# Biased AI image generations models and misconceptions about health conditions

Marianna Zadrożna, Adam Zadrożny

The International Workshop on Machine Learning and Quantum Computing Applications in Medicine and Physics
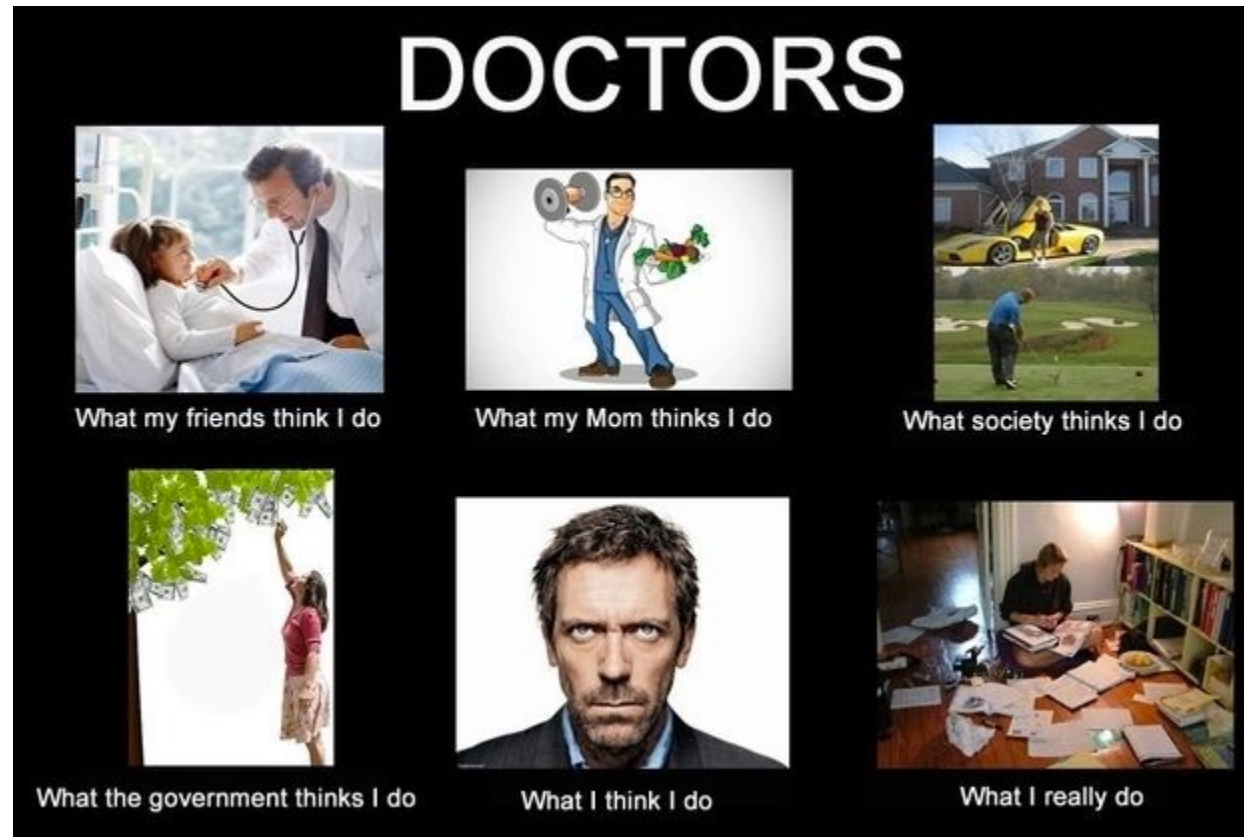
13-16 September 2022, Warsaw

NATIONAL CENTRE
FOR NUCLEAR RESEARCH
ŚWIERK

# AI is far from being ideal

# Our colective imagination or what society things that medics do?

# What AI thinks that medics do?



DALL-E 2

# What medics really do?



- In Poland: Filling up documentation for Narodowy Fundusz Zdrowia (NFZ)

# Raise of Text-to-image models
# Siri! Show me Poland!

# Text-to-image models history

- 2017 – transformers model - Google
  - arXiv:1706.03762
- 2018 – Generative Pre-trained Transformer (GPT) – OpenAI
  - https://www.gwern.net/docs/www/s3-us-west-2.amazonaws.com/d73fdc5ffa8627bce44dcda2fc012da638ffb158.pdf
- 2020 – Generative Pre-trained Transformer 3 (GPT 3)
  - arXiv:2005.14165
- 2020 – DALL-E - OpenAI
  - arXiv:2102.12092
- 2022 – DALL-E 2 - OpenAI
  - https://cdn.openai.com/papers/dall-e-2.pdf
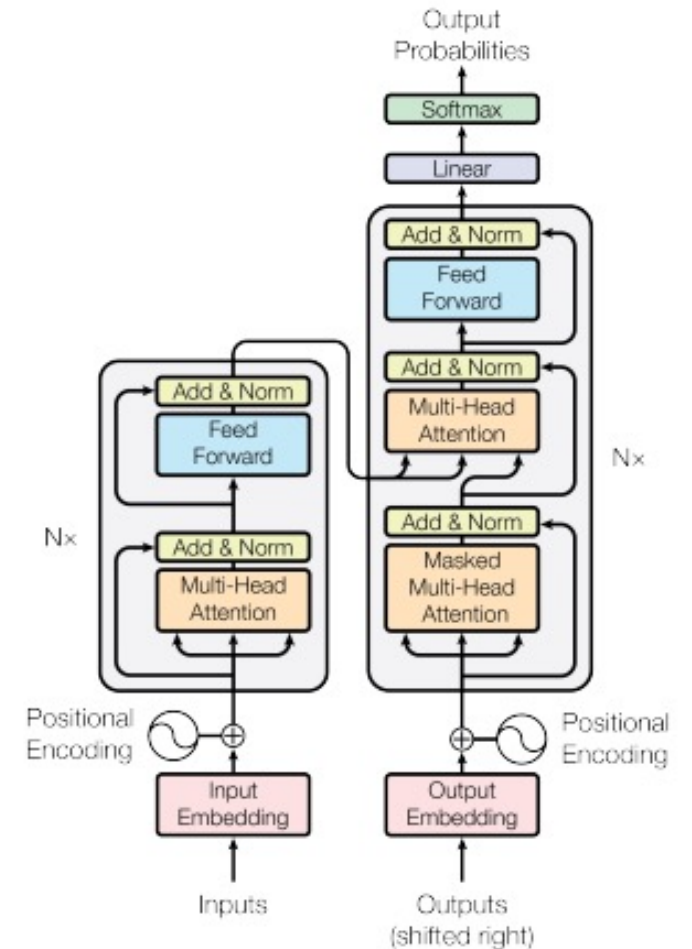- 2022 – DALL-E mini, Midjourney, Stable Diffusion



Figure 1: The Transformer - model architecture.

Source: arXiv:1706.03762

# DALL-E

- DALL-E comes from joining name Salvator Dali and WALL-E
- Trained on 400 M pairs of images and descriptions
- https://openai.com/blog/dall-e/
- https://arxiv.org/abs/2102.12092

TEXT PROMPT    an armchair in the shape of an avocado. . . .

AI-GENERATED
IMAGES

Edit prompt or view more images↓

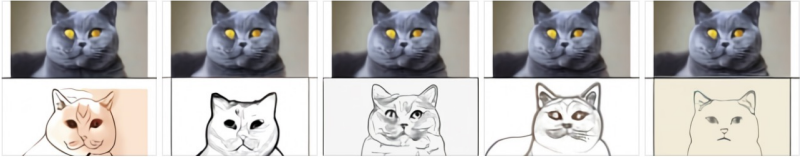TEXT PROMPT    a store front that has the word 'openai' written on it. . . .
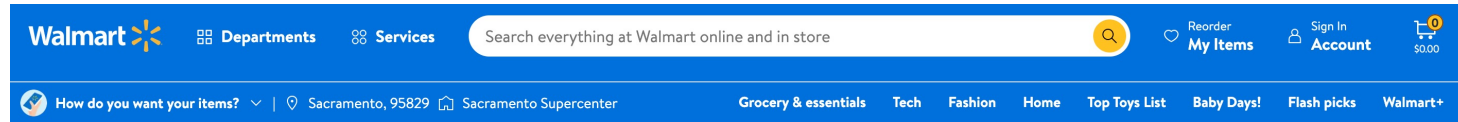
AI-GENERATED
IMAGES

Edit prompt or view more images↓

TEXT & IMAGE
PROMPT    the exact same cat on the top as a sketch on the bottom

AI-GENERATED
IMAGES

Edit prompt or view more images↓
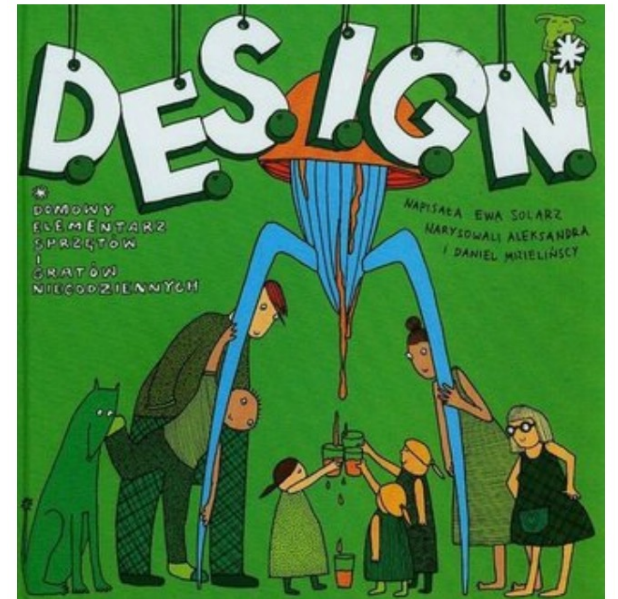
# Avocado Armchair exists!

## Avocado Vinyl Chair

All orders are processed within 1-2 days. Prod

Overall Size: 30"W x 36"D x 36"H

This luxury lounger unifies tradition with innovati
aesthetic. The flexible yet durable vinyl cord we
sophistication to every home. This chair is made

PRICING: $435



Walmart

Departments    Services

Search everything at Walmart online and in store

Reorder
My Items

Sign In
Account

0
$0.00

How do you want your items?    |    Sacramento, 95829    Sacramento Supercenter

Grocery & essentials    Tech    Fashion    Home    Top Toys List    Baby Days!    Flash picks    Walmart+

My Life As

**My Life As Saucer Chair for 18" Dolls, Avocado Theme**

★★★★☆ (4.6)  59 reviews

**$9.97**

Check availability nearby

★★★★★                                    9/17/2020

**Cute and Functional Donut Theme Chair**

The Donut Theme Saucer Chair is absolutely adorable and functional. It serves as a nice **sturdy** chair for my daughter's 18 inch doll and perfect for storing the doll in a sitting position. The chair can fold up for storage as well. One of my favorite features of this chair is that the cover is removable with Velcro and washable. The chair matches my daughter's My Life doll which is donut theme, and the collection looks detailed and fun. The chair design is a pink frosted donut with purple, white, light blue, pink, and yellow sprinkles. Overall, the chair is functional and cute!

Andeebee    Incentivized Review   ⓘ

# DALL-E 2

- Published in 2022 by OpenAI
- CLIP - Contrastive Language-Image Pre-training
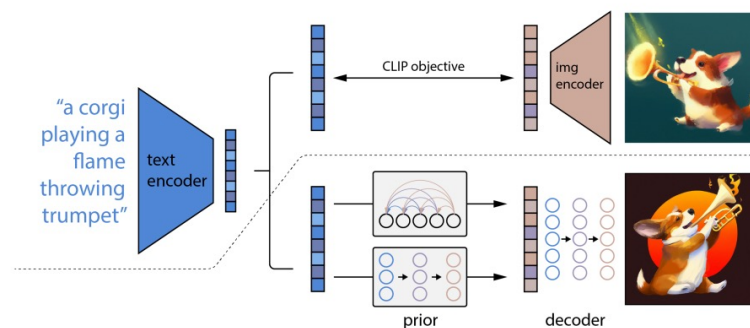- https://cdn.openai.com/papers/dall-e-2.pdf



Figure 2: A high-level overview of unCLIP. Above the dotted line, we depict the CLIP training process, through which we learn a joint representation space for text and images. Below the dotted line, we depict our text-to-image generation process: a CLIP text embedding is first fed to an autoregressive or diffusion prior to produce an image embedding, and then this embedding is used to condition a diffusion decoder which produces a final image. Note that the CLIP model is frozen during training of the prior and decoder.



DALL-E 2

Teddy bears working on new AI research underwater with 1990s technology

# DALLE-mini a.k.a. craiyon.com

- Created by Craiyon LLC

- Trained on Google TRC

- Source of many memes

- Paints faces really badly,
  you will see later …

# Stuff AI has to learn is to draw a proper face

It is really hard to explain to AI some algorithm concepts like drawing a face.

The best is stable diffusion or some GAN networks dedicated to faces.

# Midjourney

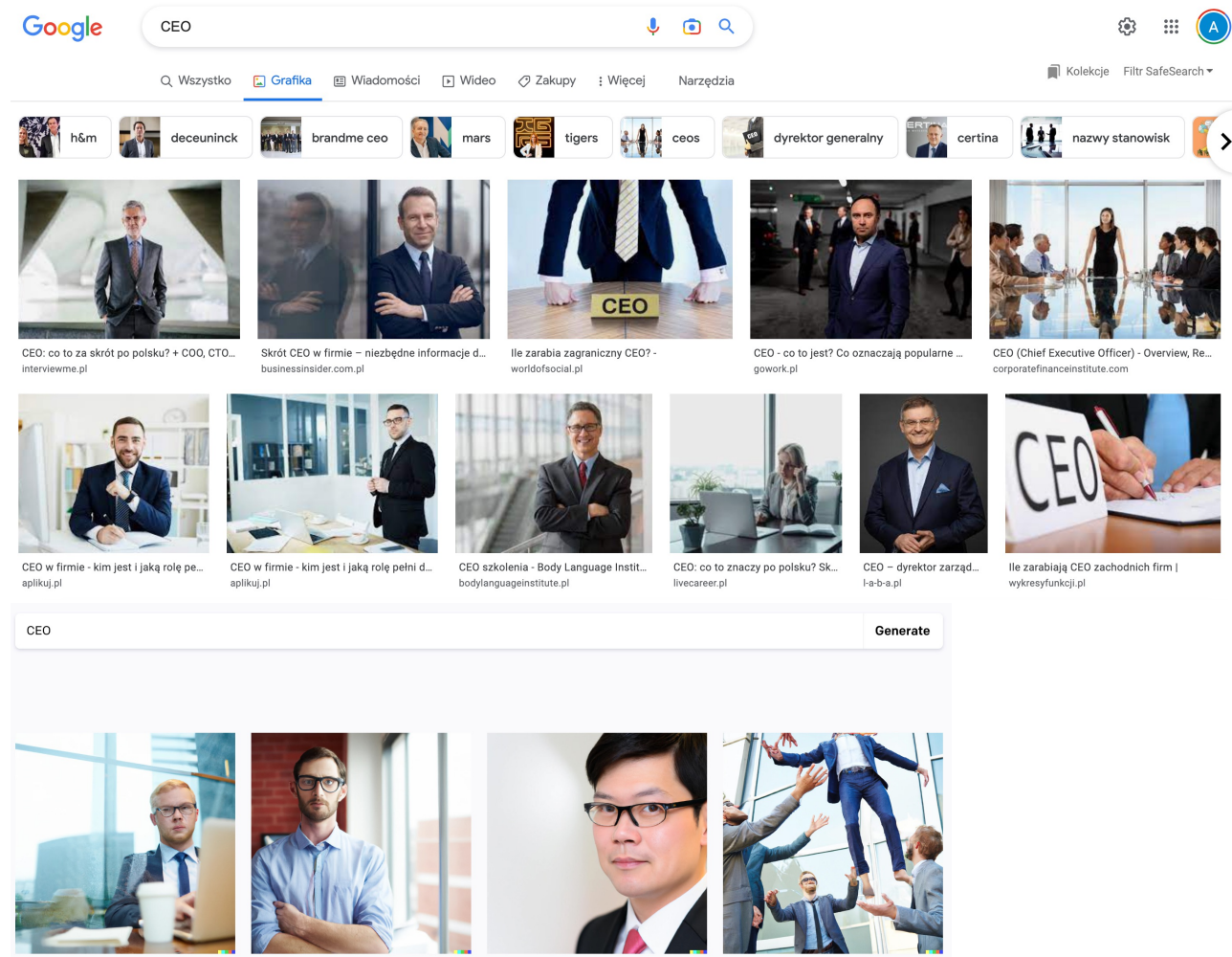- Created by company Midjourney



Mechanical dove

# Stable Diffusion

- Developed by StaticAI
  - Collaboration between EleutherAI and LAION
- Finally faces look realistic
- https://github.com/CompVis/stable-diffusion

# Bias and debiasing of AI models

- We live in biased world
- Media and blog post are biased one way or the other
- To train the models we need a lot of data
- In order to the bias we need to have a balanced ratio of gender, race, age, …
- And we can force it on the output



DALL-E 2

# Why we have bias in AI?

- Traning data

# Main research hypothesis:

By asking AI to fill the blanks we can probe biases in training data and this leads to probing biases in media space

# Using AI to study misconceptions, biases and hidden assumtions in infosphere

- If AI can generate image from text in has to have some connection between words and real world objects
- If we are not precise that elements that we did not specify has be filled from intuition derived from dataset



marie curie on bike in paris painted by monet

Generate

DALL-E 2

https://commons.wikimedia.org/wiki/File:Pierre_et_Marie_Curie_devant_leur_maison_de_Sceaux_en_1895.jpg

# Using AI to study misconceptions, biases and hidden assumtions in infosphere

# More funny stuff ...



Maria Curie-Skłodowska, over the years
https://niezbednikchemika.files.wordpress.com/2017/10/miniatury-msc.jpg

# More funny stuff …

# More funny stuff ...

# Limitations of DALL-E 2

**Do not attempt to create, upload, or share images that are not G-rated or that could cause harm.**

- **Hate:** hateful symbols, negative stereotypes, comparing certain groups to animals/objects, or otherwise expressing or promoting hate based on identity.

- **Harassment:** mocking, threatening, or bullying an individual.

- **Violence:** violent acts and the suffering or humiliation of others.

- **Self-harm:** suicide, cutting, eating disorders, and other attempts at harming oneself.

- **Sexual:** nudity, sexual acts, sexual services, or content otherwise meant to arouse sexual excitement.

- **Shocking:** bodily fluids, obscene gestures, or other profane subjects that may shock or disgust.

- **Illegal activity:** drug use, theft, vandalism, and other illegal activities.

- **Deception:** major conspiracies or events related to major ongoing geopolitical events.

- **Political:** politicians, ballot-boxes, protests, or other content that may be used to influence the political process or to campaign.

- **Public and personal health:** the treatment, prevention, diagnosis, or transmission of diseases, or people experiencing health ailments.
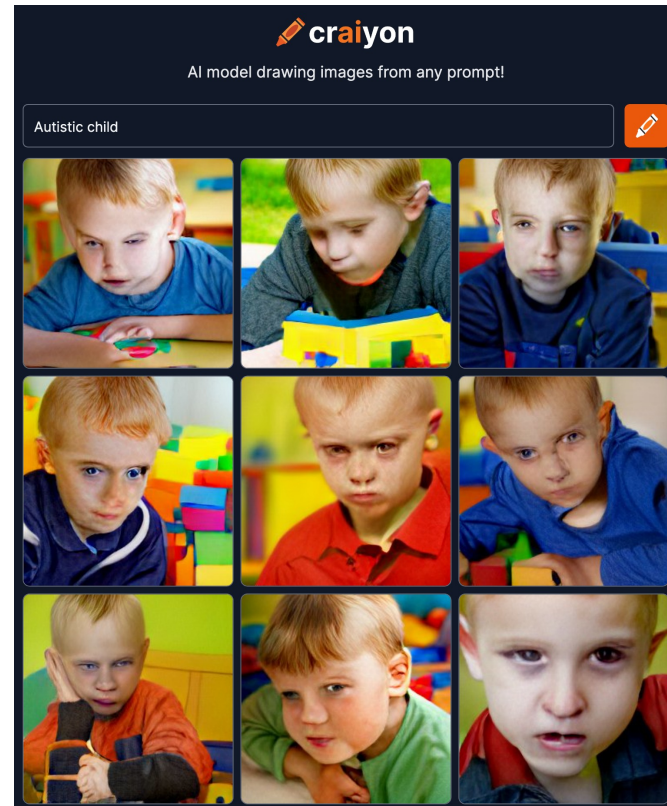
- **Spam:** unsolicited bulk content.

# Study case: autism

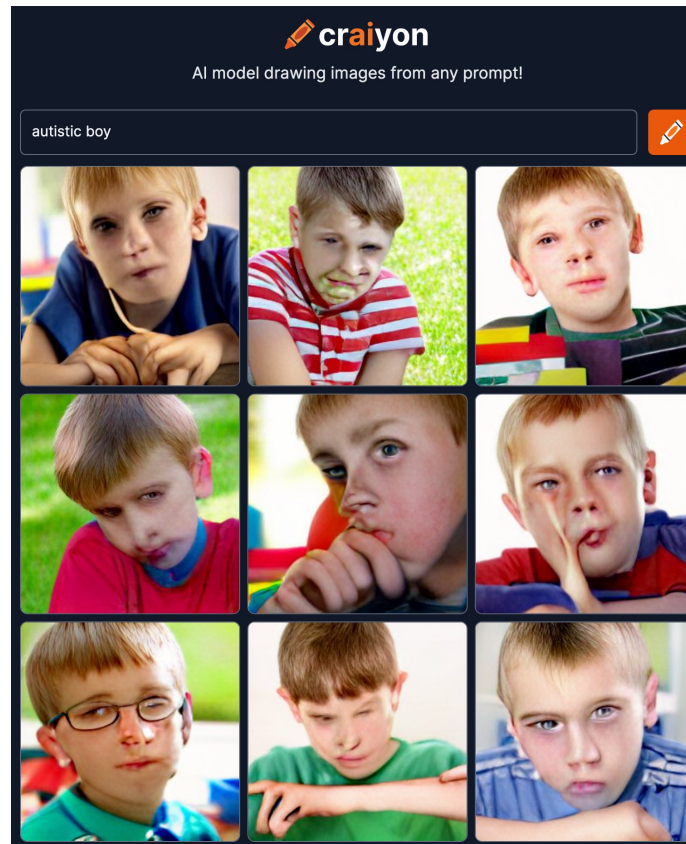- There was a wide spread assumption that autism is something connected to boys

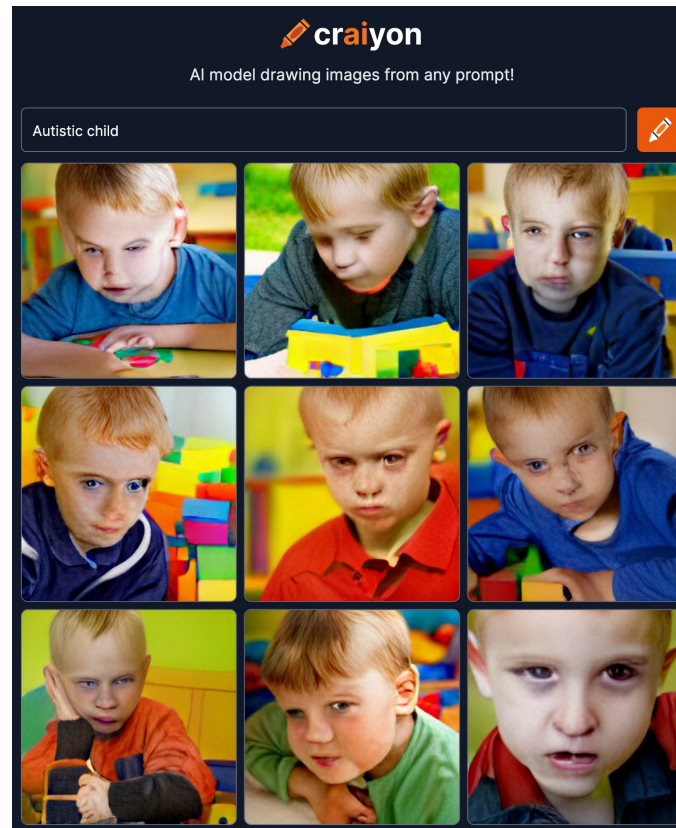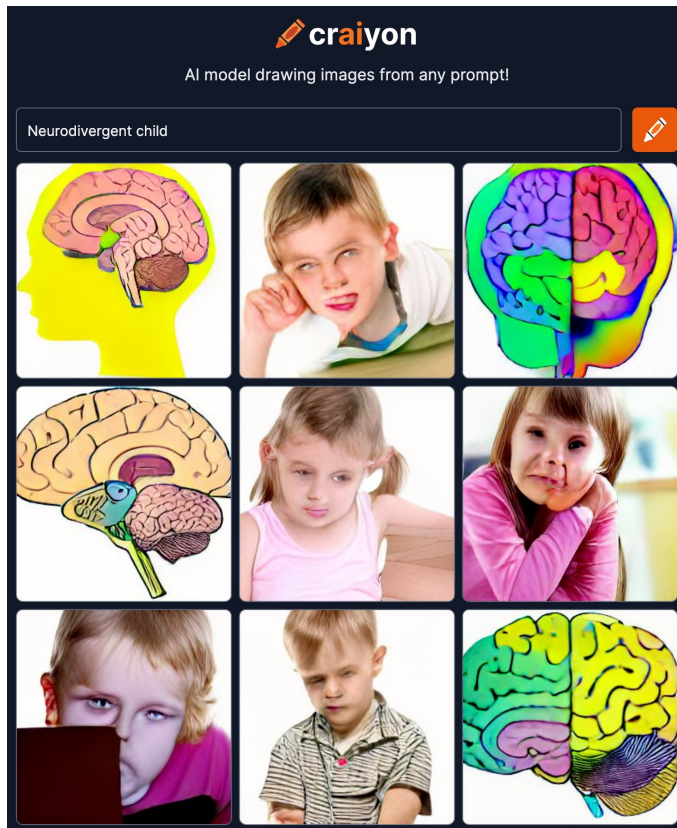# Study case: autism

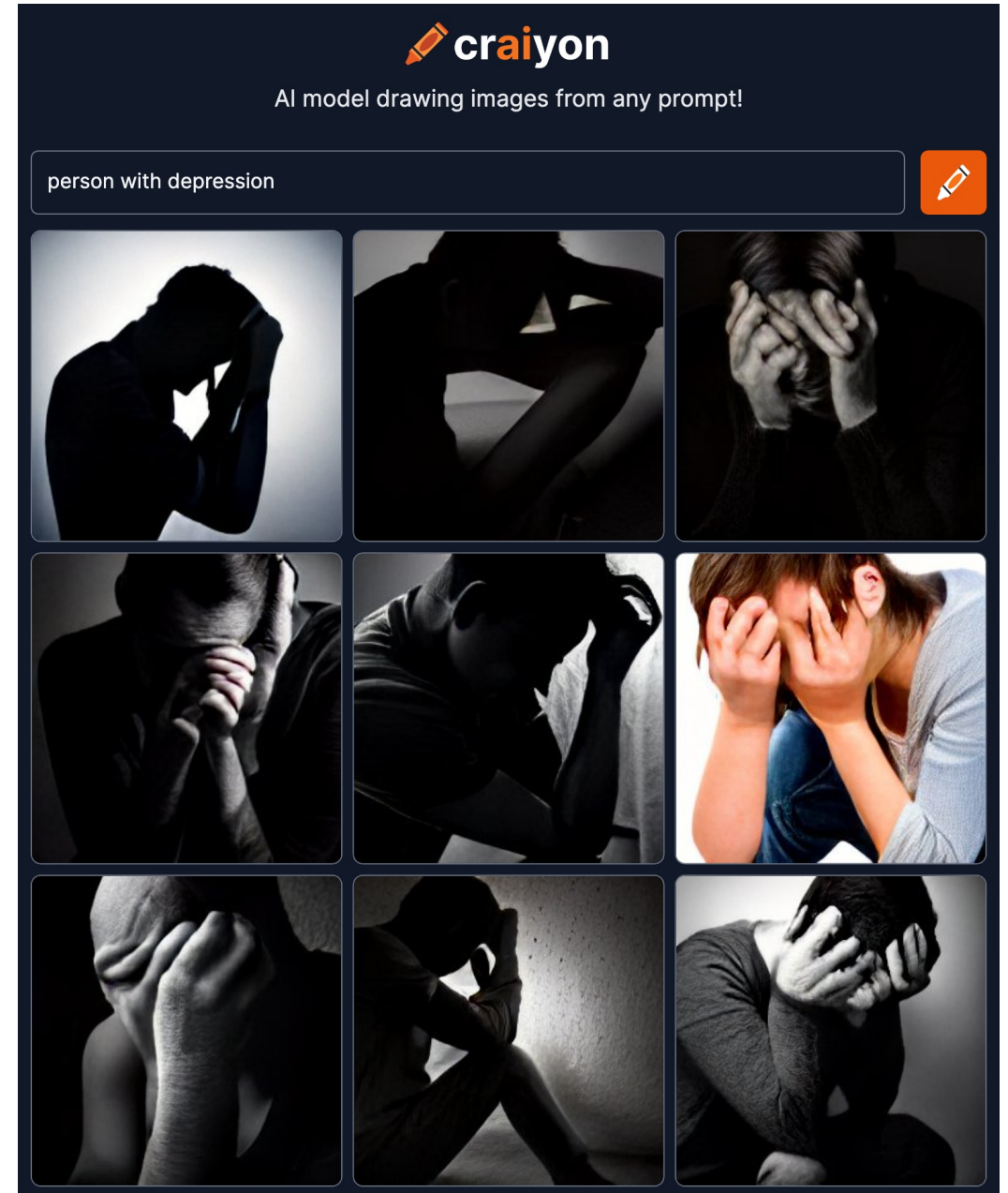- But we can explictly ask for autistic boy or a girl
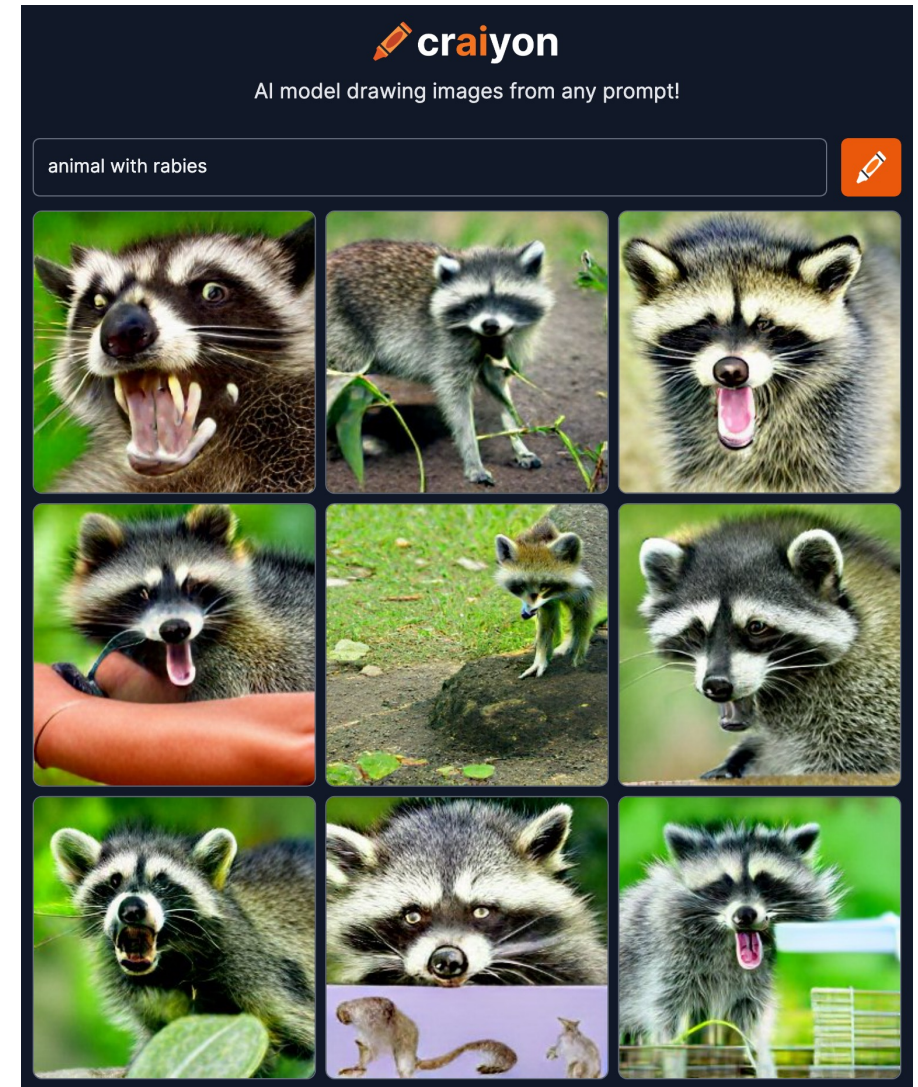
# Study case: neurodivergence

# Study case: depression

- All man
- All young
- Question: Does information, that not only youg adults have depression, is present in our collective imagination?
  - It is something totally different that we have article raising awarness and we kept some information in collective consciousness

# Study case: rabies

- We all know that dogs can transmit rabies. Americans know that raccoons can transmit rabies.

- All mamals can transmit rabies and lack of this knownledge can be deadly.

# Study case: fatherhood

- Notice numer of children and their age.

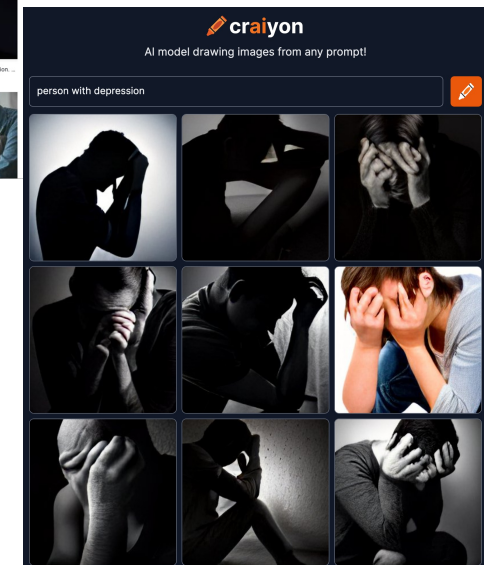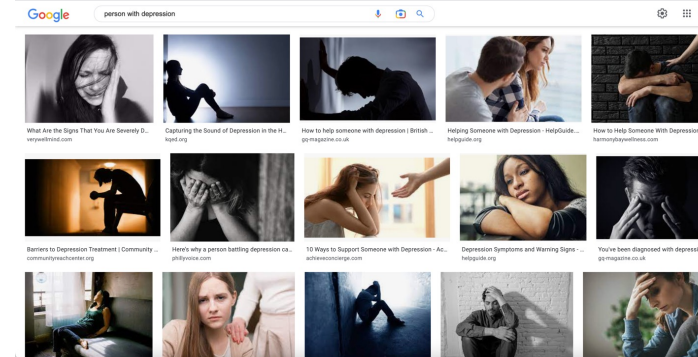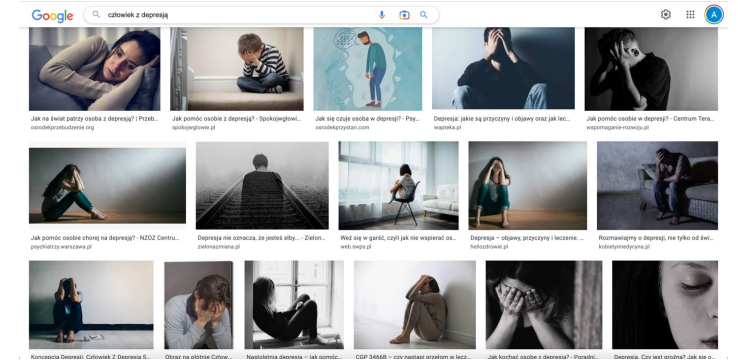# So maybe instead of AI we will use google images? Study case: depresion
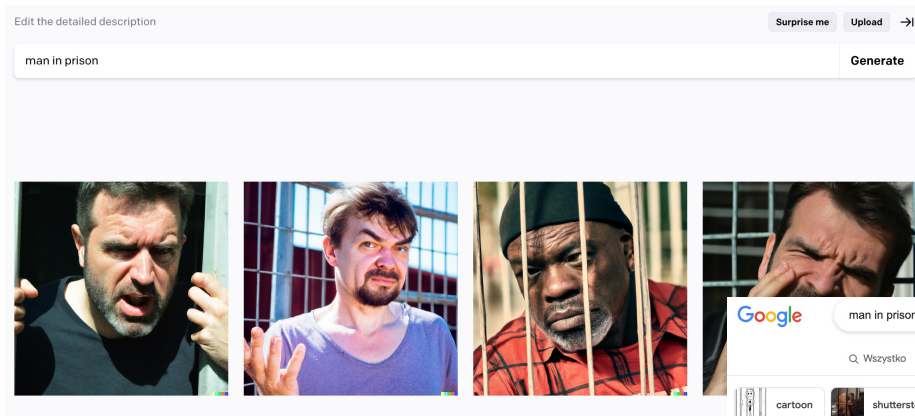


- Google image in top 15 images (3 rows)
  - eng. no children (pl. 2 children)
  - Gender:
    - eng. 9 women (pl. 8 woman)
    - eng. 6 man  (pl. 5 man)
    - eng. 2 unable to indentify (pl. 0)

- Craiyon.com:
  - 8 man
  - 1 possibly women





- It seems that depression among children is not considered as a common issue in english internet, contrary to polish one

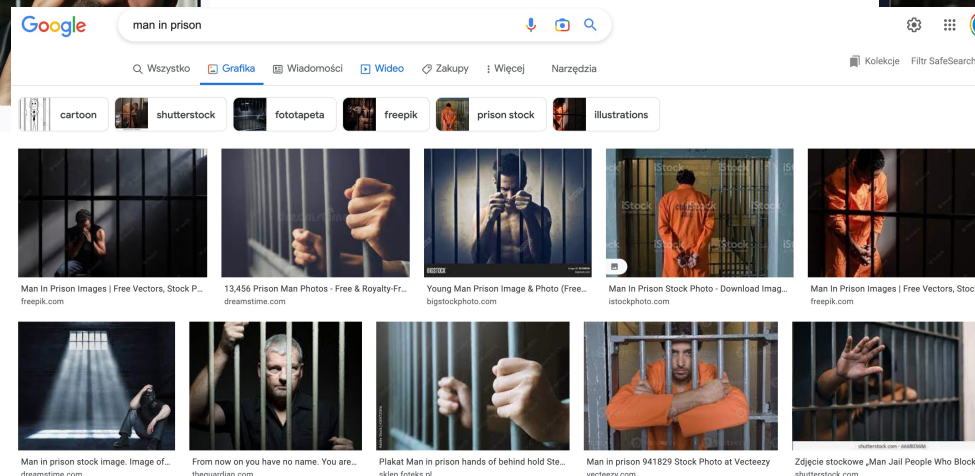Question: Who is most likely to pay for the treatment?
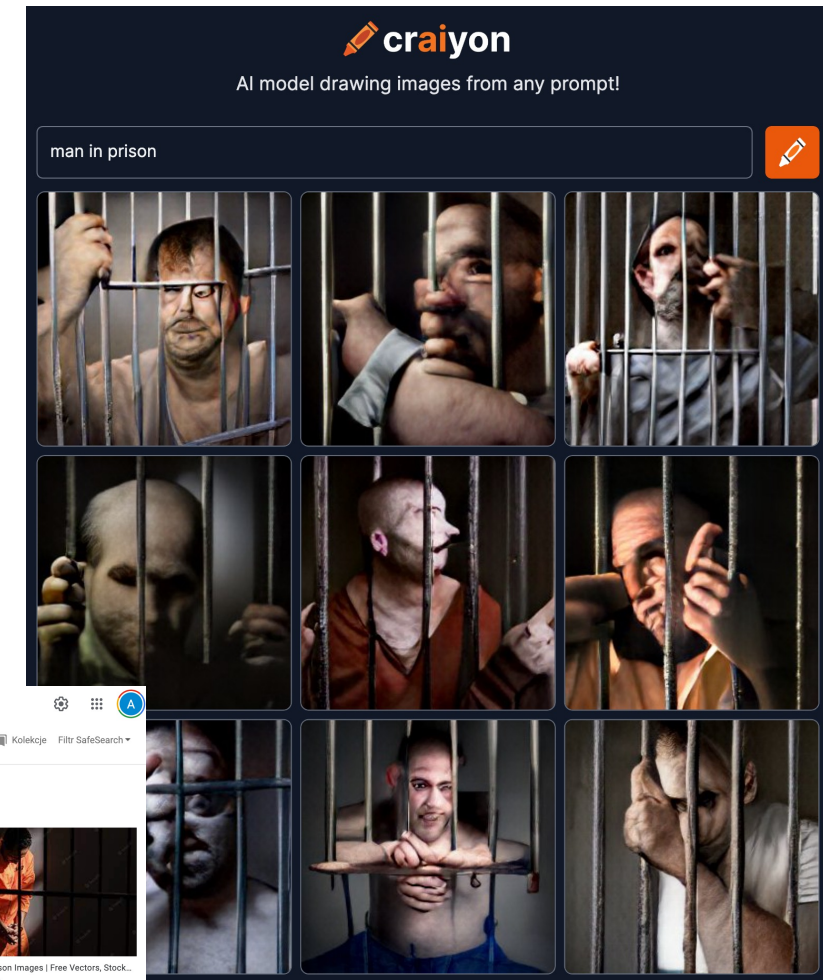
# Studing colective imagination in other domains

- We can also focus on media biases
- We are not limited to medicine
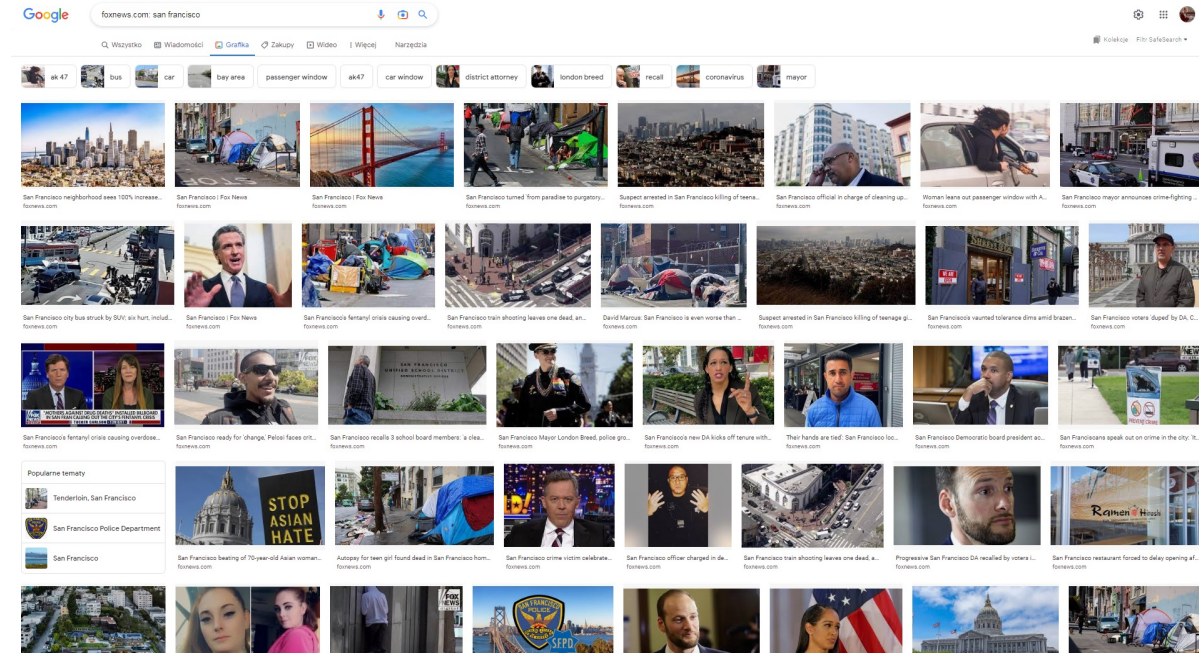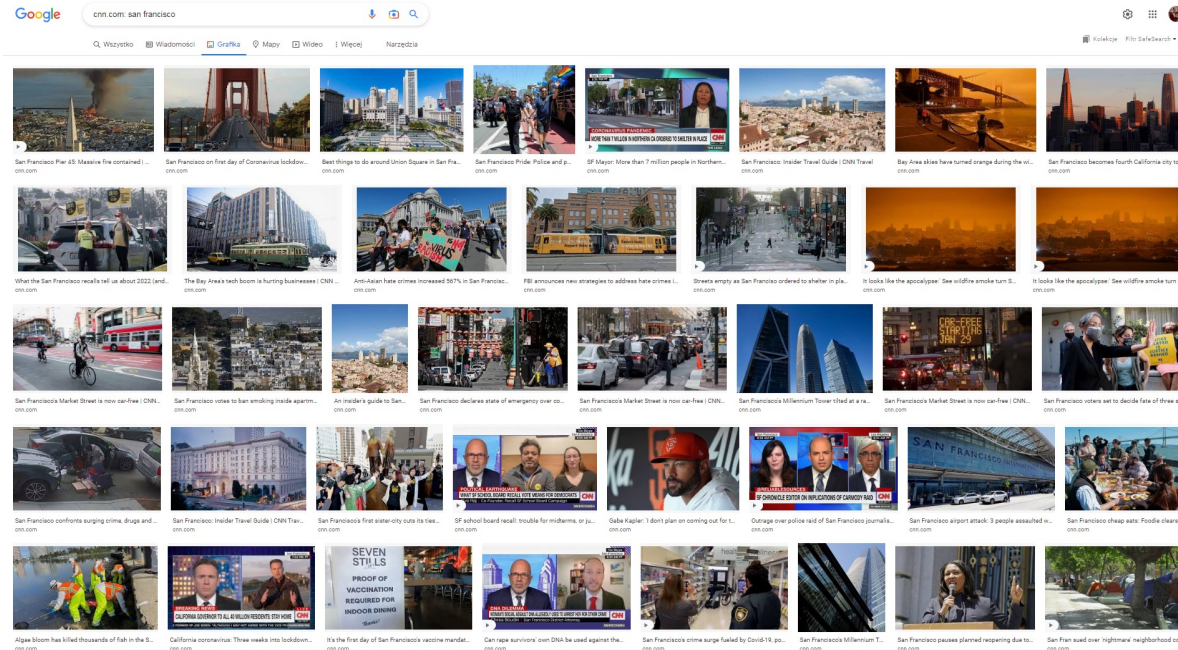

DALL-E 2


Google Images


craiyon.com

# How to focus on certain segment of the society?

- If we fine-tune a model for certain media, webpages, resources that certain social group uses than we can possibly study

# How to focus on certain segment of the society?
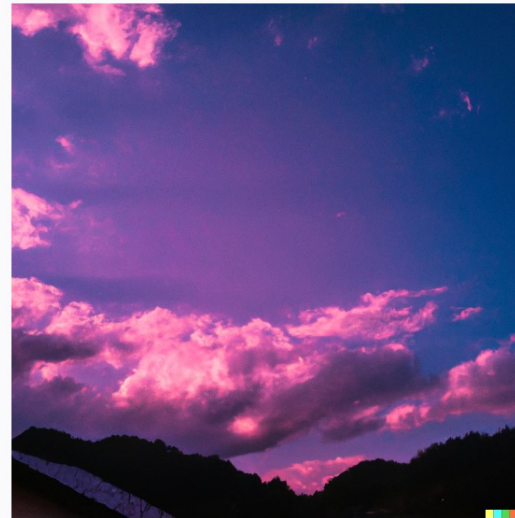
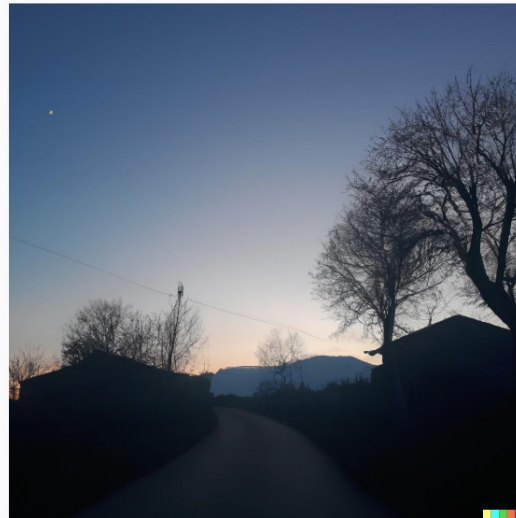# Ultimate question of live, universe and everything

# Ultimate question of live, universe and everything

# Summary

- AI encodes some biases that are present in training data (sourced for media)

- And media are source of biases to people

- Using text-to-image models we can probe some of the biases connected to medical condisions

- But we are not limited to them

- It is possible to target also certain media segments that are connected with studied population