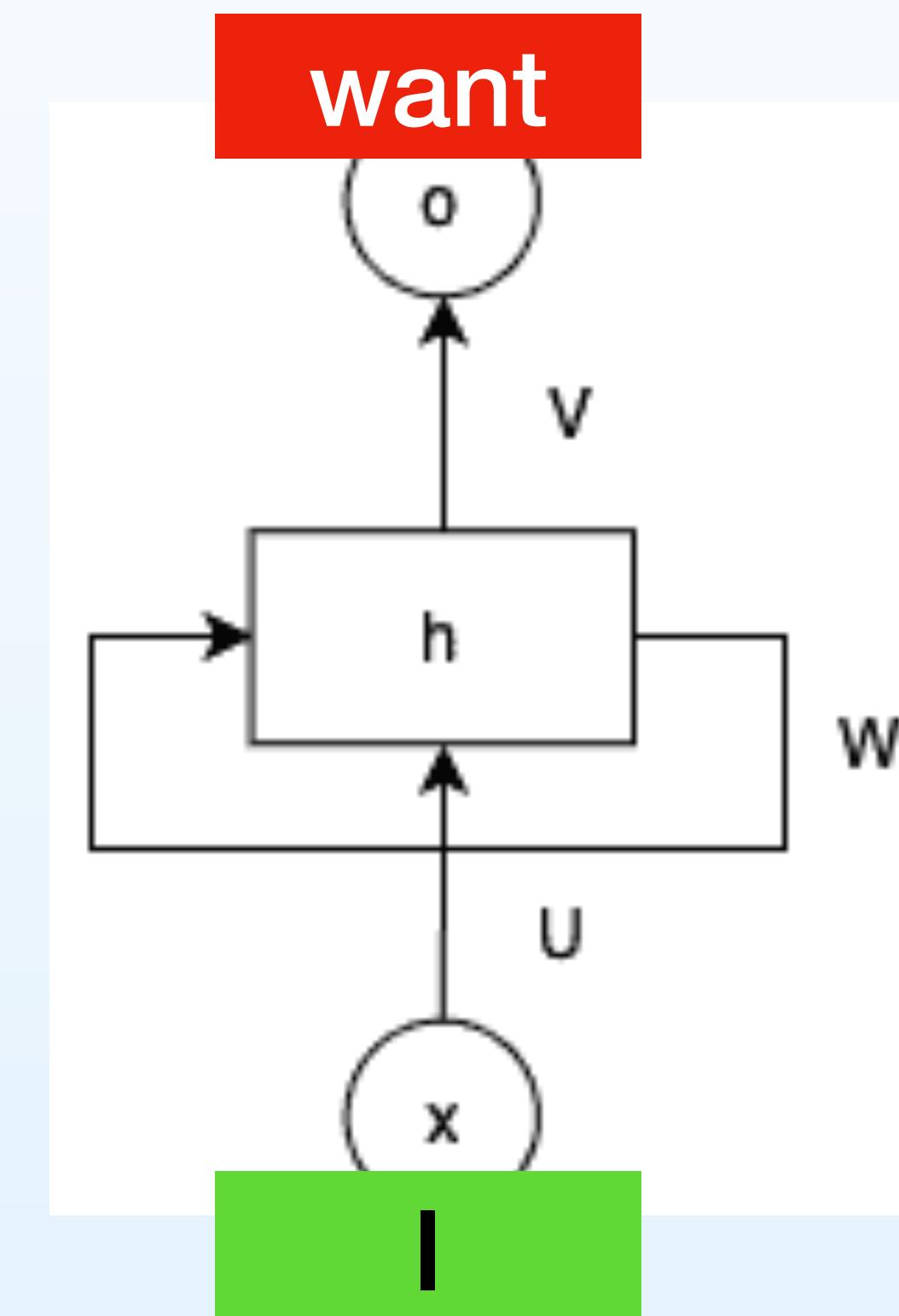
Gated  
recurrent  
units  
(2014)

Attention  
Mechanism  
(2014)



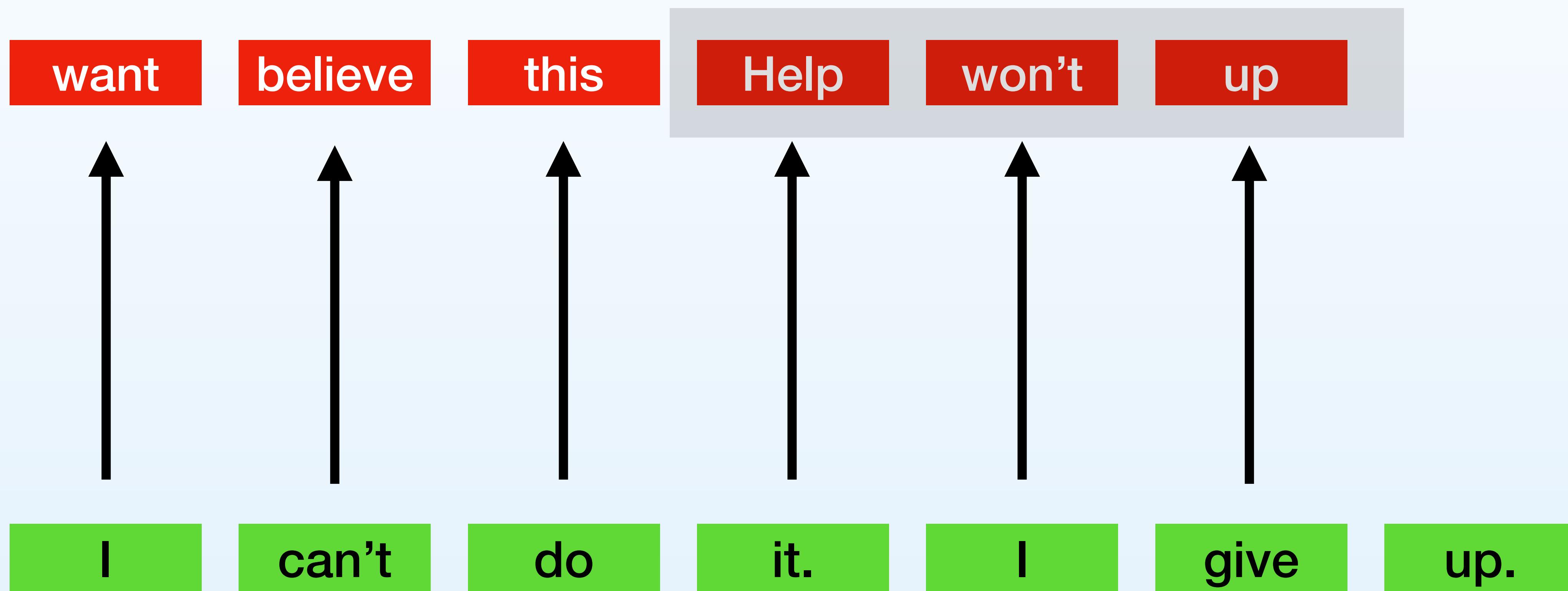
# Example

Let's say that we build an RNN to try and predict the next word that will be said in the hopes that it can make full sentences.



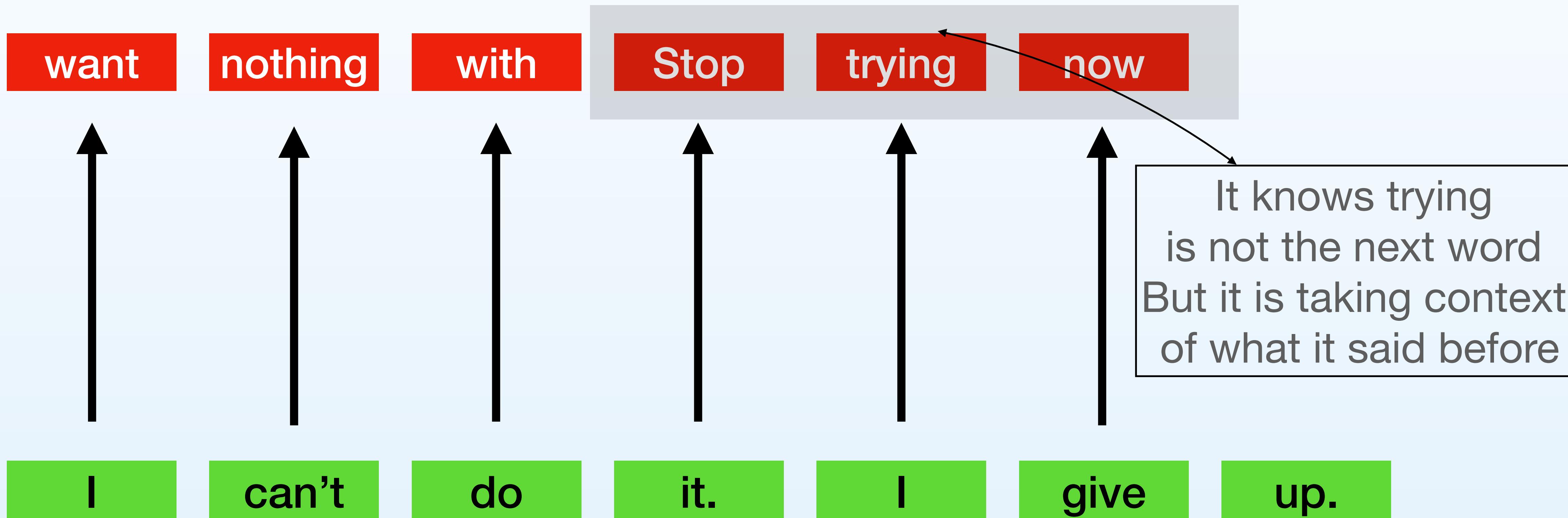
# Example - RNN

Let's say that we build an RNN to try and predict the meaning of the sentence by the next word at a time, that will be said in the hopes that it can make full sentences.



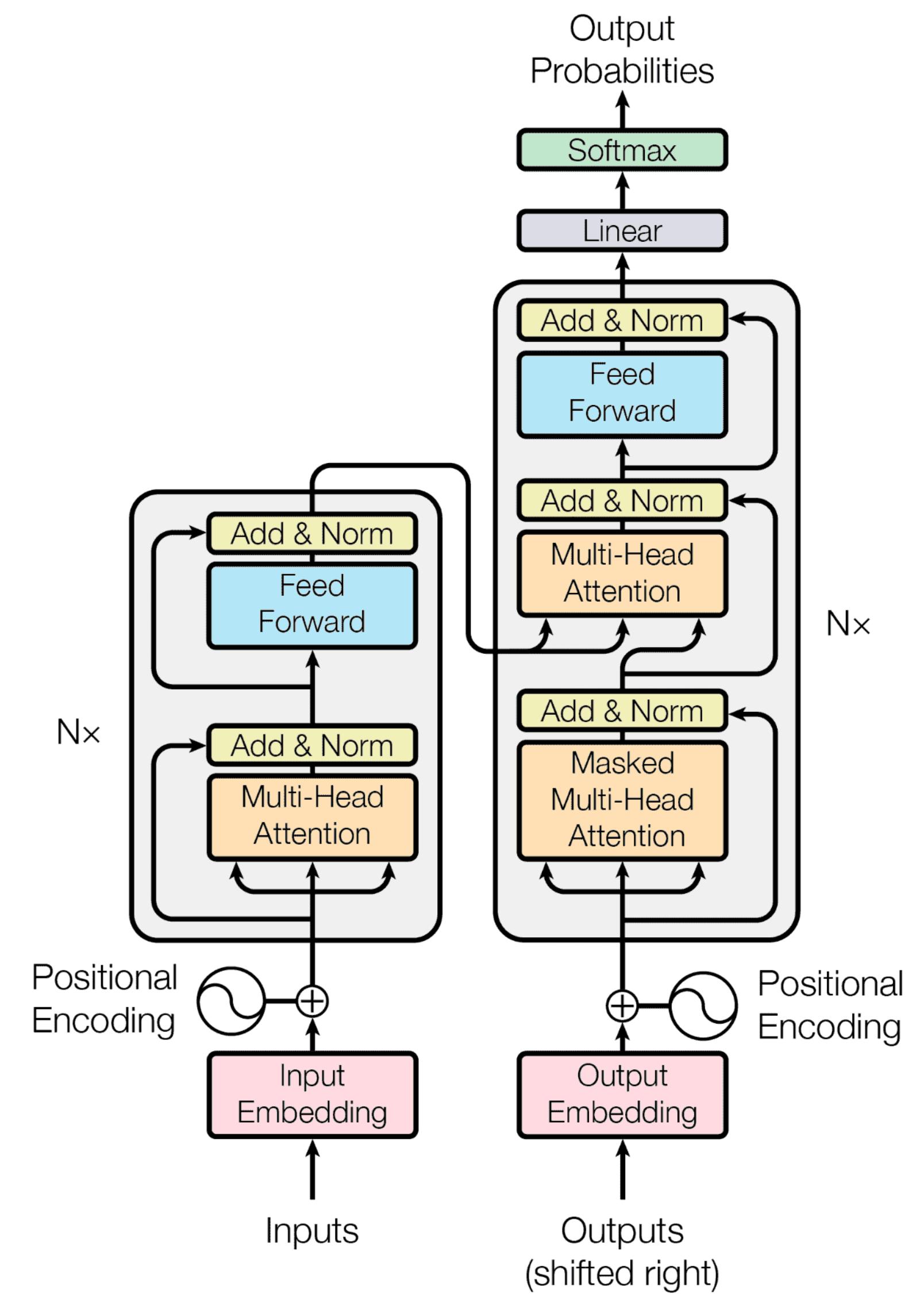
# Example - Transformer

The transformer has the context of what it said before



# Inside of a Transformer (2017)

This is called an attention model. It is a deep learning version of RNN

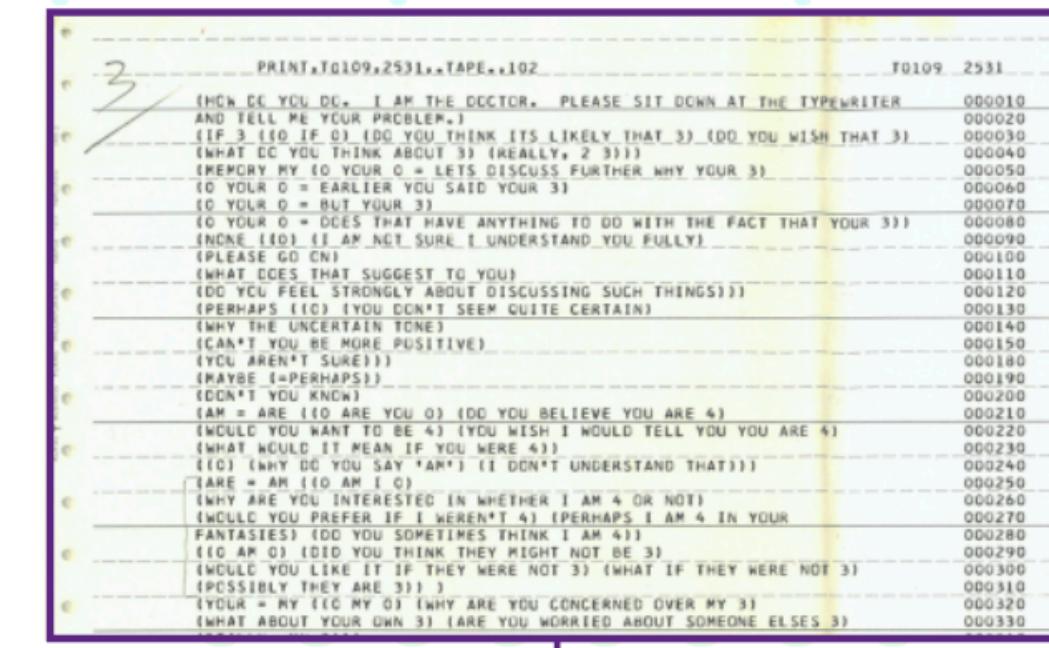


# History Of LLMs

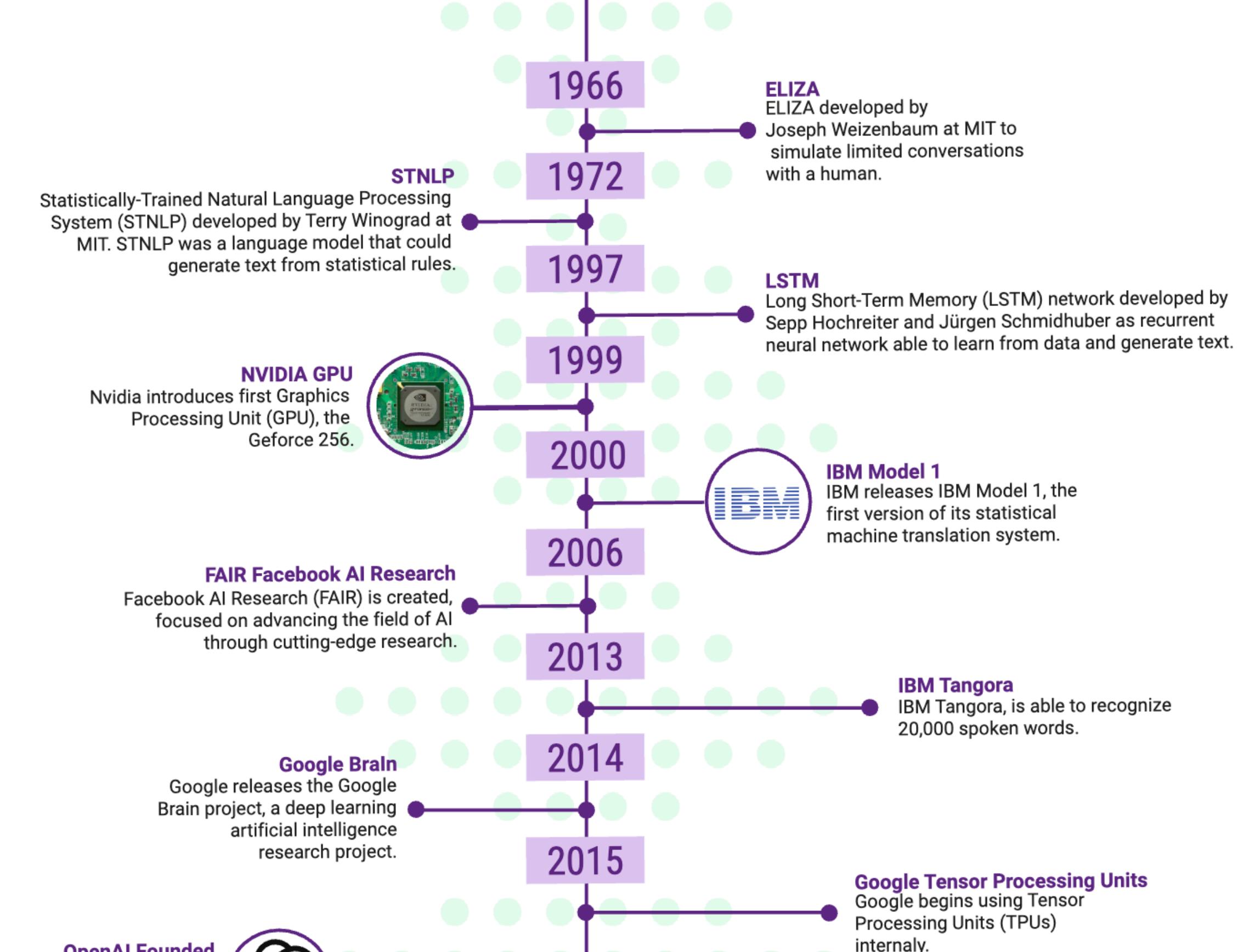
# History of Large-Language Models

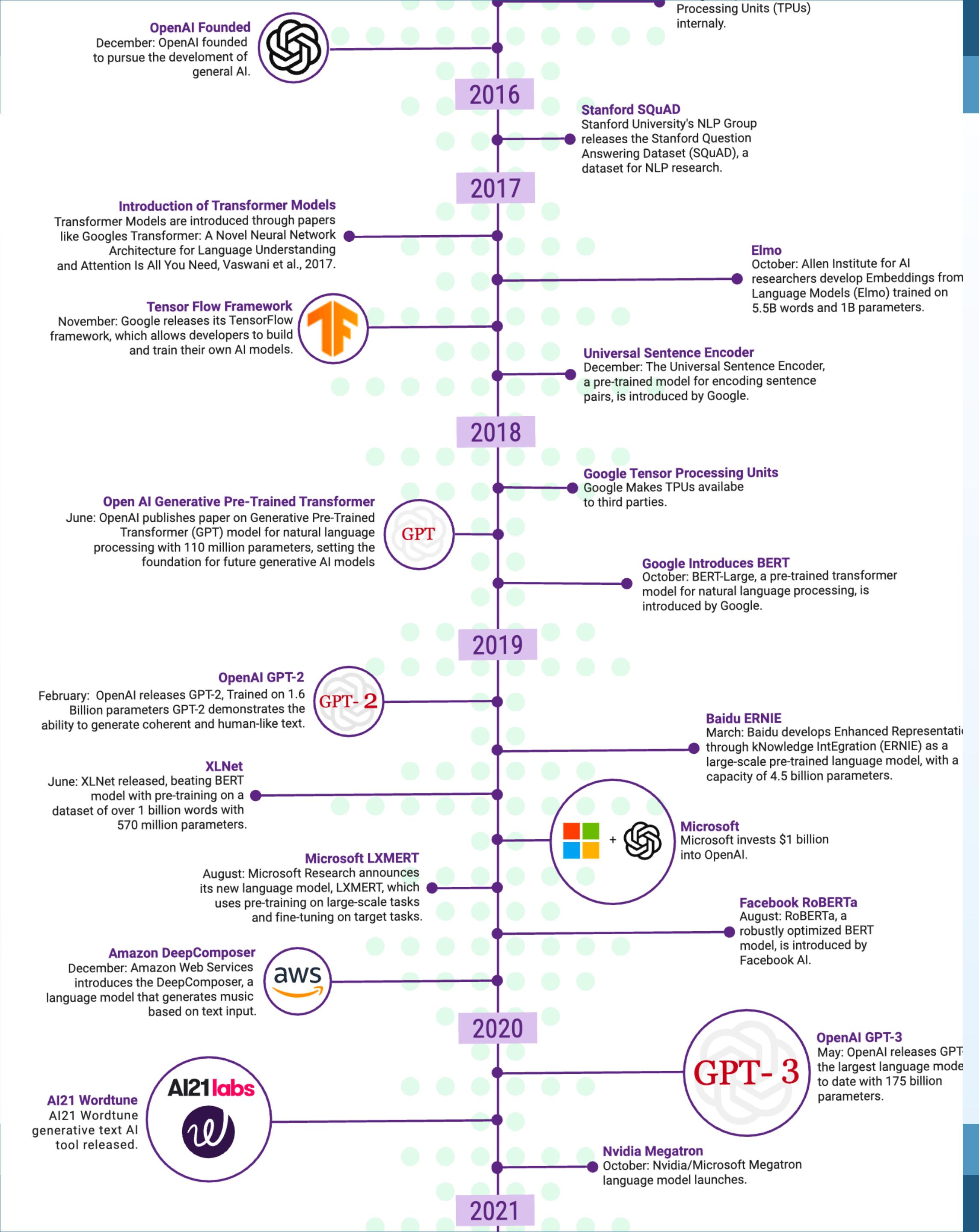
There is a difference  
between large-language  
models and Transformers

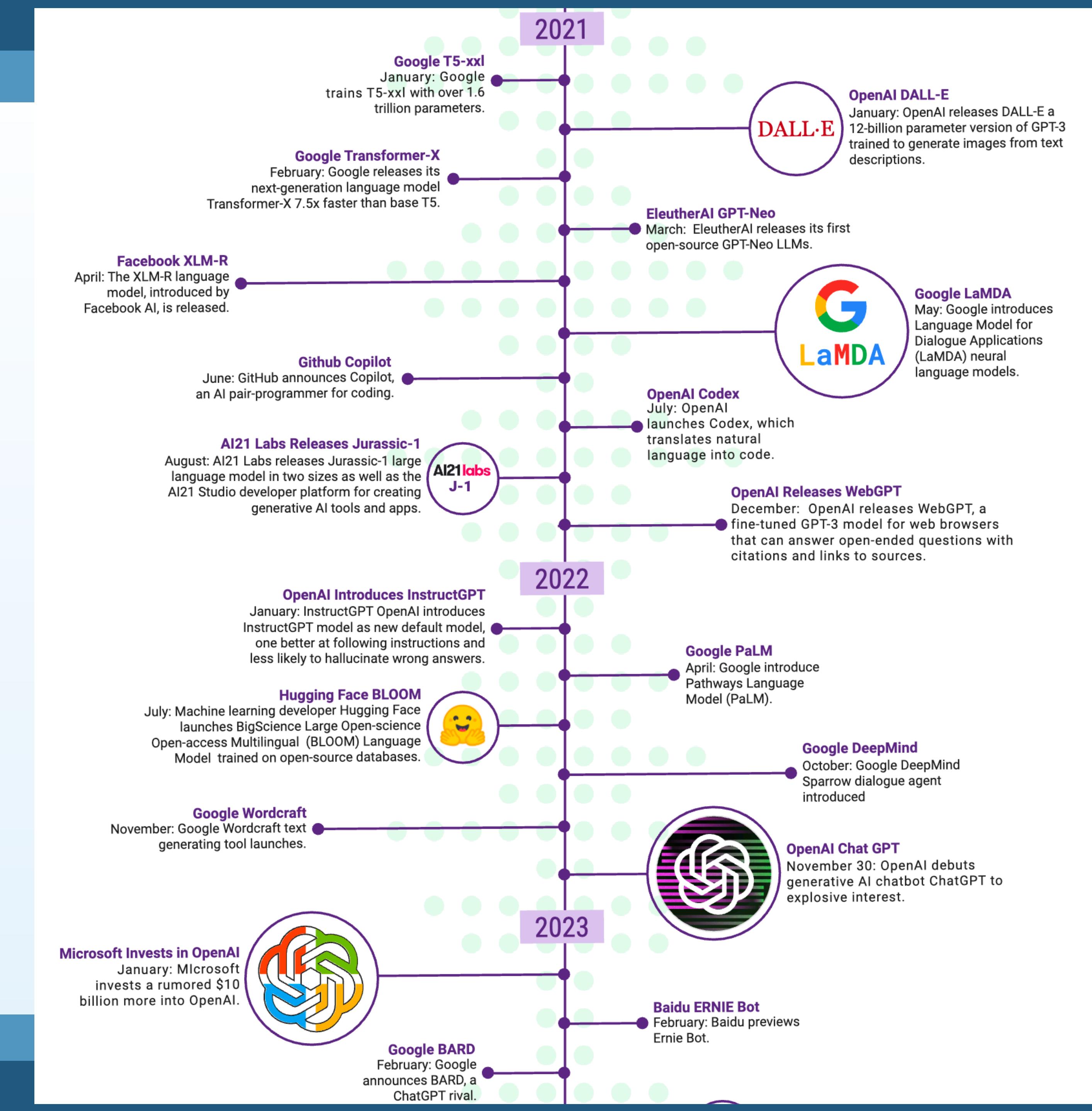
## LARGE LANGUAGE MODELS - PIVOTAL INNOVATIONS

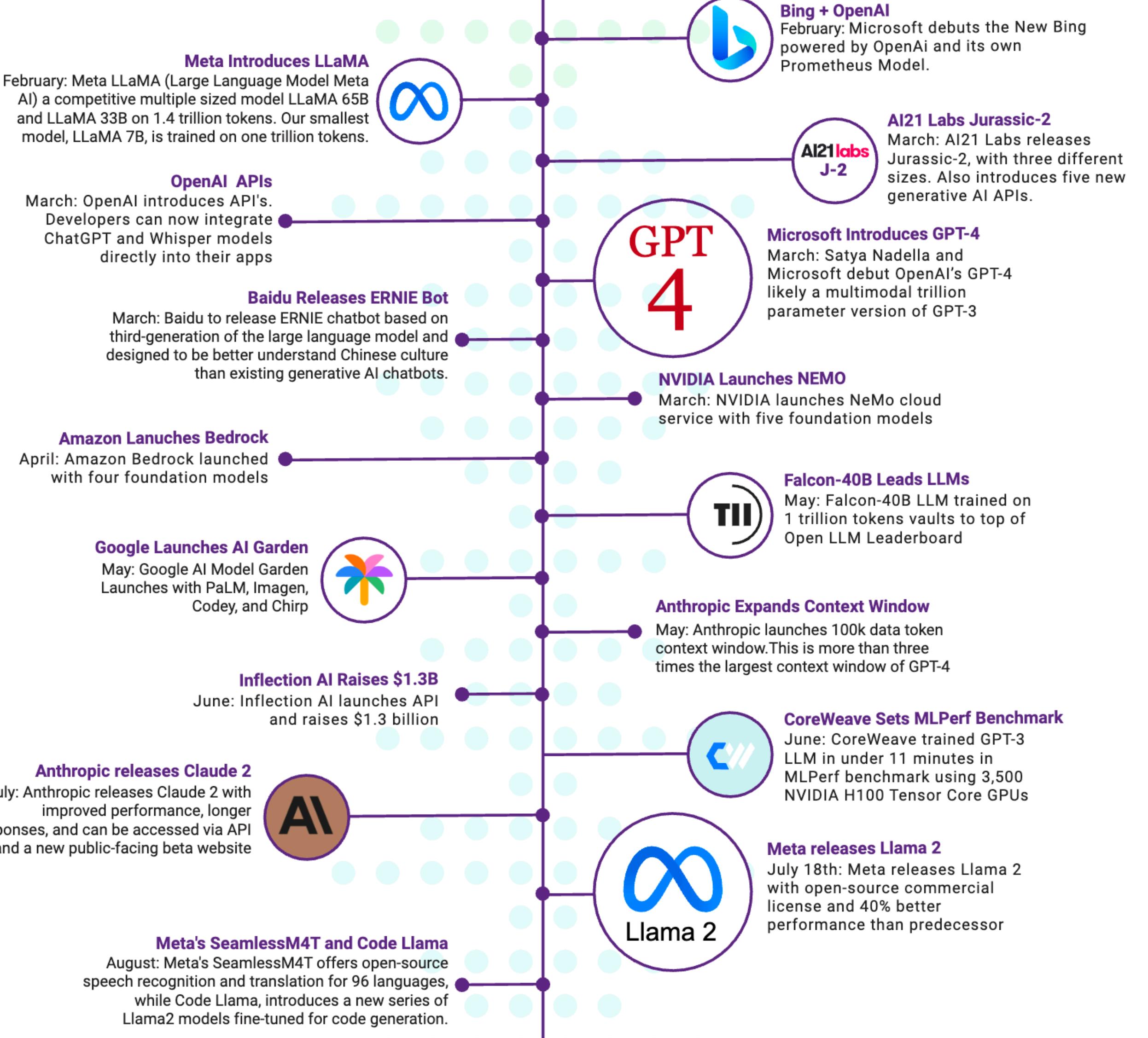


Original Eliza Script

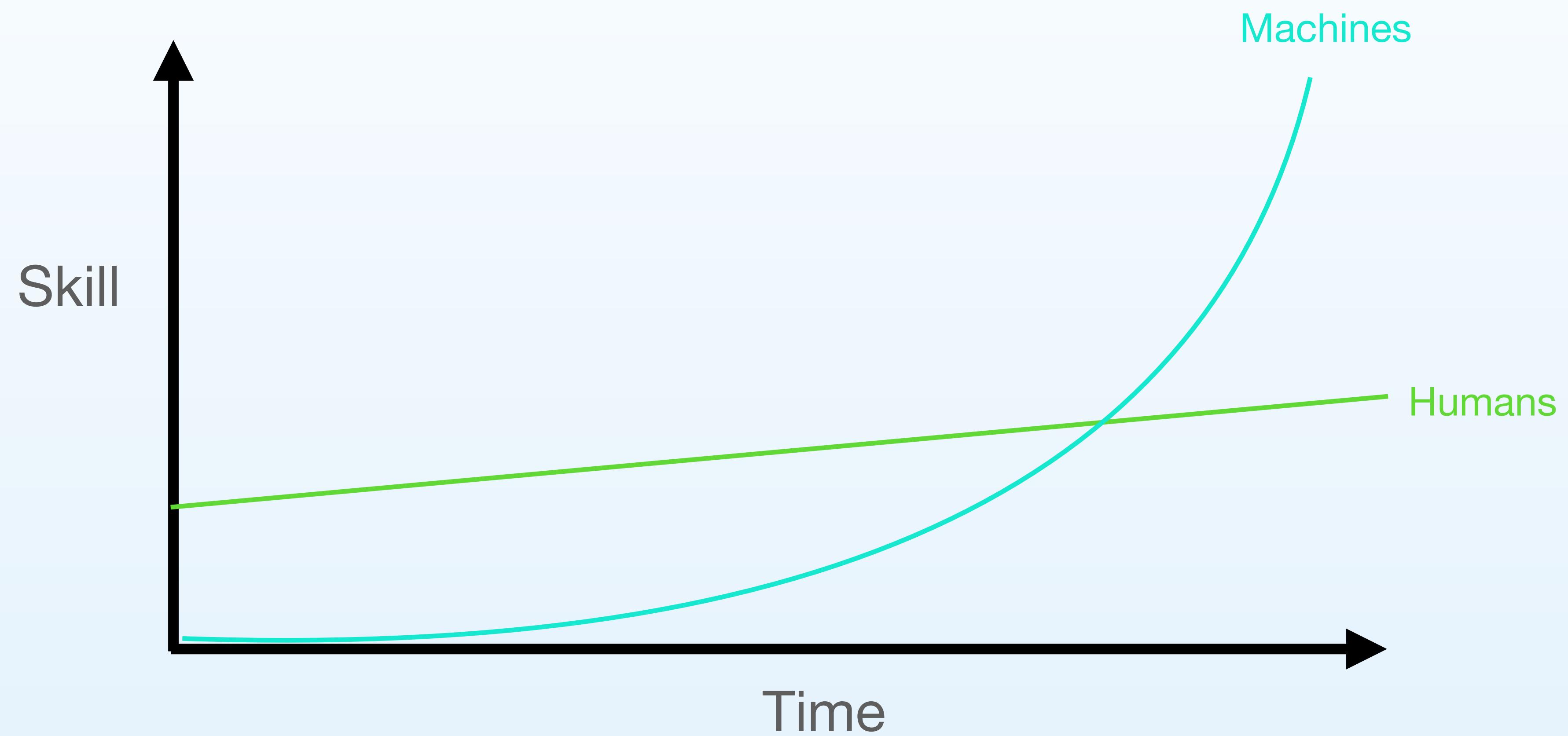








# A.I. is Accelerating but Humans are not



# Existential Conversation

Will A.I. replace jobs?

## Case Study #1 : Automated Elevators

### Elevator Automation [ edit ]

The elevator strikes were instrumental to the automation of the elevator. As elevators were a dangerous machine that could only be comfortably operated by elevator operators, manufacturers began adding safety features and allowing the elevator to run on its own.<sup>[19]</sup> New features included emergency phones, emergency stop buttons, and alarms.<sup>[19]</sup>

In the 1950s, [Otis Elevator](#) created many automatic elevators and created advertisements centered on their automation. Over time, most elevators did not need operators, and were running on their own.<sup>[19]</sup>

Resistance of Elevators caused automation to come quicker

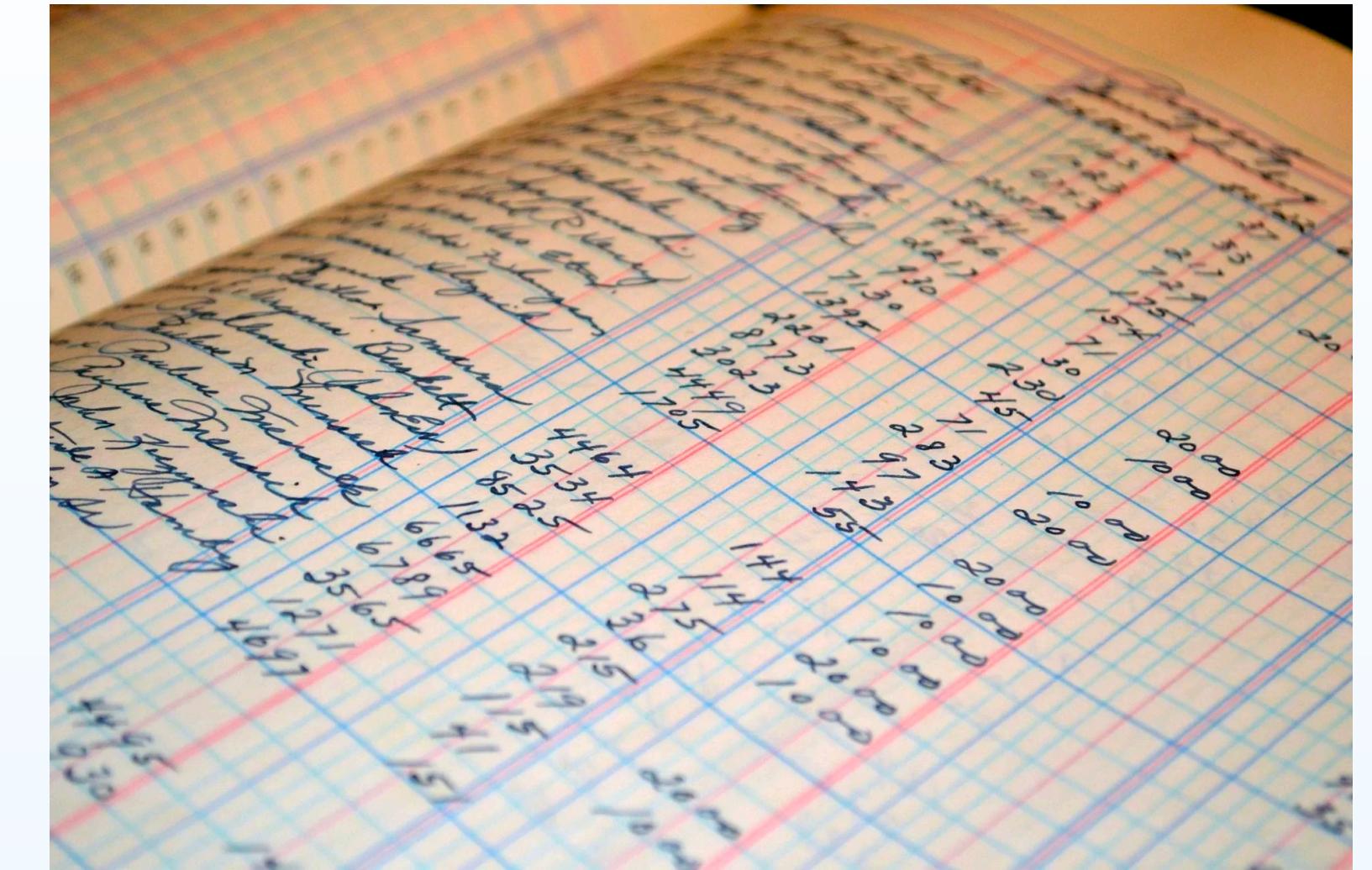


# Existential Conversation

Will A.I. replace jobs?

Case Study #2 : Electronic Spreadsheets VisiCalc

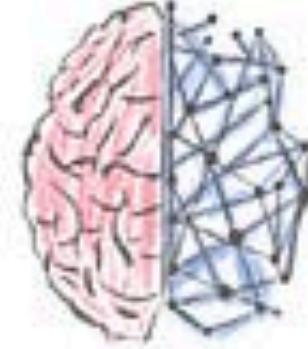
There was ~ 1 million accountants (about the same number of people who do a form of driving for a living). When the advent of electronic spreadsheet arose



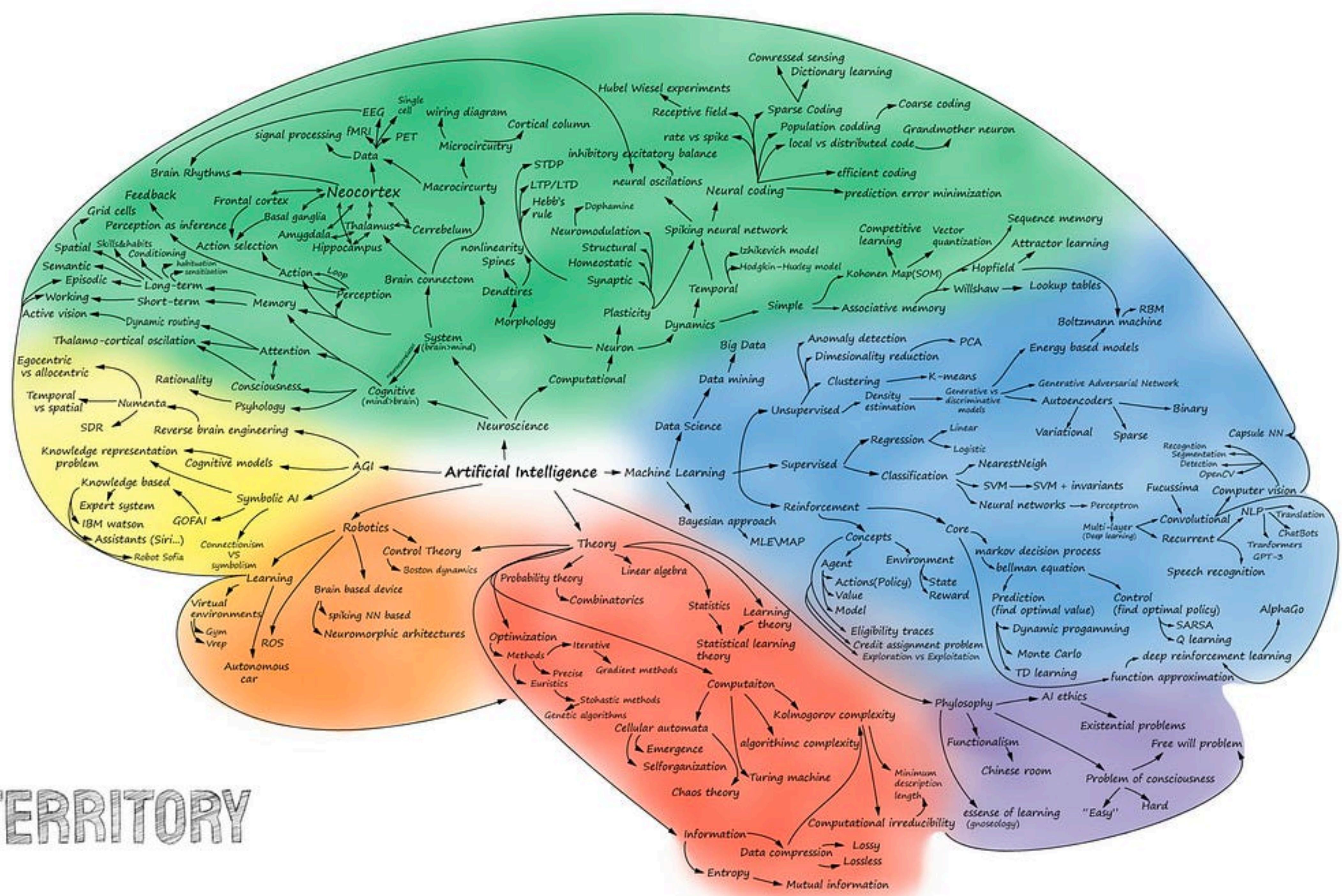
C11	(L)	TOTAL	25
1	A	B	C
2	ITEM	NO.	UNIT
3	MUCK RAKE	43	12.95
4	BUZZ CUT	15	101.25
5	TOE TONER	250	4.95
6	EYE SNUFF	2	9.90
7			
8			SUBTOTAL 13155.50
9			9.75% TAX 1282.66
10			<b>TOTAL 14438.16</b>
11			
12			
13			
14			
15			
16			
17			
18			
19			
20			

# Existential Conversation

"A.I. will not replace your job but someone using A.I. will"



# AI TERRITORY



# Many companies follow the hype but do not know anything about the topic

TECH / ARTIFICIAL INTELLIGENCE

## Forty percent of ‘AI startups’ in Europe don’t actually use AI, claims report /

Companies want to take advantage of the AI hype

According to the survey from London venture capital firm MMC, 40 percent of European startups that are classified as AI companies don't actually use artificial intelligence in a way that is “material” to their businesses. MMC studied some 2,830 AI startups in 13 EU countries to come to its conclusion, reviewing the “activities, focus, and funding” of each firm.



**“Companies that people assume and think are AI companies are probably not.”**

people assume and think are AI companies are probably not.”

“In 40 percent of cases we could find no mention of evidence of AI,” MMC head of research David Kelnar, who compiled the report, told Forbes. Kelnar says that this means “companies that

# Custom LLMs

# There are many types of LLMs

## 1. General-purpose LLMs

These models are trained on large amounts of general web text and aim to be useful for a wide range of tasks.

Examples include:

- GPT-3 — Developed by OpenAI, GPT-3 has 175 billion parameters and can generate text, answer questions and perform other tasks.
- BERT — Developed by Google, BERT uses Transformer architecture and is widely used for downstream NLP tasks.

## 2. Domain-specific LLMs

These models are trained on text from a specific domain like science, medicine or law. They tend to perform better for domain-specific tasks. Examples are:

- BioBERT — Trained on biomedical literature and used for biomedical NLP tasks.
- SciBERT — Trained on scientific text and effective for scientific information extraction and question answering.

## 3. Multilingual LLMs

These models support multiple languages and are trained on text data from different languages. Examples are:

- XLM — Developed by Facebook, XLM supports 100 languages and aims to be useful for cross-lingual tasks.
- Multilingual BERT — Supports 104 languages and can perform tasks like sentiment analysis and named entity recognition in any of those languages.

## 4. Few-shot LLMs

These models are designed to perform well even when fine-tuned with small amounts of labeled data. Examples are:

- GPT-3
- T5 — Developed by Google, T5 can achieve high performance with as little as 10 examples for training.

## 5. Task-specific LLMs

These models are tailored for a specific NLP task like summarization, question answering, translation, etc. Examples are:

- BART — Facebook's model for text generation tasks like summarization and question generation.
- ALBERT — Google's model for question answering and sentence classification.

# Open Source & Closed Source LLMs

Many LLMs are open source and

## Open Source LLMs

Open Source LLMs are language models whose source code is publicly available and can be freely accessed, used, modified, and distributed by anyone. Open source models encourage collaboration, transparency, and community involvement. Developers, researchers, and enthusiasts can contribute to the development, improvement, and customization of these models. The open-source nature allows for greater innovation, knowledge sharing, and collective development efforts.

## Closed Source LLMs

Closed Source LLMs, on the other hand, are language models whose source code is not publicly available. They are developed and maintained by organizations or companies that typically keep the underlying code proprietary and closed to the public. These models are often developed as commercial products and may require licenses or subscriptions for their use. The specific details of their architecture, training data, and algorithms are generally not accessible to the public.

# 1000's of open source LLMs are available for download and training

Many LLMs are open source and downloadable locally or accessible over the internet

The screenshot shows the Hugging Face platform interface. At the top right is a logo featuring a yellow emoji of a smiling face with hands raised, followed by the text "Hugging Face". The main search bar has the placeholder "Filter Tasks by name". Below the search bar are several navigation tabs: "Tasks", "Libraries", "Datasets", "Languages", "Licenses", and "Other". A "NEW" badge is visible above a button labeled "Create Assistants in HuggingChat". To the right of the search bar, there is a sidebar titled "Models 469,541" which lists various model cards. The cards are organized into categories: Multimodal, Computer Vision, Natural Language Processing, Audio, Tabular, and Reinforcement Learning. Each card includes the model name, a brief description, and a "View Model" button.

**Create Assistants in HuggingChat**

**The AI community building the future.**

The platform where the machine learning community collaborates on models, datasets, and applications.

Tasks   Libraries   Datasets   Languages   Licenses   Other

Filter Tasks by name

Multimodal

- 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
- Sentence Similarity

Audio

- Text-to-Speech
- Automatic Speech Recognition
- Audio-to-Audio
- Audio Classification
- Voice Activity Detection

Tabular

- Tabular Classification
- Tabular Regression

Reinforcement Learning

- Reinforcement Learning
- Robotics

Models 469,541

meta-llama/Llama-2-7b

stabilityai/stable-diffusion-v1-4

openchat/openchat

llyasviel/ControlNet-v1.1

cerspense/zeroscope\_v2\_X

meta-llama/Llama-2-13b

tiiuae/falcon-40b-instruct

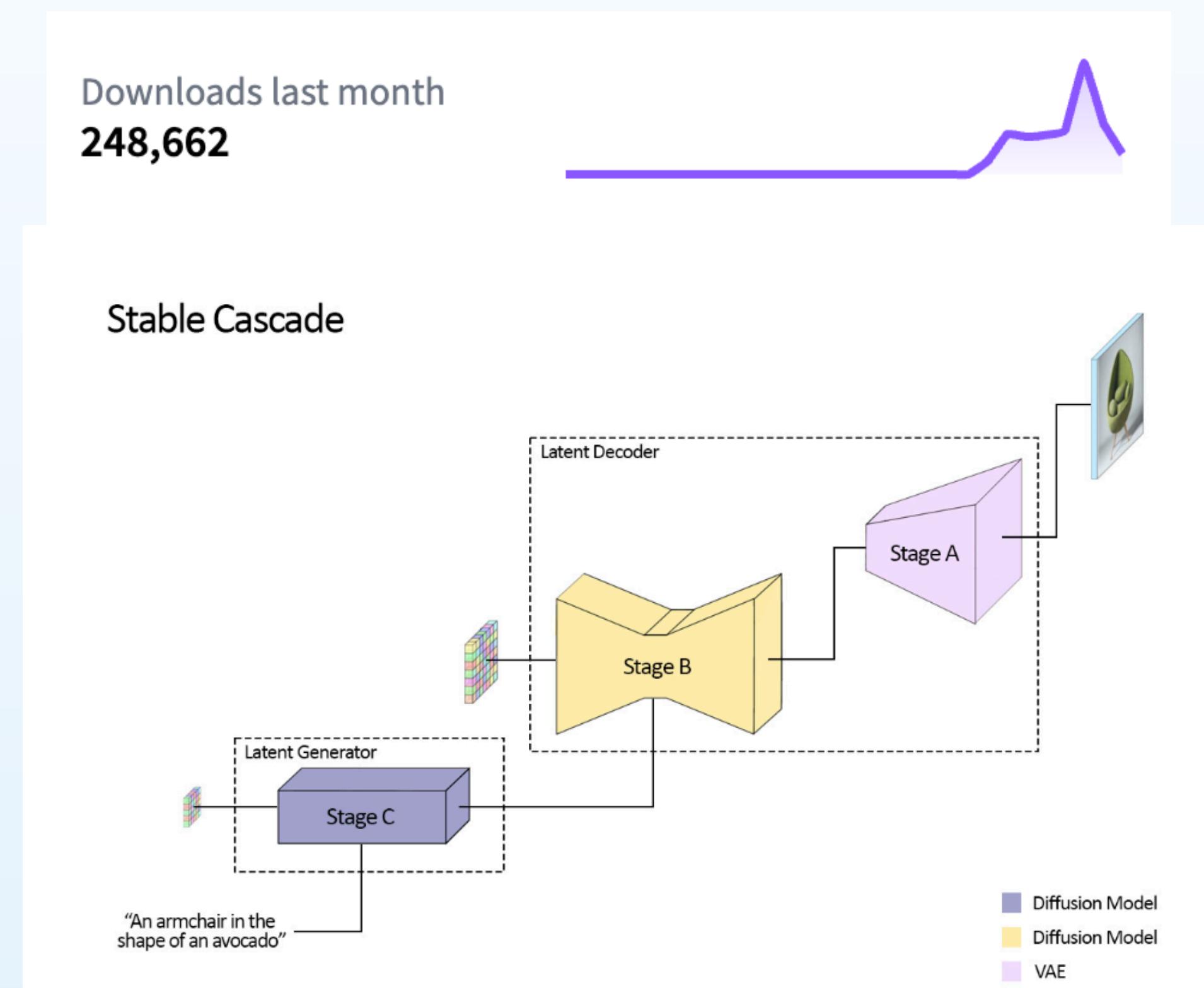
WizardLM/WizardCoder-15B-V1

CompVis/stable-diffusion-v1-4

stabilityai/stable-diffusion-2

Salesforce/xgen-7b-8k-inst

# These are downloadable to your local machine and you can run and deploy them



Models 516,341

Filter by name

new Full-text search

↑ Sort: Trending

s. stabilityai/stable-cascade

Text-to-Image • Updated 2 days ago • ↓ 249k • ❤ 749

g google/gemma-7b

Text Generation • Updated about 2 hours ago • ↓ 188 • ❤ 330

CohereForAI/ay-a-101

Text2Text Generation • Updated 1 day ago • ↓ 7.76k • ❤ 383

g google/gemma-7b-it

Text Generation • Updated about 2 hours ago • ↓ 234 • ❤ 182

BioMistral/BioMistral-7B

Text Generation • Updated 2 days ago • ↓ 891 • ❤ 179

ByteDance/SDXL-Lightning

Text-to-Image • Updated 11 minutes ago • ❤ 165

briaai/RMBG-1.4

Image-to-Image • Updated 9 days ago • ↓ 108 • ❤ 566

LargeWorldModel/LWM-Text-Chat-1M

Text Generation • Updated 10 days ago • ↓ 443 • ❤ 128

metavoiceio/metavoice-1B-v0.1

Text-to-Speech • Updated 11 days ago • ↓ 2.99k • ❤ 577

mistralai/Mixtral-8x7B-Instruct-v0.1

Text Generation • Updated Dec 15, 2023 • ↓ 1.06M • ❤ 2.95k

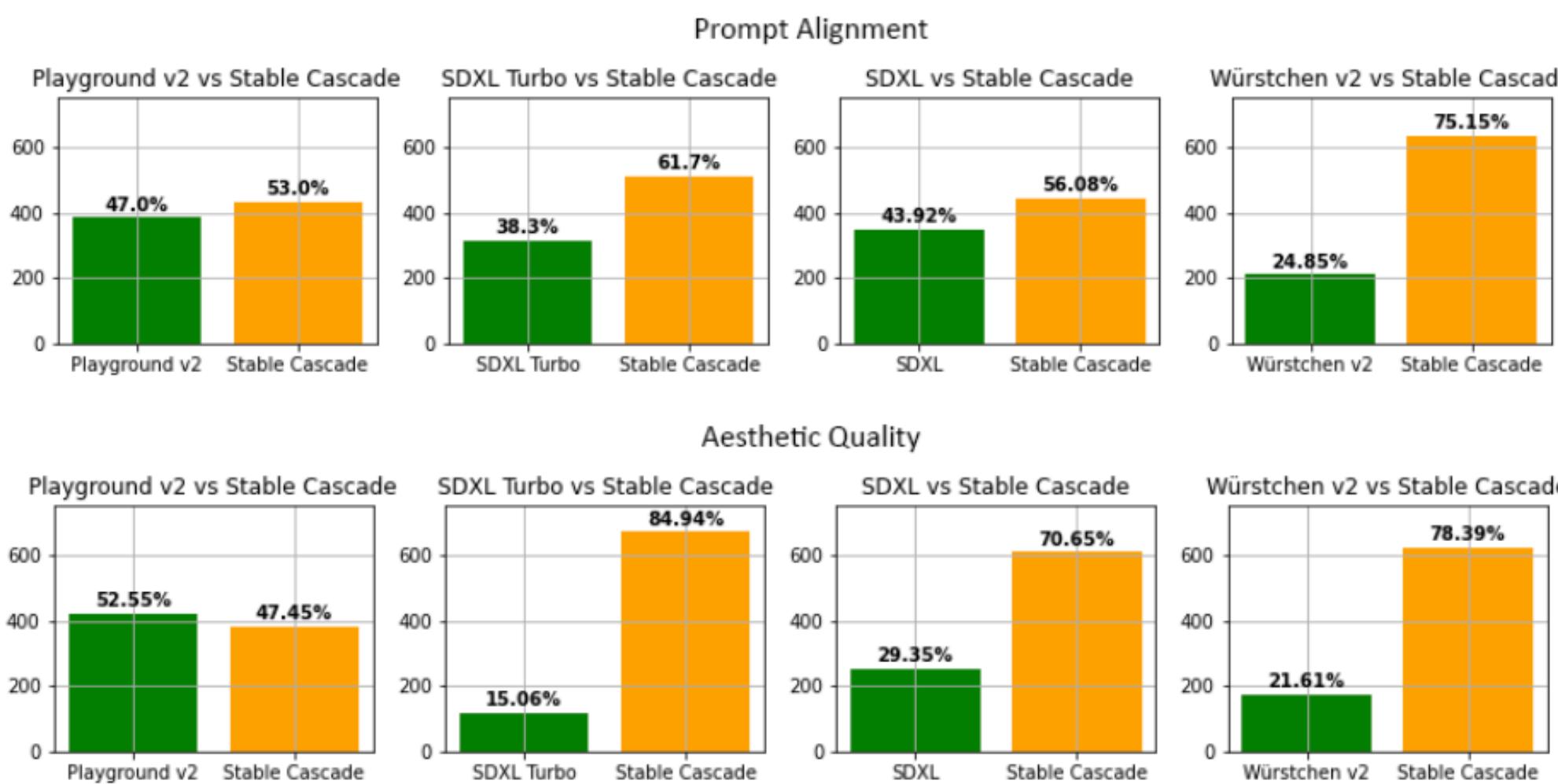
# Each of these were built for a given purpose

## Code Example

⚠ Important: For the code below to work, you have to install diffusers from this branch while the PR is WIP.

```
pip install git+https://github.com/kashif/diffusers.g...
```

### Evaluation



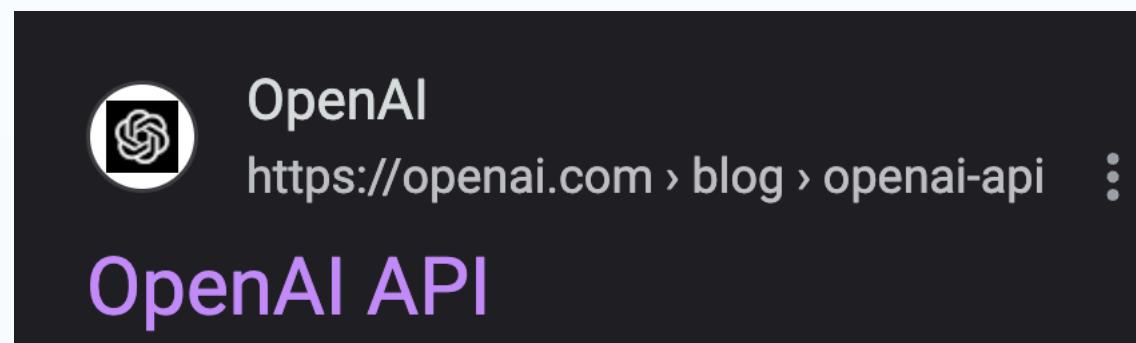
```
import torch
from diffusers import StableCascadeDecoderPipeline, S...
device = "cuda"
num_images_per_prompt = 2

prior = StableCascadePriorPipeline.from_pretrained("s...
decoder = StableCascadeDecoderPipeline.from_pretrained("...

prompt = "Anthropomorphic cat dressed as a pilot"
negative_prompt = ""
```

# OpenAI API

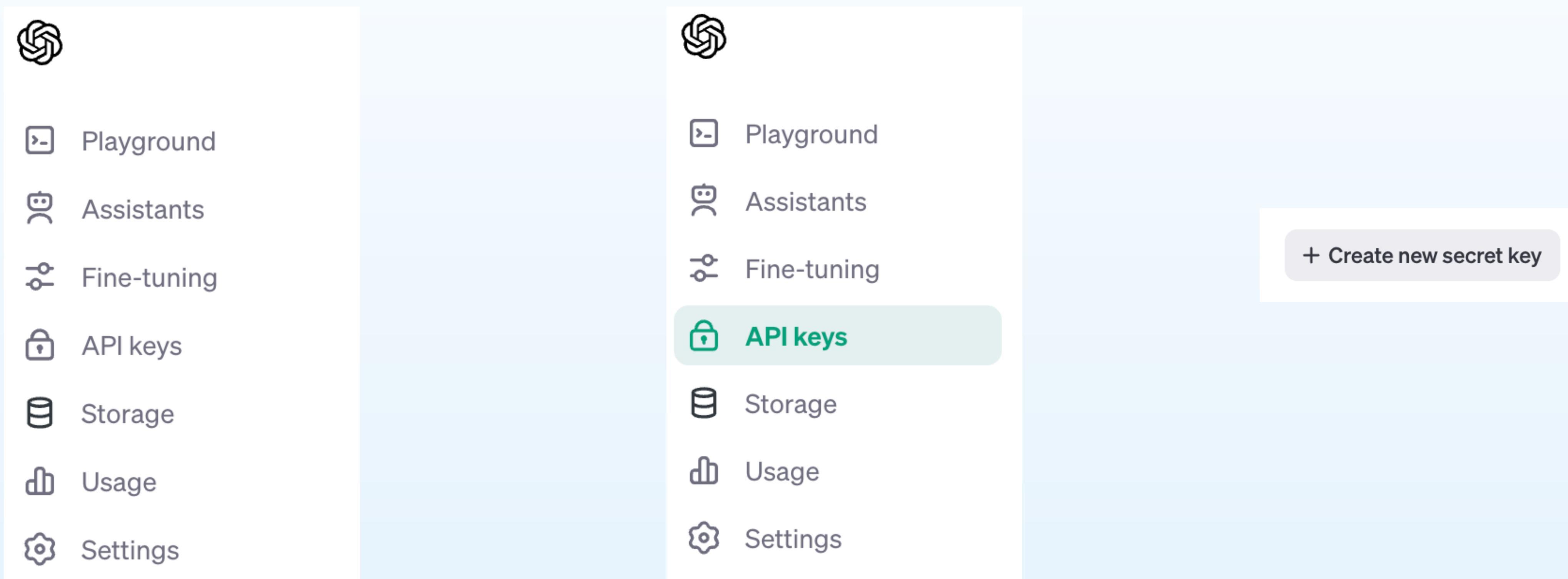
# OpenAI API



A screenshot of an orange-themed web page. At the top left is the OpenAI logo and the word "OpenAI". Below the logo is a small "Blog" link. The main title "OpenAI API" is centered in a large, bold black font. Below the title is a paragraph of text: "We're releasing an API for accessing new AI models developed by OpenAI." At the bottom are two buttons: a white "Sign up" button with a blue border and a blue "Explore the API" button with white text.

# OpenAI API

## Get API Keys



# OpenAI API

Import and install packages

```
pip install openai
```

```
import openai  
import os
```

```
openai.api_key = 'sk-m6xgK4'
```

# OpenAI API

Create a client and get a response

```
client = openai.OpenAI(  
    api_key=openai.api_key,  
)  
  
chat_completion = client.chat.completions.create(  
    model="gpt-4", # Use the appropriate model for your task  
    messages=[  
        {"role": "system", "content": "I will be passing you a t  
        {"role": "user", "content": text}  
    ],  
)  
  
chat_completion.choices[0].message.content
```

# OpenAI API Docs

[https://  
platform.openai.com/  
docs/introduction](https://platform.openai.com/docs/introduction)

## GET STARTED

Overview

### Introduction

Quickstart

Models

Tutorials

Changelog

## CAPABILITIES

Text generation

Function calling

Embeddings

Fine-tuning

Image generation

Vision

Text-to-speech

Speech-to-text

Moderation

## Create image

POST <https://api.openai.com/v1/images/generations>

Creates an image given a prompt.

### Request body

**prompt** string Required

A text description of the desired image(s). The maximum length is 1000 characters for `dall-e-2` and 4000 characters for `dall-e-3`.

**model** string Optional Defaults to `dall-e-2`

The model to use for image generation.

**n** integer or null Optional Defaults to 1

The number of images to generate. Must be between 1 and 10. For `dall-e-3`, only `n=1` is supported.

**quality** string Optional Defaults to `standard`

The quality of the image that will be generated. `hd` creates images with finer details and greater consistency across the image. This param is only supported for `dall-e-3`.

**response\_format** string or null Optional Defaults to `url`

The format in which the generated images are returned. Must be one of `url` or `b64_json`. URLs are only valid for 60 minutes after the image has been generated.

# In-Class Example

# Anyone can submit information to wikipedia. Build a wikipedia fact-checker

- 1.) Create an OpenAI API account
- 2.) Use the wikipedia api to get a function that pulls in the text of a wikipedia page
- 3.) Build a chatgpt bot that will analyze the text given and try to locate any false info
- 4.) Make User input functionality for this bit to return a report of any potentially false info via wikipedia