

Before we get started:

The image illustrates the process of opening a Google Drive link in Google Colab. It is divided into three main sections:

- Mobile Interface:** A screenshot of a mobile app showing an "In-call messages" window. A message from "You" at 12:30 PM contains a Google Drive link: https://drive.google.com/file/d/1Xqac_Y8vDpAFzA_KI_SOLDDBZxs9F5YHh/view?usp=sharing. A hand icon points to the link.
- Action Button:** A dark button labeled "Open with Google Colab" is shown above the mobile interface. A hand icon points to it, and a large black arrow points from the button towards the Colab interface.
- Colab File Menu:** A screenshot of the Google Colab "File" menu is shown. The menu items are: "Locate in Drive", "Open in playground mode", "New notebook", "Open notebook" (with "Ctrl+O" shortcut), "Upload notebook", "Rename", "Move", "Move to trash", "Save a copy in Drive" (highlighted with a hand icon), "Save a copy as a GitHub Gist", and "Save a copy in GitHub".

Hand-drawn arrows and hand icons indicate the flow of the process: clicking the link in the message, opening the Colab button, and selecting "Save a copy in Drive" in the Colab menu.

om SC

From A to Z. From zero t

NLP from scratch

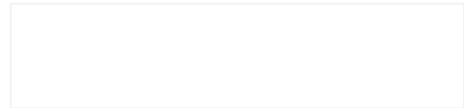


Learn natural language processing. From A to Z. From zero to hero. Fast.

What the Heck is an LLM?

Monthly Webinar

www.nlpfromscratch.com



Housekeeping



Camera on if comfortable doing so



This meeting will not be recorded



Stay muted unless speaking



Be professional

Who am I?

- Data Scientist
- Career consultant (SapientNitro, PwC, Accenture)
- Trainer
- Human





ChatGPT



Capabilities

Remembers what user said earlier in the conversation

Allows user to provide follow-up corrections

Trained to decline inappropriate requests



Limitations

May occasionally generate incorrect information

May occasionally produce harmful instructions or biased content

What is a Large Language Model?

ChatGPT is an example of a large language model (LLM), a type of *deep learning model* trained with hundreds of millions or billions of parameters on very large bodies of text. Large language models currently represent the state of the art in natural language processing (NLP) applications.

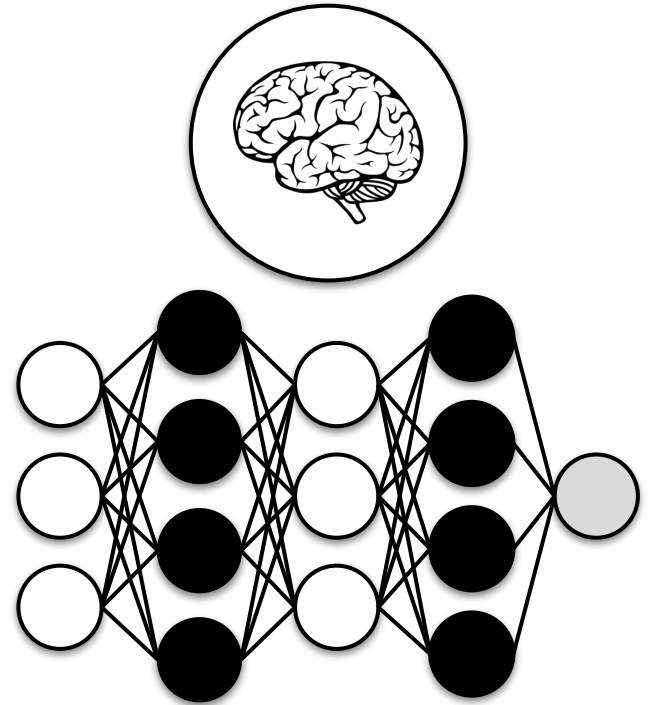
While we're here, ChatGPT is not sentient, nor is it an example of an Artificial General Intelligence (AGI).

Let's take a step back...



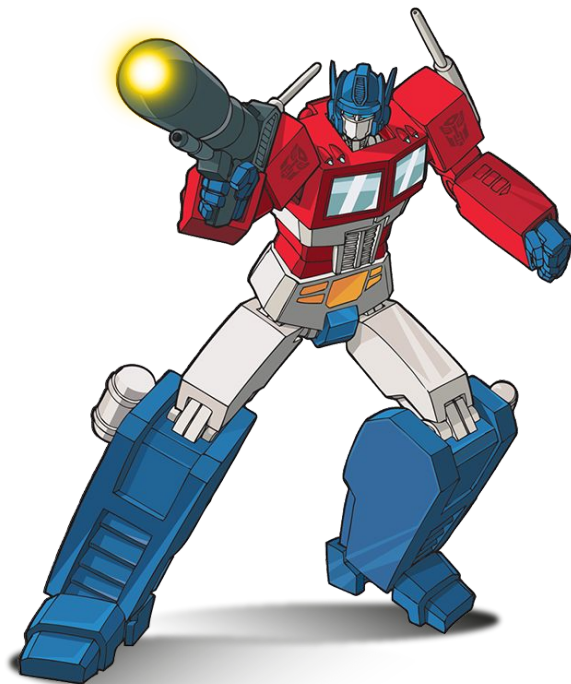
What is Deep Learning?

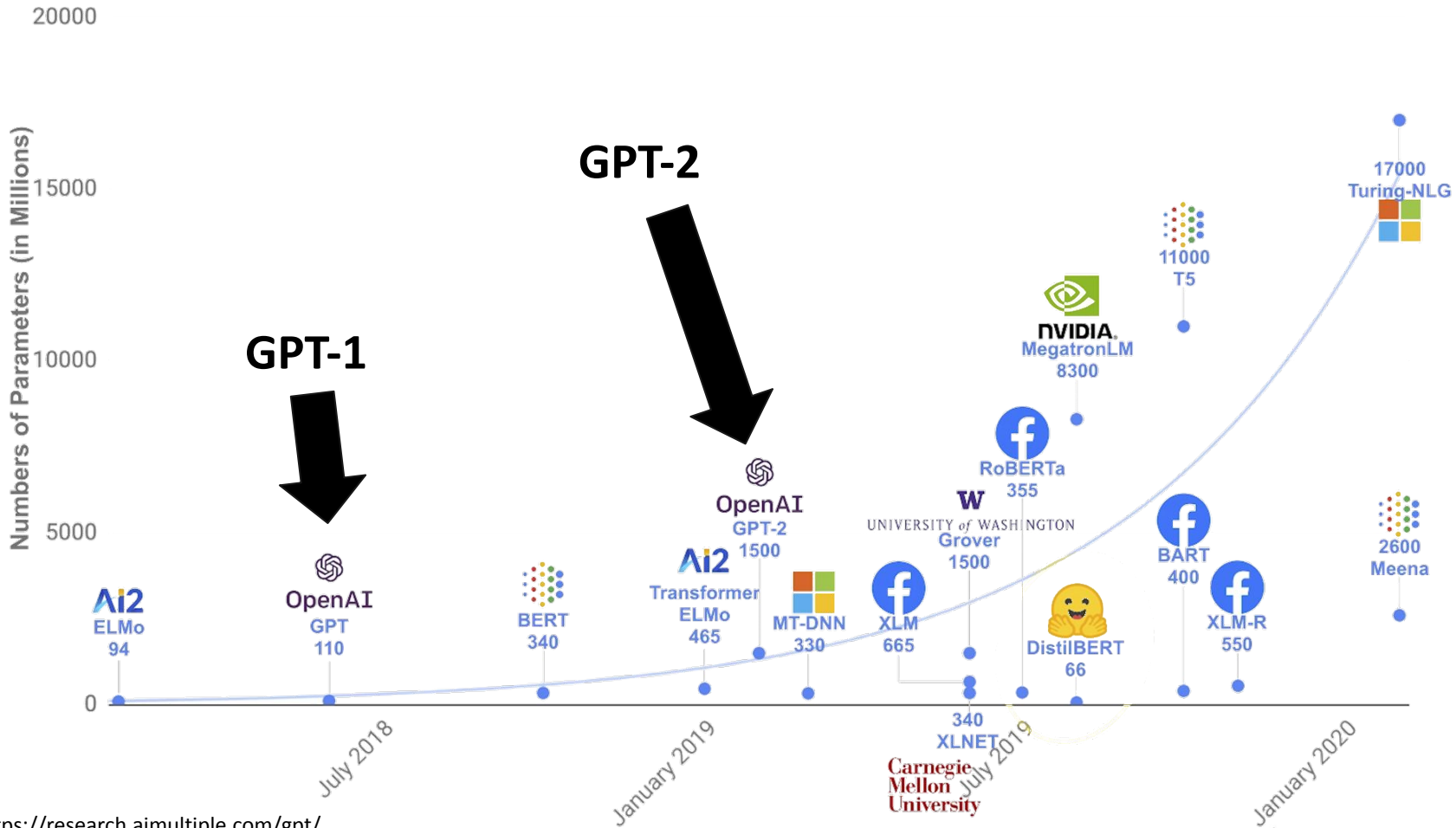
- Deep learning is the branch of machine learning dealing with neural networks - models which mimic the structure of the human brain by having individual “nodes” communicate with each other and pass information along
- Neural networks “learn” by fitting the model parameters or *weights*, by optimizing them against training data and a target objective. This could be text data where the task is to predict the most likely next letter, or classifying types of images.
- LLMs are deep learning models which are massive in size; it is understood that most have billions, hundreds of billions, or even trillions of parameters. It is also understood that these types of models are trained on very large datasets.




The Transformer Architecture

- Groundbreaking paper "Attention is All You Need" from Google researchers (Vaswani et al, 2017) introduced Transformer architecture
- Original application in machine translation but now general purpose and applied to a myriad of other tasks
- Represents the state of the art for LLMs and also applied in domains outside of language (image generation) - virtually all new models based on this architecture
- Popularized by OpenAI and the Generative Pretrained Transformer (GPT) series of models





<https://research.aimultiple.com/gpt/>

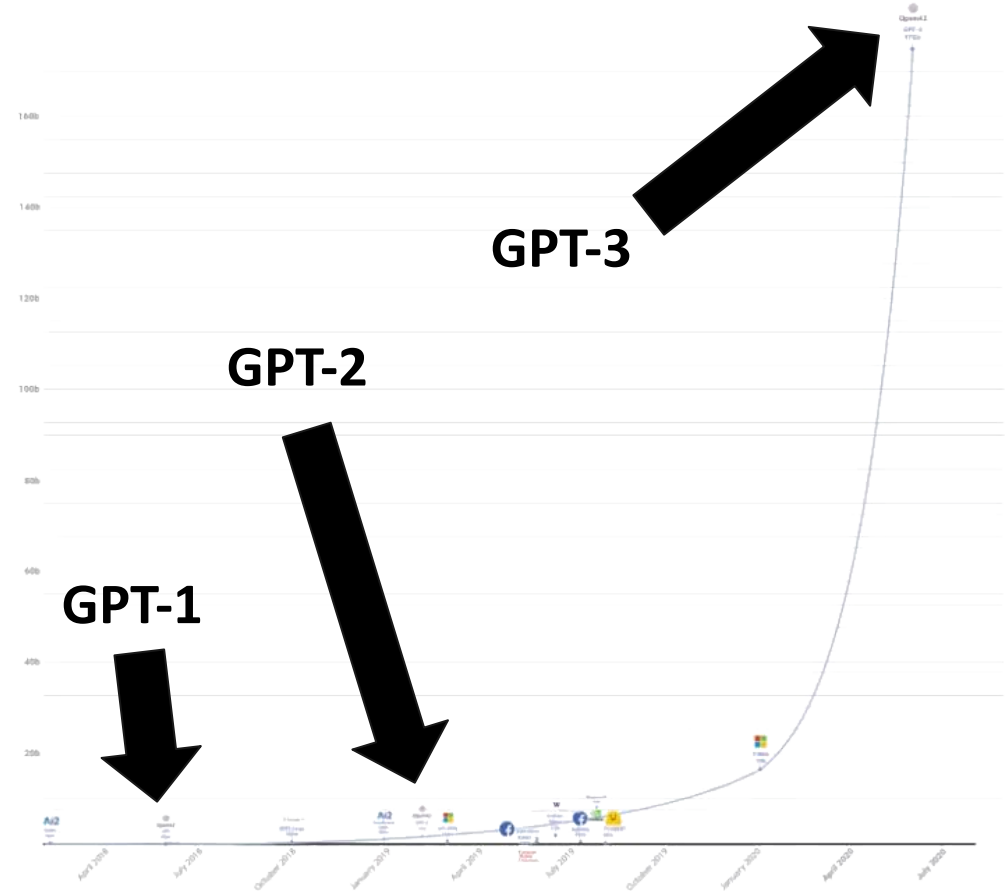
NLP from scratch 

To the Moon?


GPT-2 (2019):
1.5B parameters

GPT-3 (2020):
175B parameters

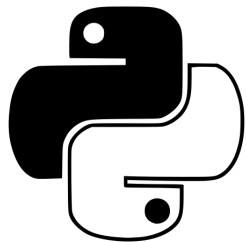
GPT-4 (March 14, 2023):
~1.8T (?) parameters



<https://research.aimultiple.com/gpt/>

NLP from scratch 

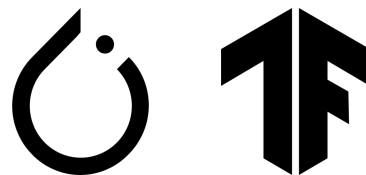
Tools of the Trade



Python 3



Google Colab
/ Jupyter



Hugging Face

Deep Learning and LLM

Deep Learning Frameworks



- Google product
- Graph-based computation, GPU training
- Other deployment options (Tensorflow Lite, TF.js)
- Easy with integration of Keras into TF 2.x



- Meta product
- Graph-based computation, GPU training
- Pytorch Mobile for embedded, no web (ONNX?)
- OOP dev focus (ML eng), Lightning equivalent to Keras

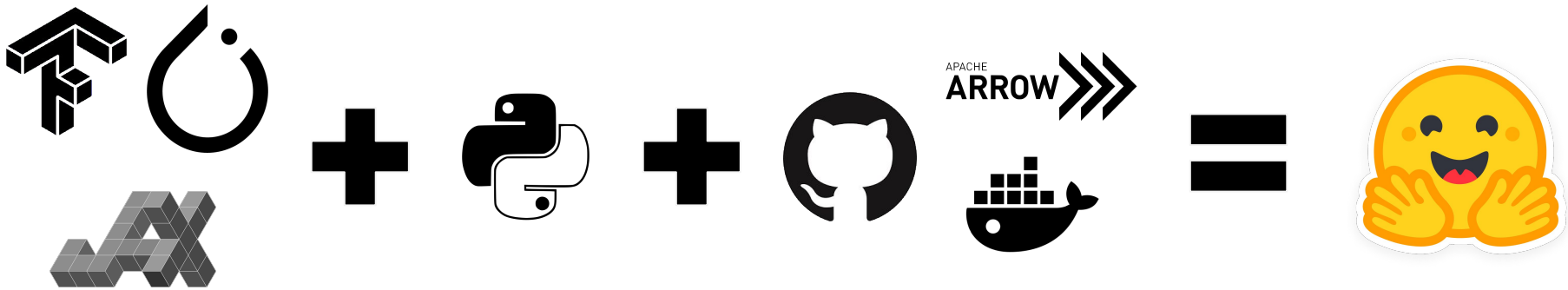
Hugging Face


Hugging Face is a software company founded in 2013 and based in New York city. As of August 2023, the company is in Series 'D' funding with a valuation of \$4.5B and backing from companies such as Salesforce, Google, Amazon, IBM, Nvidia, AMD, and Intel.

While this name refers to the company, it also refers to the software and platform they develop for working with large language models and data in the natural language processing and other domains.

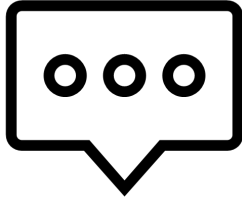
The `datasets` library allows working with data hosted on the platform, and the `transformers` library for working with models of this type. There are also other libraries for working with specialized types of models (e.g. `diffusers` for diffusion models) and data processing and model optimization.



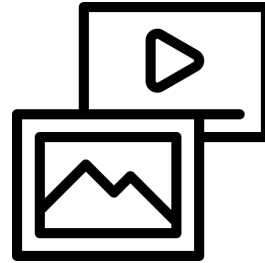


NLP from scratch 

LLM Use Cases



**Generative
Text**



**Synthetic
Media
Generation**



**Automatic
Speech
Recognition**

Generative Text: GPT

Undoubtedly, the most popularly known generative text model is that of the Generative Pretrained Transformer (GPT) by OpenAI.

As we've seen, there have been a series of GPT models of increasing size and trained on increasingly large and more complex datasets.

While GPT-3 remains proprietary and only available to use through the OpenAI API, the weights of GPT-2 are publicly available and can also be accessed through Hugging Face.

Let's take a look at generating text with GPT-2.

openai/gpt-2

Code for the paper "Language Models are Unsupervised Multitask Learners"



14
Contributors

9
Used by

20k
Stars

5k
Forks

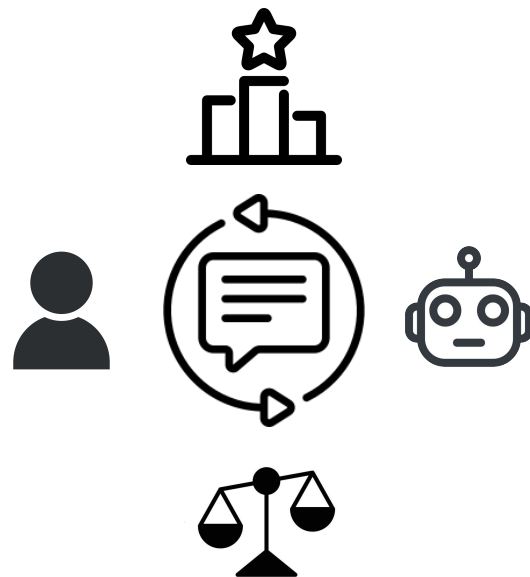


What makes ChatGPT so convincing?

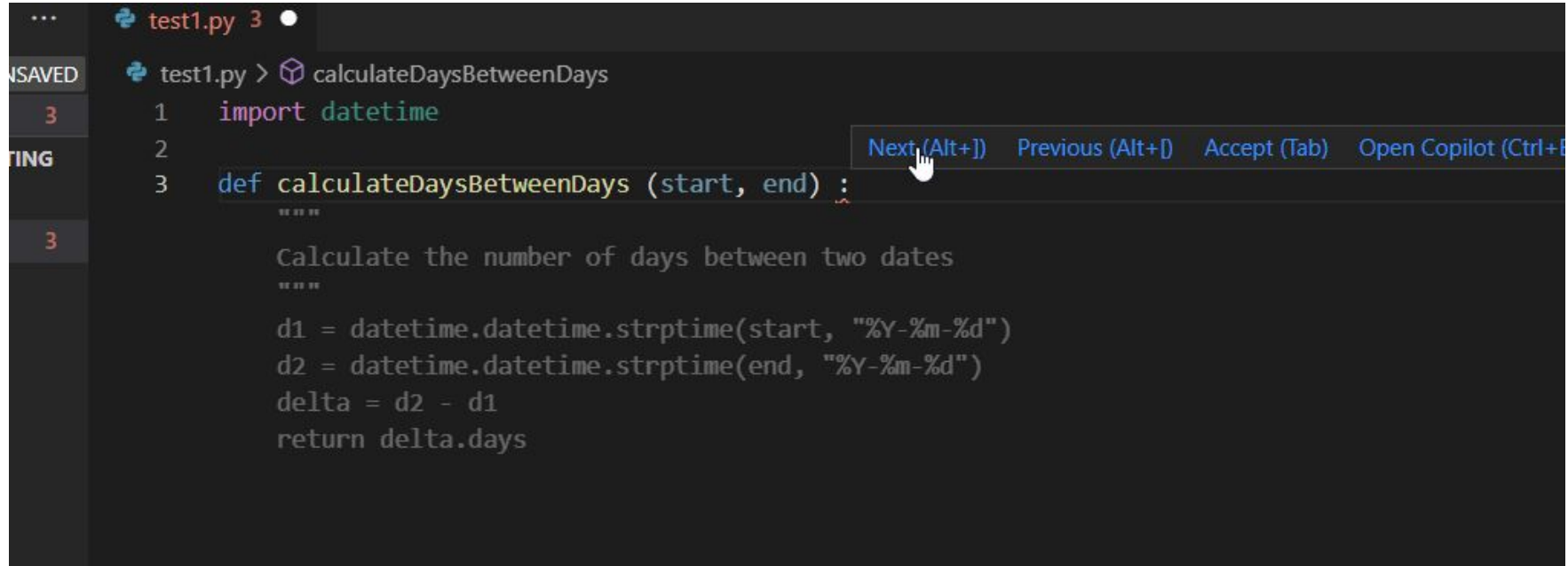
In addition to being “pre-trained” on a very large corpus of publicly available text data and having a very large model size, ChatGPT and other models are also fine-tuned and further refined using reinforcement learning from human feedback.

This method of keeping “humans in the loop” by providing feedback to the model of preferred responses, allows the models to display learn the abilities they display for things like detailed question answering and summarization.

The complexity added is that an additional reward model must be trained and incorporated into the overall model development process.



Use Case: Code Completion (Github Copilot)



The screenshot shows a code editor window with a file named 'test1.py' open. The editor is displaying a Python function definition for 'calculateDaysBetweenDays'. The function signature is 'def calculateDaysBetweenDays (start, end) :'. A tooltip is visible over the colon, showing navigation options: 'Next (Alt+J)', 'Previous (Alt+D)', 'Accept (Tab)', and 'Open Copilot (Ctrl+I)'. The function body contains a docstring and code to calculate the number of days between two dates using 'datetime.datetime.strptime' and 'delta.days'.

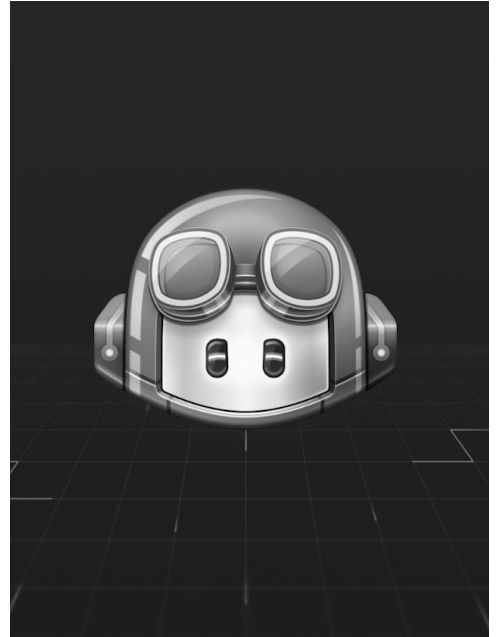
```
test1.py 3
UNSAVED test1.py > calculateDaysBetweenDays
3 1 import datetime
2
3 def calculateDaysBetweenDays (start, end) :
    """
    Calculate the number of days between two dates
    """
    d1 = datetime.datetime.strptime(start, "%Y-%m-%d")
    d2 = datetime.datetime.strptime(end, "%Y-%m-%d")
    delta = d2 - d1
    return delta.days
```

Adoption

“When we first launched GitHub Copilot for Individuals in June 2022, more than 27% of developers’ code files on average were generated by GitHub Copilot. Today, GitHub Copilot is behind an average of 46% of a developers’ code across all programming languages...”

- *GitHub Product Blog, February 2023*

<https://github.blog/2023-02-14-github-copilot-now-has-a-better-ai-model-and-new-capabilities/>



Generating Images: Stable Diffusion

One of the most popularly known models for media generation is that of Stable Diffusion, created as part of research by the CompVis group at the University of Munich and funded by Stability AI.

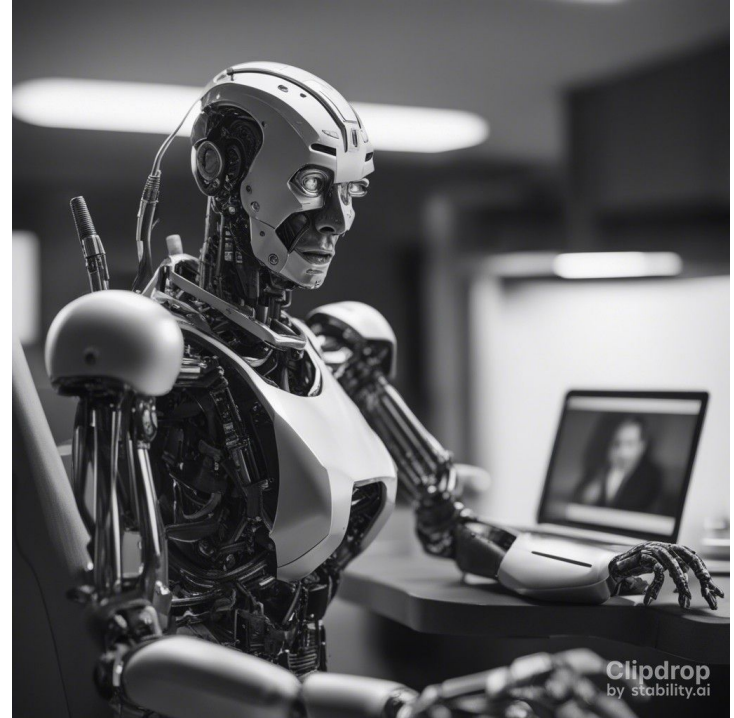
Stable Diffusion is a type of latent diffusion model, a neural network that maps images and text to a latent space and does image generation through a denoising (diffusion) process similar to that of generative adversarial networks.

This family of models gained considerable notoriety given Stability's decision to publicly release the code and weights under a license, making the model freely available to use and also making generative image capabilities widely available.



Stable Diffusion XL (SDXL)

- Released July 2023 by researchers at Stability AI, the successor to Stable Diffusion 2.1
- 3x in size to (core of) original model
- Additional refiner model (image-to-image) for denoising used in a supplementary fashion after base model for high fidelity outputs
- Available through Clipdrop (paid) and on Hugging Face spaces (free, various)
- Now near real-time image generation “as you type” with SDXL Turbo



Use Case: AI Photo Editing (Adobe Firefly)



Use Case: Model Texture Generation (Meshy)

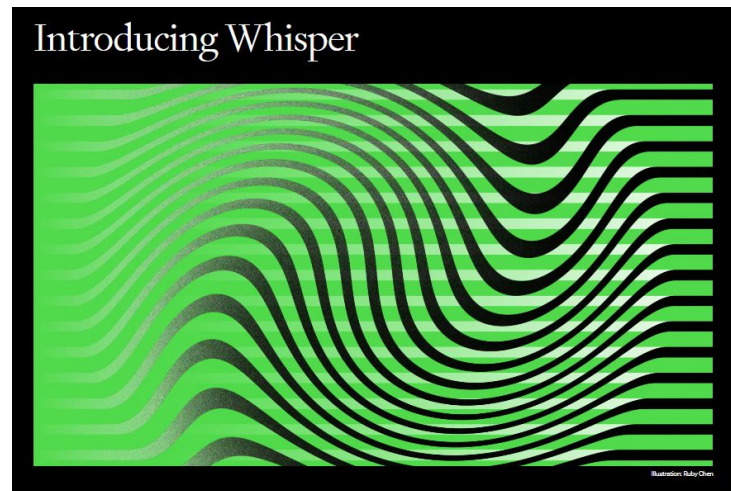


Speech Transcription: Whisper

OpenAI's Whisper is a series of multilingual models for high performance transcription of human speech (i.e. automatic speech recognition or ASR).

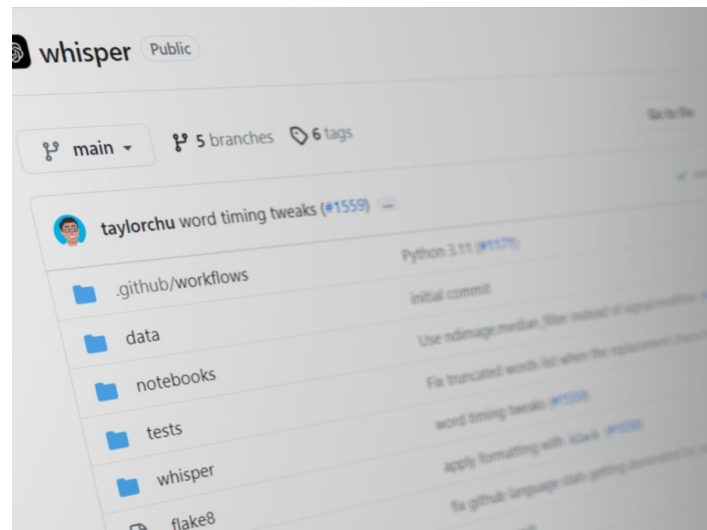
The models are a type of sequence-to-sequence transformer trained on multiple tasks (multilingual speech recognition, translation, spoken language identification, voice activity detection) but fundamentally work by first converting the audio to a spectrogram and then generating sequences of output tokens based on this input.

The company has open-sourced these models and made the weights publicly available. There is also a series of different model sizes to use, ranging from whisper-tiny (39M parameters, ~151 MB) to whisper-large (1.5B parameters, ~6.2 GB).



Building with Whisper

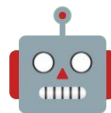
- The Whisper model is also available to be used directly through the [python package](#)
- Model weights and code are also available directly from [the github page](#)
- Available directly from [OpenAI as a service](#) (\$0.006 / minute (rounded to the nearest second))
- [Whisper V3 released](#) as part of OpenAI Dev Day (November 6th, 2023) with significant improvements across multiple languages over large V2 model
- [Distil-Whisper](#) from Hugging Face is a distilled version 6x faster, 49% smaller, and performs within 1% word error rate (WER)



Use Case: Subtitle Generation (e.g. Google Meet)



NLP4Free



<https://github.com/nlpfromscratch/nlp4free>

A Free Natural Language Processing (NLP) microcourse, from basics to deep learning

```
# Remove punctuation with regex
import re
my_review = re.sub('[^A-Za-z0-9]+', '', my_review)

# Stem
my_review = ' '.join([ps.stem(token) for token in my_review.split()])
```



LLM and Generative AI Workshops



Generative Text Models & Fine-tuning LLMs

- Intro to Hugging Face
- LLMs for generative text
- Fine-tuning models
- PEFT (LoRA) and quantization



Building GenAI Apps with OpenAI and GPT

- Intro to APIs and OpenAI
- Working with the OpenAI API
- Setting up a dev environment
- Build a streaming chat app
- Open source alternatives to GPT



Intro to Python & Natural Language Processing

- Intro to NLP
- Intro to Python
- Working with text in Python
- Manipulating text in Pandas



GenAI for Work (1 hr)

- Introduction to GenAI
- Generative AI landscape
- Everyday use cases for GenAI tools
- Prompting and prompt engineering
- Productivity with GenAI

Today ●



December 2023



Filters ▾

Sun	Mon	Tue	Wed	Thu	Fri	Sat
						2
	3 ●	4	5	6	7	8
	What the Heck is a...	Deconstruct Care...		Zero to NLP in 60...		
10	11	12	13	14	15	16
Large Language M...			What the Heck is a...		Large Language M...	
17	18	19	20	21	22	23

Sign up at nlpfromscratch.com/training!

Upcoming events
Dec 01 - Feb 2

DEC 04 What the Heck is a LLM? - F...

DEC 05 Deconstruct Care Progress

DEC 07 Zero to NLP FREE Wo

DEC 11 Large Language Models a

DEC 13 What the Heck is a LLM? - F...

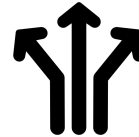
Manifesto



Knowledge is only valuable if it is useful.



The best way to learn is by doing.



Learning is a non-linear process.



Learning is exploration, not a journey.



Teaching and learning are complementary.

I would value your feedback.



NLP from scratch



www.nlpfromscratch.com