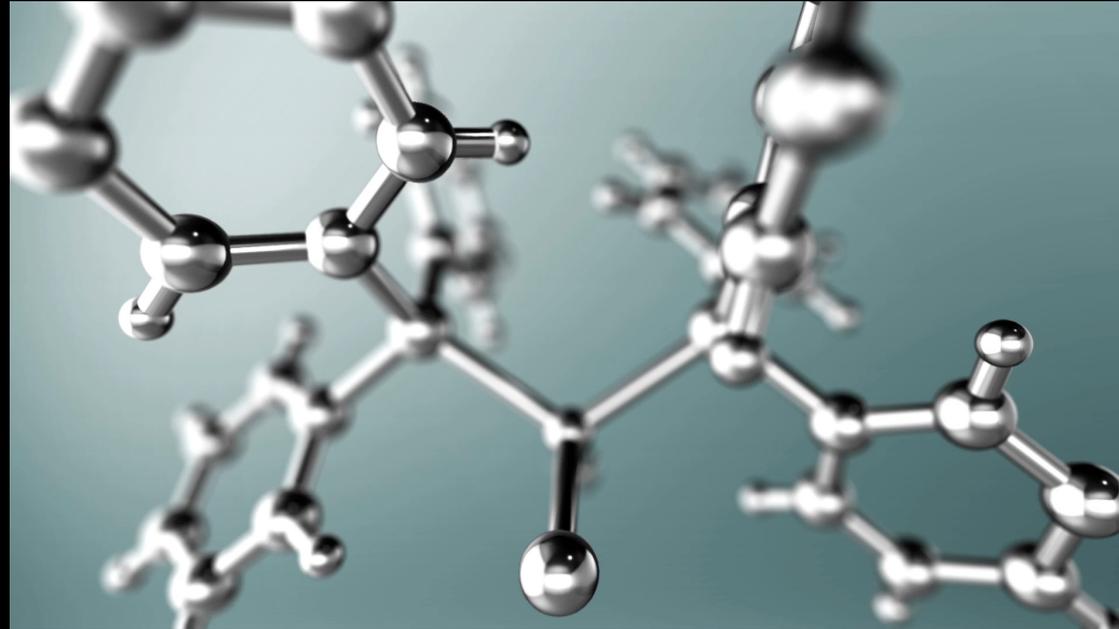


How far can we get in creating a digital replica of a philosopher?



ANNA STRASSER, MATTHEW CROSBY, ERIC SCHWITZGEBEL

OVERVIEW



1

MOTIVATION

Why do philosophers get involved with AI?
Can a language model represent a philosopher?

2

TECHNIQUE

What is a GPT-3?
How to fine-tune GPT-3 & use a fine-tuned GPT-3?

3

LIMITS & RISKS

Limitations of neuronal networks
Legal and ethical questions

4

WHAT WE DID

Fine-tuning | prompt engineering | evaluating

5

RESULTS

Distinguishing & evaluating machine-generated output

MOTIVATION

Why do philosophers get involved with AI?

motivating factors

- understanding how the human mind works
- proving that humans are special
- building new thinking tools
- ...

WHENEVER A MACHINE SUCCESSFULLY SOLVES A TASK
– of which we think that humans solve it because of their intelligence –

then *no* 'intelligence' in the true (deep) sense of the word seems to be necessary for solving this task

BUT

Investigating multiple realizations of solving tasks is a good starting point for further philosophical investigations.

Can we build machines with which we can have interesting conversations?

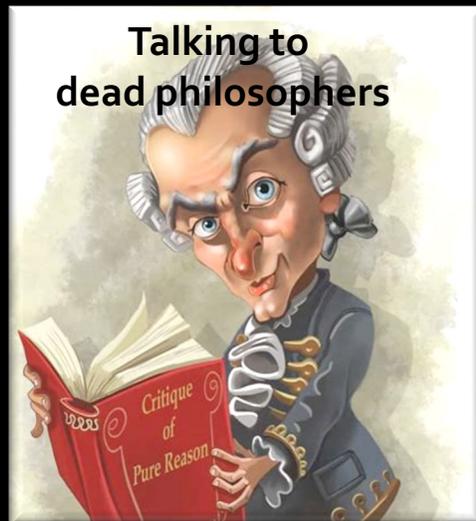
1

AMBITIOUS GOAL:

We want to find out how far we can get in creating a digital replica of a philosopher.

→ investigating of how the best model can be built & used

→ exploring the limits and risks which are accompanied by the creation of digital replicas



several pilot experiments with GPT-3's Currie engine

- fine-tuning with Kant's work in English translation
- fine-tuning with a collection of philosophical blog posts

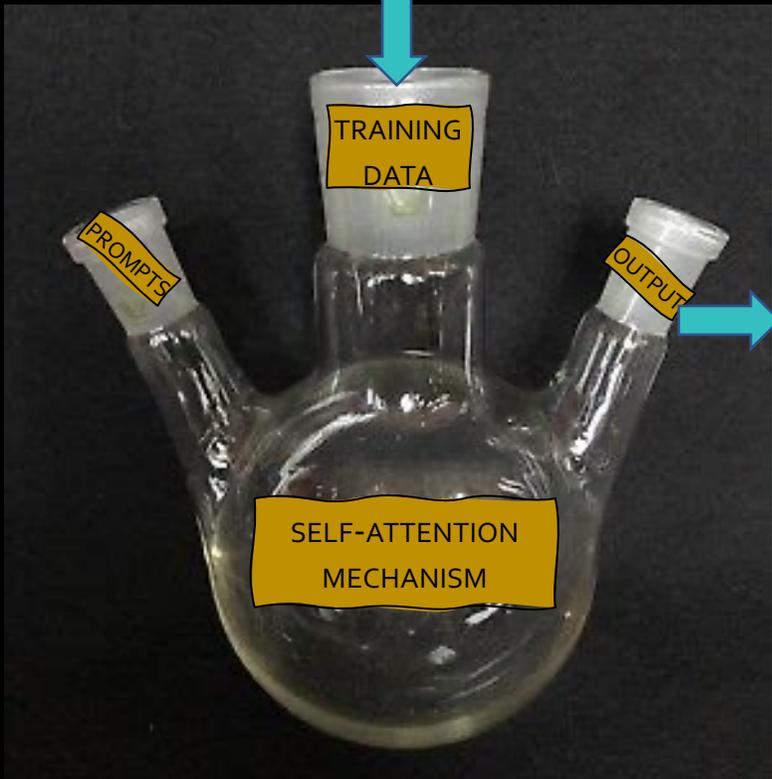


What is a GPT-3?

a neural network trained to predict the next likely word in a sequence

Pre-trained

- 499 billion tokens*
(Common Crawl / WebText / Books / Wikipedia)



Generative

- can generate long sentences
- not just yes or no answers or simple sentences

Transformer

- calculating the probability of the next word appearing surrounded by the other ones

Generative Pretrained Transformer

- a 175 billion parameter language model which shows strong performance on many NLP tasks

*1 token = significant fractions of a word (on average 0,7 words per token)



Applications using GPT-3

AI|Writer-App

allows people to correspond with replicas of historical figures via e-mail (Mayne A. <https://www.aiwriter.app>)

Applications useful for code generation

Codex model by GitHub Copilot using GPT-3 to translate conventional language into formal computer code (Langston 2021)

GPT-3 can be used to convert common natural language expressions into legal language and vice versa

Online documentation of Open-AI.

chatbots can be trained to speak in the tone of particular corporations

Online documentation of Open-AI.

For an overview of applications, see <https://beta.openai.com/examples>

According to an opinion piece in *MIT Technology Review*, GPT-3 is "**shockingly good – and completely mindless**" (Heaven 2020)

Beware of limits & risks

GPT-3 IS A DEEP LEARNING SYSTEM

- deals only with text
- limited input & output sizes (2048 linguistic tokens, roughly 1500 words)
- lacks any form of memory
 - we cannot expect GPT-3 to succeed in text-related tasks which require a larger amount of context knowledge than can be captured by the limited input size
- problems with reliability & interpretability
 - we cannot presuppose that all outputs of a GPT-3 will be acceptable



OFTEN, WE CANNOT PREDICT HOW SUCH A SYSTEM WILL REACT TO NEW INPUT



BIASES IN THE DATA LEAD THE MODEL TO GENERATE STEREOTYPED OR PREJUDICED CONTENT

- BE AWARE OF BIASES & OFFENSIVE LANGUAGE OF THE TRAINING DATA
- REMAIN SKEPTICAL BECAUSE OUTPUTS CAN BE SUBTLY FLAWED OR UNTRUE



NEVER USE THEM IN CONTEXTS WHERE AN INCORRECT OUTPUT IS ETHICALLY QUESTIONABLE

ABOUT LIMITS & RISKS
Bommasani et al. (2021), Floridi & Chiriatti (2020)



More Limitations

Discrete language tasks

- notable gaps in reading comprehension & comparison tasks

Text synthesis

- starts repeating itself

Tasks that empirically benefit from bidirectionally

- fill-in-the-blank tasks

tasks presuppose looking back and comparing two pieces of content, tasks that require re-reading or carefully considering a long passage

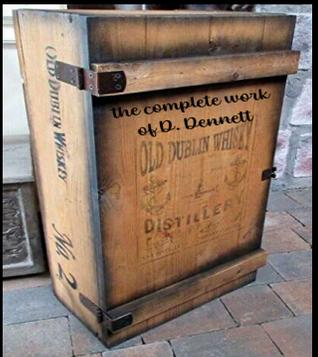
Expensive & inconvenient

- very computing-power hungry model
- become “overqualified” for some specific tasks

**ANALYZING ALL KINDS OF SHORT-COMINGS MAY TELL US SOMETHING ABOUT OURSELVES.
MY DIRECT OBJECTIVE IS NOT TO DEVELOP A SUPER-INTELLIGENT MACHINE.**

What we did

TRAINING DATA:
13 books & 269 articles / interviews



Name	Änder...	datum	Größe	Art
(110)	19.11.21	7 KB	Text	
(111)	19.11.21	21 KB	Text	
(112)	19.11.21	50 KB	Text	
(113)	19.11.21	69 KB	Text	
(114)	19.11.21	59 KB	Text	
(115)	19.11.21	24 KB	Text	
(116)	19.11.21	42 KB	Text	
(117)	19.11.21	24 KB	Text	
(118)	19.11.21	40 KB	Text	
(119)	19.11.21	24 KB	Text	
(1110)	19.11.21	20 KB	Text	
(1111)	19.11.21	28 KB	Text	
(1112)	19.11.21	20 KB	Text	
(1113)	19.11.21	28 KB	Text	
(1114)	19.11.21	34 KB	Text	
(1115)	19.11.21	28 KB	Text	
(1116)	19.11.21	40 KB	Text	
(1117)	19.11.21	67 KB	Text	
(1118)	19.11.21	62 KB	Text	
(1119)	19.11.21	28 KB	Text	
(1120)	19.11.21	10 KB	Text	
(1121)	19.11.21	477 KB	Text	
(1122)	19.11.21	752 KB	Text	
(1123)	19.11.21	28 KB	Text	
(1124)	19.11.21	10 KB	Text	
(1125)	19.11.21	28 KB	Text	
(1126)	19.11.21	18 KB	Text	
(1127)	19.11.21	18 KB	Text	
(1128)	19.11.21	62 KB	Text	
(1129)	19.11.21	14 KB	Text	
(1130)	19.11.21	52 KB	Text	
(1131)	19.11.21	20 KB	Text	
(1132)	19.11.21	40 KB	Text	
(1133)	19.11.21	50 KB	Text	
(1134)	19.11.21	6 KB	Text	
(1135)	19.11.21	20 KB	Text	
(1136)	22.11.21	80 KB	Text	
(1137)	22.11.21	91 KB	Text	
(1138)	22.11.21	6 KB	Text	
(1139)	22.11.21	4 KB	Text	
(1140)	22.11.21	5 KB	Text	
(1141)	22.11.21	12 KB	Text	

txt

max. 1024 tokens

```

{"prompt":"","completion":"<paragraph of text of 1-n.txt>"}
{"prompt":"","completion":"<paragraph of text of 1-n.txt>"}
{"prompt":"","completion":"<paragraph of text of 1-n.txt>"}

```

jasonl training data



FINE-TUNING

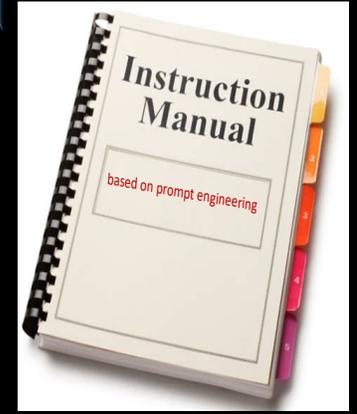
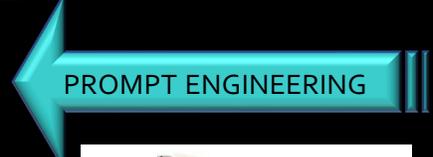
1828 PROMPTS → 3,275,000 TOKENS

OPEN ENDED GENERATION

- leave the prompt empty
- at least a few thousand examples

Anna Strasser

FORMAT OF OUR PROMPTS
Interviewer: [text of the question]
Dennett:



JUST A DROP IN THE OCEAN

