



# AI-driven tools

## Applications in research, pitfalls and perspectives

Mauricio Moreira-Soares, 31.01.2023

Oslo Centre for Biostatistics and Epidemiology, Norway

# AI-Driven tools

- Softwares that apply AI to interact with humans and execute specific tasks
  - Language models like GPT-3 and ChatGPT

# AI-Driven tools

- Softwares that apply AI to interact with humans and execute specific tasks
  - Language models like GPT-3 and ChatGPT
  - Image recognition software that can identify objects or people in photos

# AI-Driven tools

- Softwares that apply AI to interact with humans and execute specific tasks
  - Language models like GPT-3 and ChatGPT
  - Image recognition software that can identify objects or people in photos
  - Chatbots that can respond to customer inquiries and carry out simple tasks

# AI-Driven tools

- Softwares that apply AI to interact with humans and execute specific tasks
  - Language models like GPT-3 and ChatGPT
  - Image recognition software that can identify objects or people in photos
  - Chatbots that can respond to customer inquiries and carry out simple tasks
  - Recommendation systems

# AI-Driven tools

## The State-of-the-Art

- Programming and data science
  - Github copilot, Akkio
- Conversational
  - ChatGPT (OpenAI), GODEL (Microsoft), Meena (Google)
- Art
  - Midjourney, DALL-E

# AI-Driven tools

What problem do they solve?

# AI-Driven tools

## What problem do they solve?

- It is all about changing how humans communicate with machines



Imagine if we needed to punch cards nowadays...

# GPT-3

## Generative Pre-trained Transformer

- Autoregressive (AR) language model with deep learning
  - The output depends on its own previous values and on a stochastic term
  - It's all about autocompleting
- Model specs
  - 175 thousand million parameters (US billion)
  - 800Gb of storage

# ChatGPT

## Generative Pre-trained Transformer

- ChatGPT is fine tuned for general purpose conversational AI
- Model specs
  - 1.5 thousand million parameters (US billion)

# InstructGPT

## Generative Pre-trained Transformer

- InstructGPT is optimised for providing step-by-step instructions
  - Cooking, building, DIY
- Model specs
  - 1.3 thousand million parameters (US billion)

# The data used

- Publicly available web pages and proprietary data up to 2021

- Wikipedia
- Common Crawl
- Books and articles

Dataset	Tokens	Assumptions	Tokens per byte	Ratio	Size
	(billion)		(Tokens / bytes)		(GB)
Web data	410B	-	0.71	1:1.9	570
WebText2	19B	25% > WebText	0.38	1:2.6	50
Books1	12B	Gutenberg	0.57	1:1.75	21
Books2	55B	Bibliotik	0.54	1:1.84	101
Wikipedia	3B	See RoBERTa	0.26	1:3.8	11.4
Total	499B			753.4GB	

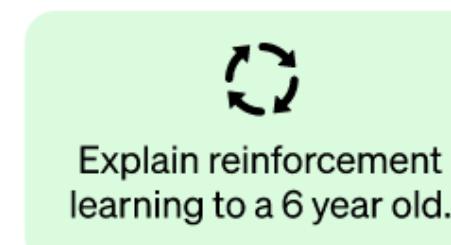
Source: <https://lifearchitect.ai/chatgpt/>

# ChatGPT

Step 1

**Collect demonstration data and train a supervised policy.**

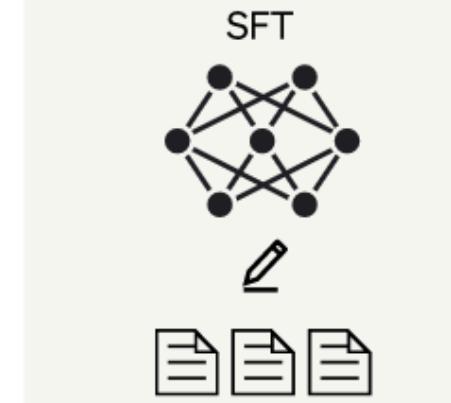
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



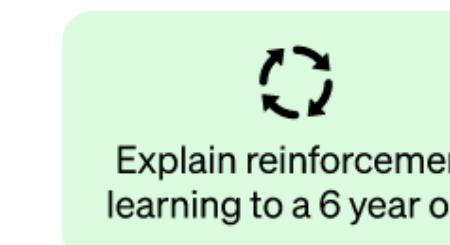
This data is used to fine-tune GPT-3.5 with supervised learning.



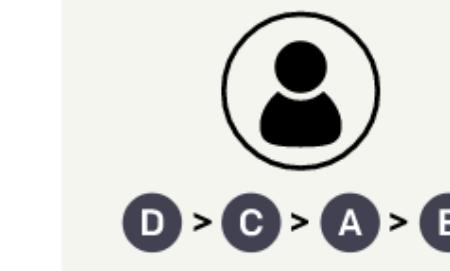
Step 2

**Collect comparison data and train a reward model.**

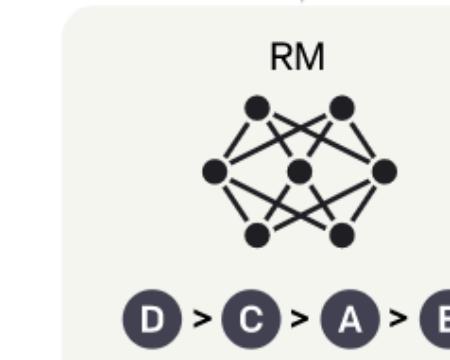
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



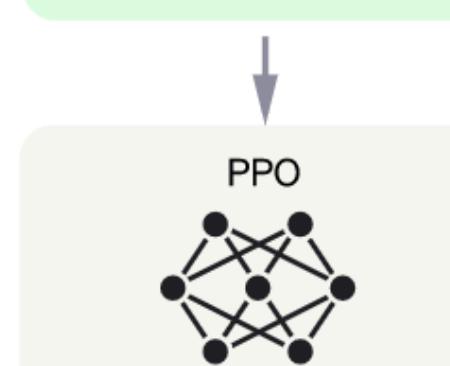
Step 3

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

A new prompt is sampled from the dataset.



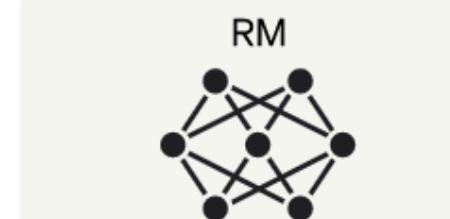
The PPO model is initialized from the supervised policy.



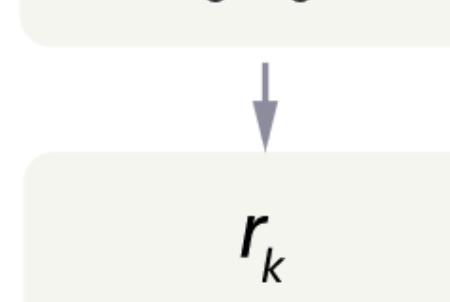
The policy generates an output.



The reward model calculates a reward for the output.

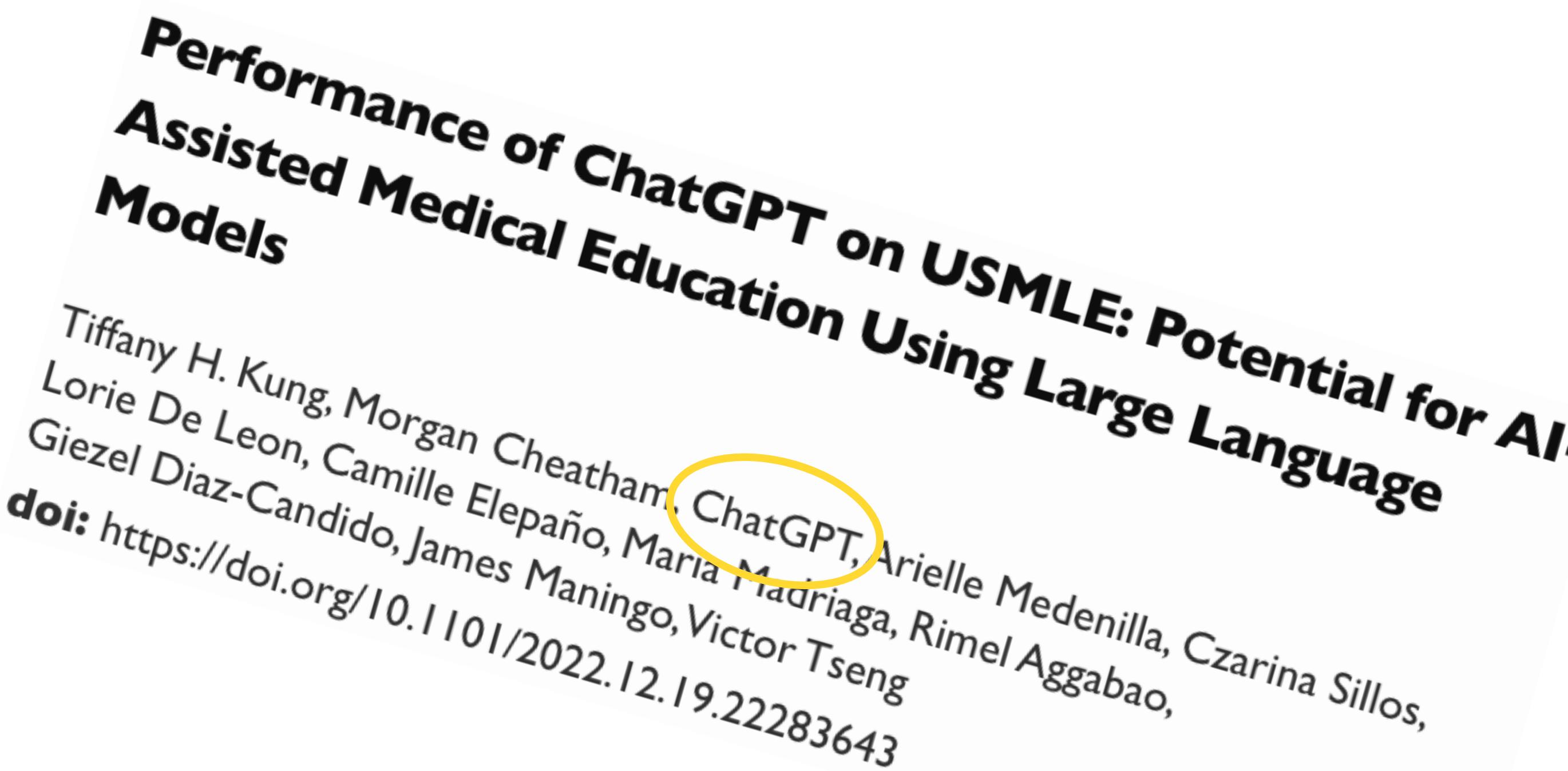


The reward is used to update the policy using PPO.



# Why should we care?

- These tools are already imposing challenges in society, including academia



# Why should we care?

- These tools are already imposing challenges in society, including academia



# Potential applications in research

- Data science
  - Data harmonisation
  - Data processing

# Potential applications in research

- Data science
  - Data harmonisation
  - Data processing
- Programming
  - Documentation generation
  - Unittest generation
  - Code translation
  - Coding itself

# Potential applications in research

- Data science
  - Data harmonisation
  - Data processing
- Writing
  - Simplifying complex concepts
  - Rephrasing
  - Kicking-off drafts
- Programming
  - Documentation generation
  - Unittest generation
  - Code translation
  - Coding itself

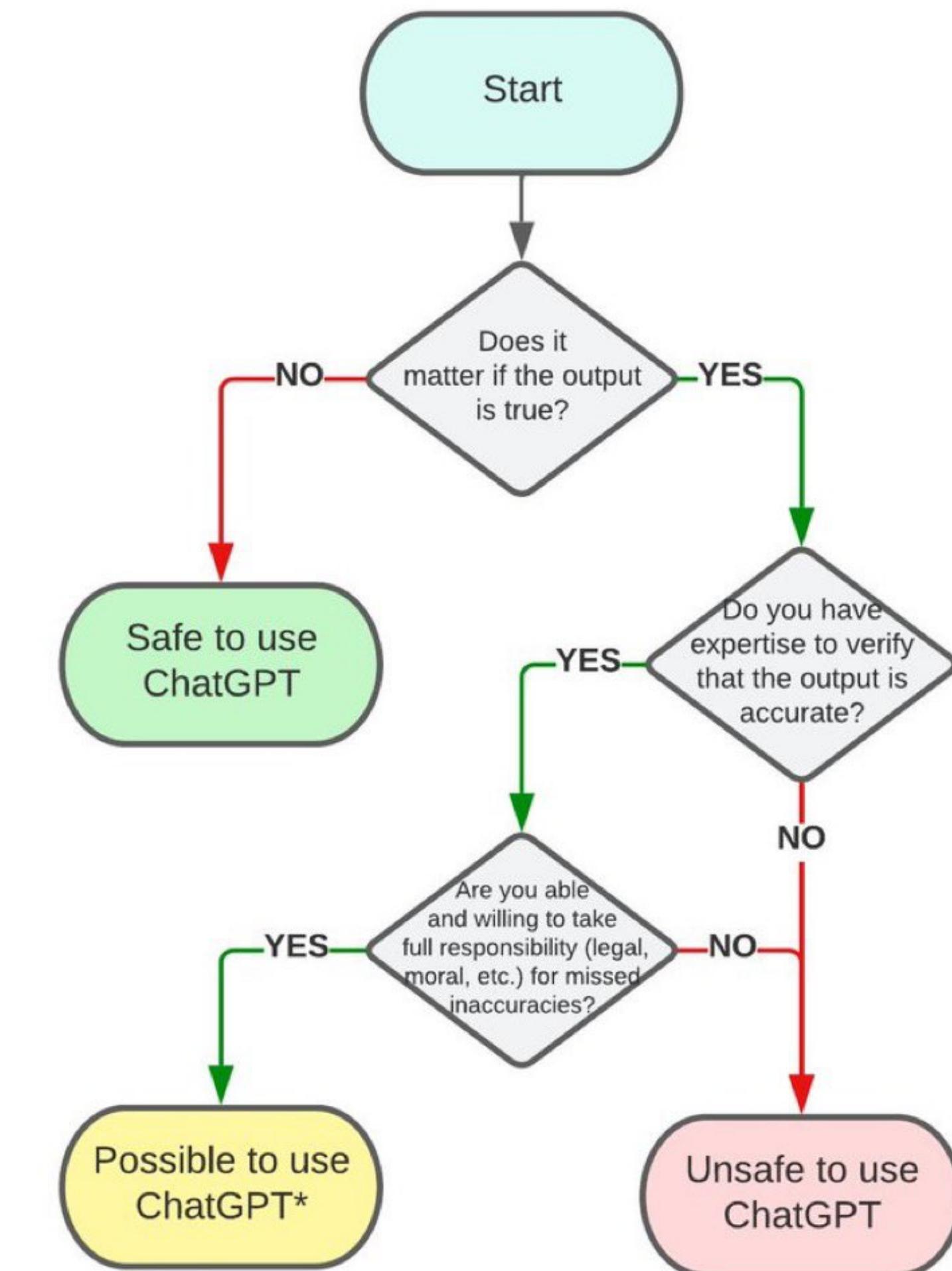
# Potential applications in research

- Data science
  - Data harmonisation
  - Data processing
- Writing
  - Simplifying complex concepts
  - Rephrasing
  - Kicking-off drafts
- Programming
  - Documentation generation
  - Unittest generation
  - Code translation
  - Coding itself



# Is it safe to use ChatGPT for your task?

Aleksandr Tiulkhanov | January 19, 2023



\* but be sure to verify each output word and sentence for accuracy and common sense



# Data harmonisation



I have two datasets that I need to harmonise. One of the datasets contains the following variables: age, sex (values: 1 - male, 2 - female), education (values: 1 - primary school, 2 - secondary school, 3 - high education). The other dataset contains the following variables: age, sex (values: 0 - male, 1 - female), education (values: 5 years, 12 years, > 12 years). Can you suggest harmonisation rules using dplyr in order to merge the two datasets?



Yes, you can use the `dplyr` library in R to harmonize the two datasets. Here are some example `dplyr` code snippets that demonstrate how you could harmonize the two datasets:

1. Harmonizing the `sex` variable:

Copy code

```
# Harmonize the sex variable in dataset1
dataset1 <- dataset1 %>%
  mutate(sex = ifelse(sex == 1, "male", "female"))

# Harmonize the sex variable in dataset2
dataset2 <- dataset2 %>%
  mutate(sex = ifelse(sex == 0, "male", "female"))
```

# Consulting documentation



Show me step by step how to build a neural network for the iris data with keras



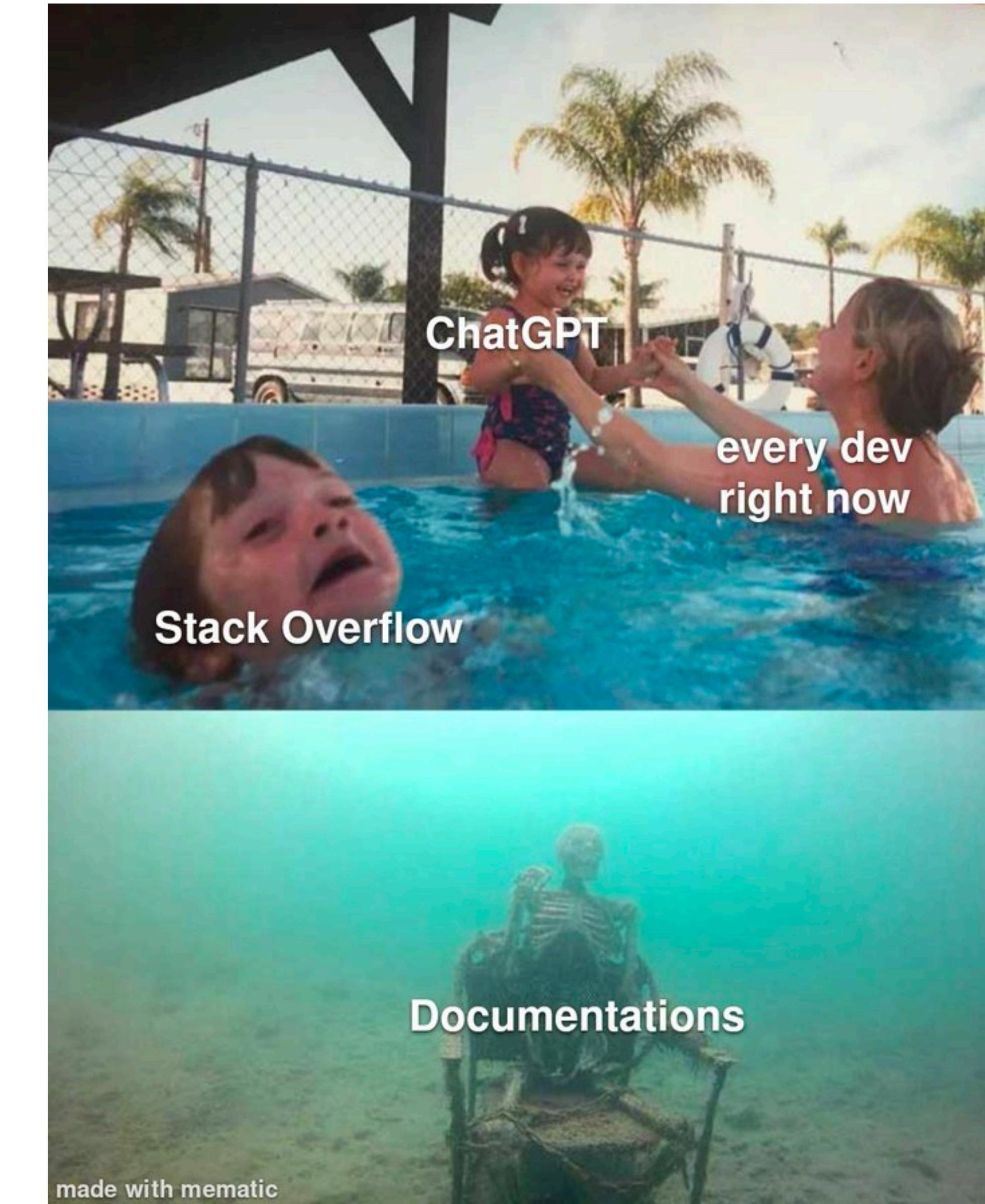
Here is a step-by-step guide to building a neural network for the Iris dataset using Keras:

1. Start by importing the necessary libraries such as numpy, keras, and sklearn. You will also need to import the Iris dataset from sklearn.

 Copy code

```
import numpy as np
from sklearn.datasets import load_iris
from sklearn.model_selection import train_test_split
from keras.models import Sequential
from keras.layers import Dense
```

2. Load the Iris dataset and split it into training and testing sets.



# Code explainability

- We can ask Github Copilot and ChatGPT for code explanations. Ex:

```
measure_rsquared_curve <- function(number_of_samples, sigma, mu, betas) {  
  rsquared <- c(1:101)  
  for (noise in 0:100){  
    # generate data  
    df <- as.data.frame(  
      MASS::mvrnorm(number_of_samples, mu = mu, Sigma = sigma )  
    )  
  
    colnames(df) <- c("x1", "x2")  
    df <- df %>%  
      mutate(LP = betas[1] + betas[2] * x1 + betas[3] * x2 +  
        rnorm(n(), mean = 0, sd = noise))  
    # measure R-squared  
    model_lm <- lm(LP ~ x1 + x2, data = df)  
    rsquared[noise + 1] <- summary(model_lm)$r.squared  
  }  
  
  return(rsquared)  
}
```

# Code explainability

## Output from Github Copilot

- Here is the explanation for the code above:
  1. The function "measure\_rsquared\_curve" takes as arguments: the number of samples, the covariance matrix, the means of the variables, and the betas of the model.
  2. Then, for each noise level (noise in 0:100), the function creates a dataframe with the variables x1 and x2, and the outcome LP that is the linear combination of the variables plus some noise (rnorm(n(), mean = 0, sd = noise)).
  3. Then, it fits a linear model to the data and measures the R-squared.
  4. Finally, it returns the vector of R-squareds for each noise level.

# Code translation

## R to python

```
def measure_rsquared_curve(number_of_samples, sigma, mu, betas):  
  
    rsquared = np.arange(101)  
    for noise in np.arange(0, 101):  
        # generate data  
        df = pd.DataFrame(  
            np.random.multivariate_normal(mu, sigma, number_of_samples)  
        )  
        df.columns = ["x1", "x2"]  
        df["LP"] = betas[0] + betas[1] * df["x1"] + betas[2] * df["x2"] + np.random.normal(0, noise)  
        # measure R-squared  
        model_lm = sm.ols(formula="LP ~ x1 + x2", data=df).fit()  
        rsquared=noise] = model_lm.rsquared  
  
    return(rsquared)
```

# Data science

Akkio Income

**Table** Upload and configure datasets

income.csv Replace 45,222 rows, 14 columns

1 Rename × 2 Remove Missing × 3 Drop columns × 4 Create Overworked ×

one-hot encode every category column with less than 6 unique values and more than 2

Cancel Apply Transform

AI Interpretation

Transform the dataset by one-hot encoding any columns that have more than two and less than six unique values. This means that each unique value in the column will be converted into its own column, with a binary value indicating whether the row has that value or not.

Clone

Previewing 100 rows Columns Added: race\_Amer-Indian-Eskimo, race\_Asian-Pac-Islander, race\_Black, race\_Other, race\_White Columns Removed: race

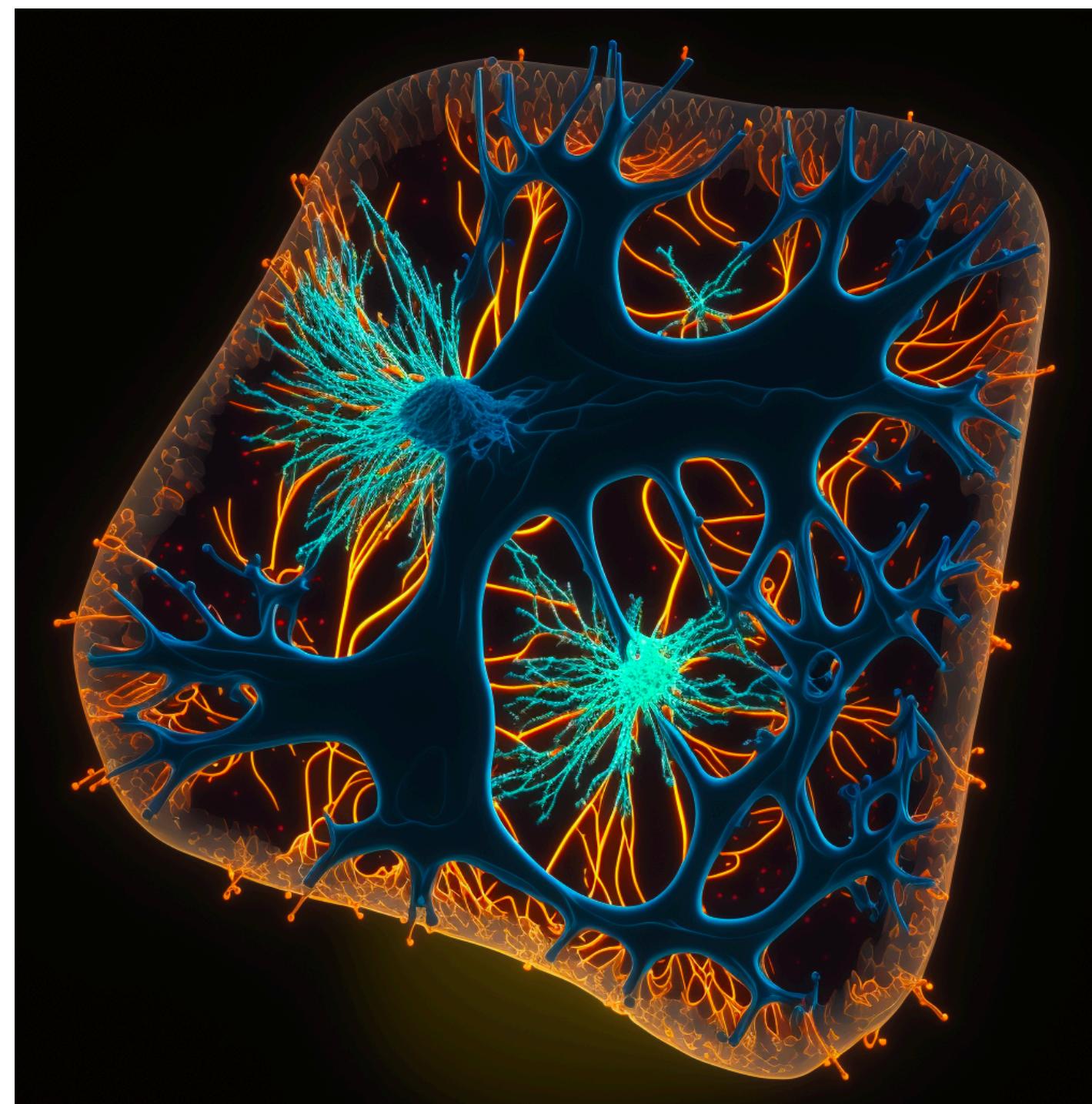
race	gender	capital-gain	capital-loss	hours-per-week	race_Amer-Indian-Eskir	race_Asian-Pac-Islan
Category	Category	Number (Integer)	Number (Integer)	Number (Integer)	Preview	Preview
COLUMN REMOVED	Male 70% Female 30%				NEW COLUMN	NEW COLUMN
1 Black	Male	0	0	40	0	0
2 White	Male	0	0	50	0	0
3 White	Male	0	0	40	0	0
4	Male	0	0	40	0	0

Example provided by Santiago: <https://twitter.com/svpino>

# Improving outreach

## Image generation with Midjourney

- **Prompt:** Adhesion modulates cell morphology and migration within dense fibrous networks



# Improving outreach

## Image generation with Midjourney

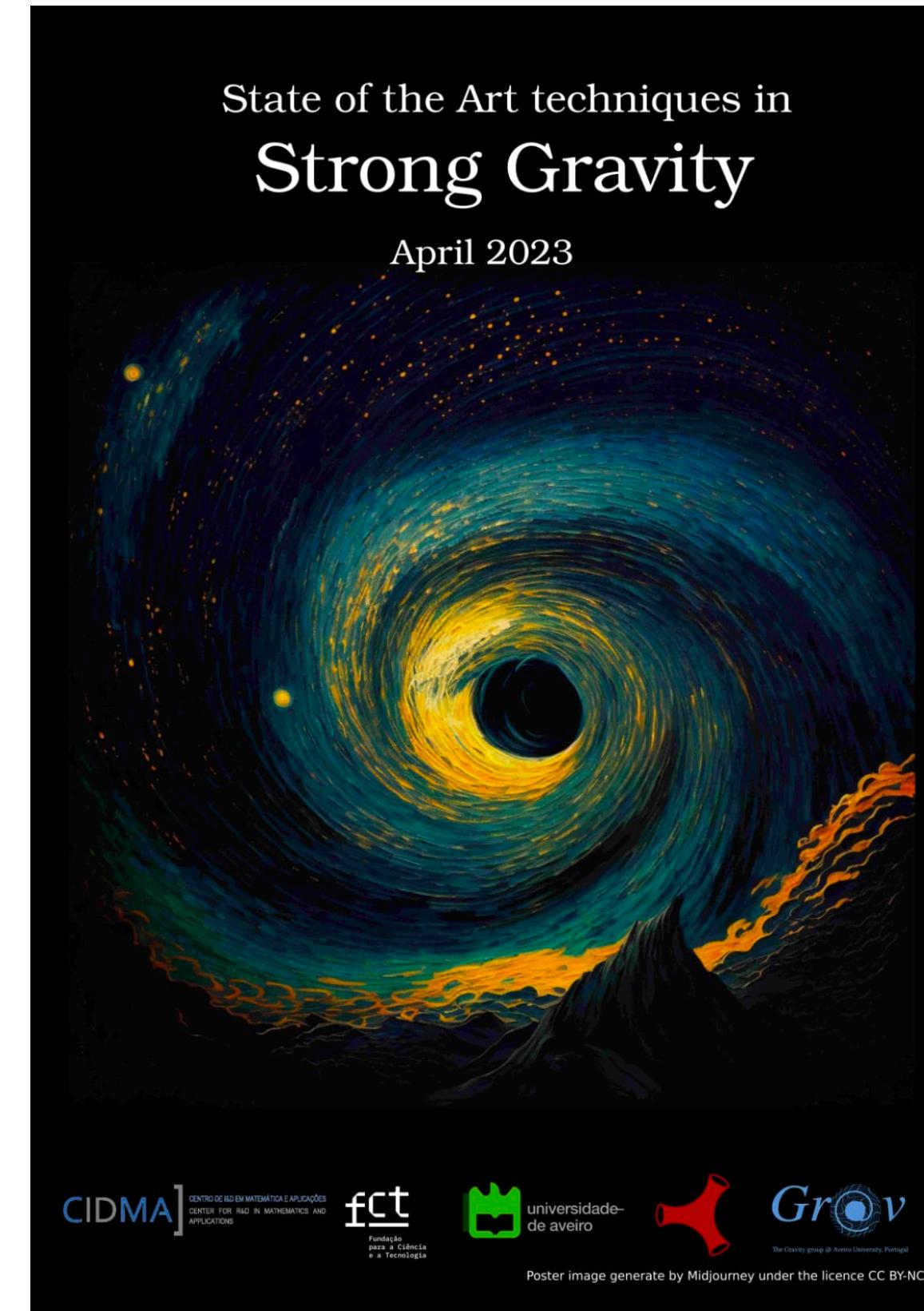
- **Prompt:** Angiogenic Factors produced by Hypoxic Cells are a leading driver of Anastomoses in Sprouting Angiogenesis—a computational study



# Improving outreach

## Image generation with Midjourney

- **Prompt:** Van gogh style painting of a black hole with a rotating disk of glowing gas



Kindly shared by Pedro V. P . Cunha

# Improving outreach

- Explain complex concepts using plain language



Explain the Damped harmonic oscillator at high education level



A damped harmonic oscillator is a system that experiences a restoring force proportional to its displacement, but also experiences a frictional force that opposes its motion. The equation of motion for a damped harmonic oscillator is given by:

$$m * x''(t) + b * x'(t) + k * x(t) = 0$$

where  $m$  is the mass of the oscillator,  $x(t)$  is the displacement as a function of time,  $x'(t)$  and  $x''(t)$  are the first and second derivatives of  $x(t)$  with respect to time,  $k$  is the spring constant, and  $b$  is the damping constant.

The solution to this equation can be represented in the form of  $x(t) = e^{rt} * (c_1 * \cos(\omega t) + c_2 * \sin(\omega t))$

# Improving outreach

- Explain complex concepts using plain language



Now explain in a short paragraph for a 6 years old



A damped harmonic oscillator is like a toy that bounces up and down, but it gets tired and slows down eventually. The toy is made up of a spring and a mass and it moves back and forth because of the force of the spring pushing it. But, sometimes there is something that makes the toy tired, like friction, and it makes the toy slow down and stop bouncing eventually. This is what we call a damped harmonic oscillator.

# Improving outreach

- Explain complex concepts using plain language



The 6 years old doesn't understand the concept of force. Can you explain it again without using this word?



Sure, A damped harmonic oscillator is like a toy that bounces up and down, but it gets tired and slows down eventually. The toy is made up of a spring and a mass and it moves back and forth because the spring is pushing it. But, sometimes there is something that makes the toy tired, like friction, and it makes the toy slow down and stop bouncing eventually. This is what we call a damped harmonic oscillator.

# Limitations

- There is no ground truth
  - The models are trained with public available information
  - They have capability to generate fluently well articulated **disinformation**
  - You need knowledge and critical thinking to judge the output
- It does not provide references or can provide made up references

# Limitations

- I asked if they used Reinforcement Learning in ChatGPT

As for the use of reinforcement learning, OpenAI has used it in some of its models such as DALL-E, but not in GPT-3 or ChatGPT. GPT-3 and ChatGPT are both pre-trained models and are fine-tuned for specific tasks, but they are not fine-tuned using reinforcement learning.

In other words, GPT-3 and ChatGPT are both fine-tuned using supervised learning techniques, but not reinforcement learning.

# Limitations

- You said that ChatGPT was not developed with reinforcement learning, but your website says the contrary. Which one is true?



I apologize for the confusion. You are correct that the OpenAI website states that GPT-3 and ChatGPT were trained using a variant of reinforcement learning called "divergence minimization."

Divergence minimization is an unsupervised learning technique that is closely related to reinforcement learning and it tries to reduce the difference between the current model's behavior and the data distribution from which it was trained on. It uses a reward function to encourage the model to generate text that is similar to the text in the training data. This can lead to the model learning to generate text that is similar to the text in the training data, and thus, it could improve its ability to interact with humans.

So, ChatGPT was not fine-tuned using reinforcement learning, but it was trained on a variant of reinforcement learning called "divergence minimization" which is closely related to it. I apologize for any confusion caused by my previous answer.

# Impact and risks for society

- LLMs can produce convincing and authoritative content
- Huge potential to be weaponised in an information war

# Impact and risks for society

- LLMs can produce convincing and authoritative content
- Huge potential to be weaponised in an information war

Press release

**UK exposes sick Russian  
troll factory plaguing social  
media with Kremlin  
propaganda**

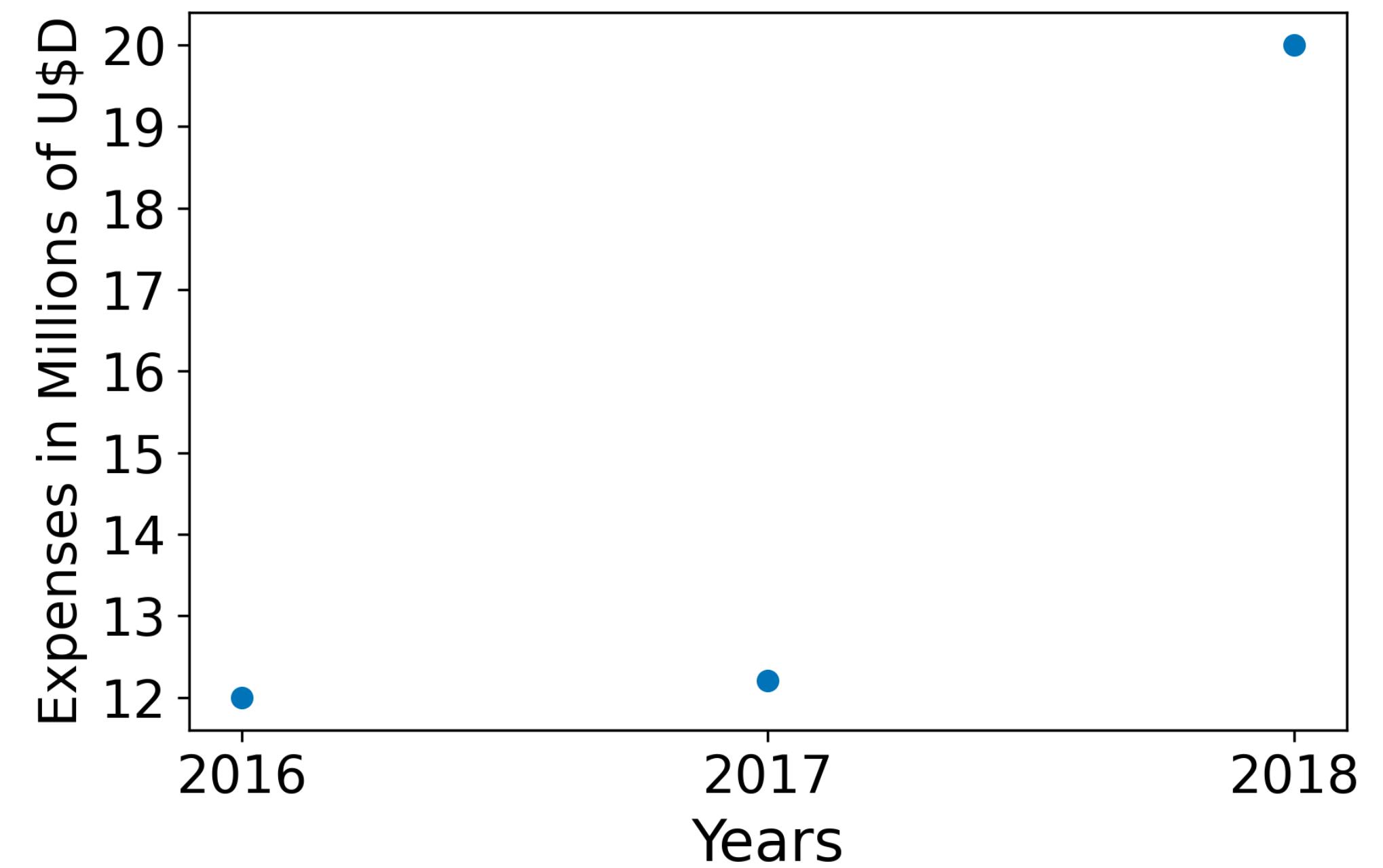
UK-funded expert research has exposed how the Kremlin is using a troll factory to spread lies on social media and in comment sections of popular websites.

# Impact and risks for society

- Information war existed before AI. Why this is important?

# Impact and risks for society

- Information war existed before AI. Why this is important?
  - It's about scale!
  - These operations have a high cost



# Understanding ChatGPT

- GPT can be implemented under 300 lines of code
  - <https://github.com/karpathy/minGPT>
- You can test other models on [huggingface.co](https://huggingface.co)



# Perspectives

- I don't believe we will be replaced by AI...

# Perspectives

- I don't believe we will be replaced by AI...
  - ... but we may be replaced by someone else using AI
- The internet changed research methodology forever

# Perspectives

- I don't believe we will be replaced by AI...
  - ... but we may be replaced by someone else using AI
- The internet changed research methodology forever
- It is likely that another change will happen soon with these tools

# Perspectives

- I don't believe we will be replaced by AI...
  - ... but we may be replaced by someone else using AI
- The internet changed research methodology forever
- It is likely that another change will happen soon with these tools
- Some valuable skills can become obsolete

# Conclusions

- AI-driven tools have the potential to make our work more efficient
- Many concerns about ethics and intellectual property are still open
- Tools like ChatGPT are not trustworthy, but there are safe use cases

# Thank you!