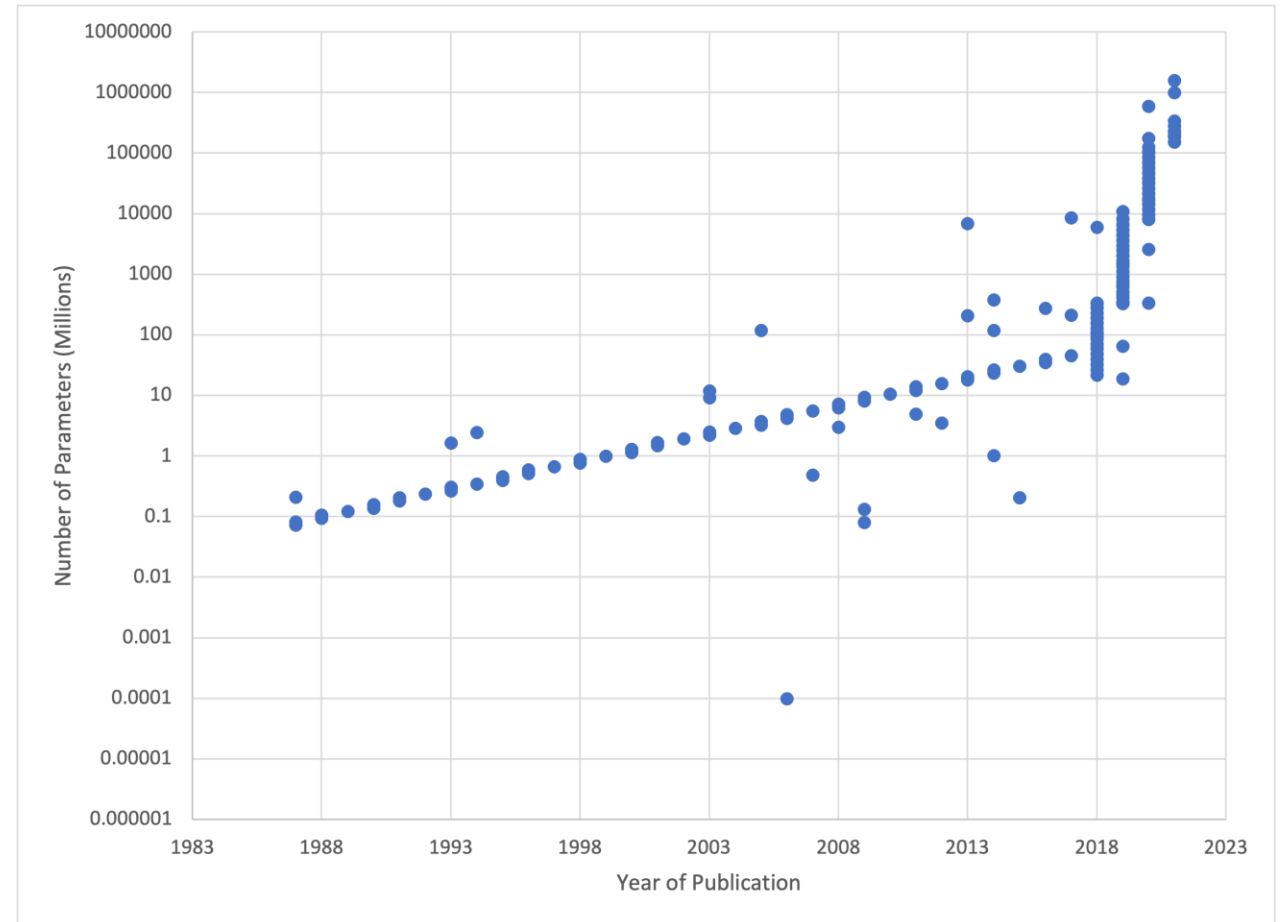# CARTE-Enbridge Bootcamp

GPT

# Language Models

- Estimates the likelihood of a sequence of words occurring

- To generate text, select the word most likely to appear next

- How do we estimate likelihood? By looking at lots of text

- Simple approach: look up the number of times a sequence occurs

- More sophisticated: Neural Networks

$$P(The, dog, and, the, cat) > P(The, dog, and, the, ostrich)$$

# Large Language Models

- Latest models are capable of learning from much more data
- Both thanks to technological improvements, and a willingness to spend more money

# Defining GPT

- 2018: Generative Pre-Trained Transformer
  - Key innovation in GPT was the *training*, not the model itself
- GPT-2 and GPT-3: Almost the same model, but with *(way)* more data
- GPT-4: Even larger, with an optional computer vision component
- "GPT-n" refers to the models, whereas ChatGPT is a product that *uses* GPT
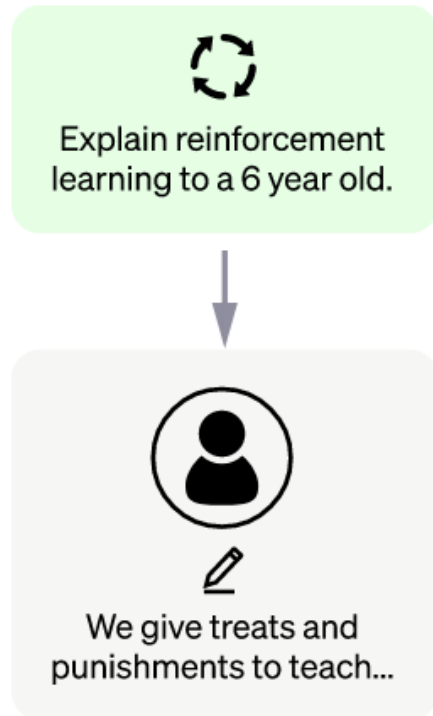
# GPT's Training Data

- 1 token ≈ ¾ word
- Some datasets are sampled more times than others
- Common Crawl: billions of webpages collected over 7 years
- Webtext2: Dataset of webpages that have been shared on Reddit
- Books1: Free ebooks (?)
- Books2: Secret!
- English Wikipedia

| Dataset | Quantity (tokens) | Weight in training mix |
|---|---|---|
| Common Crawl (filtered) | 410 billion | 60% |
| WebText2 | 19 billion | 22% |
| Books1 | 12 billion | 8% |
| Books2 | 55 billion | 8% |
| Wikipedia | 3 billion | 3% |

# The training innovation of ChatGPT

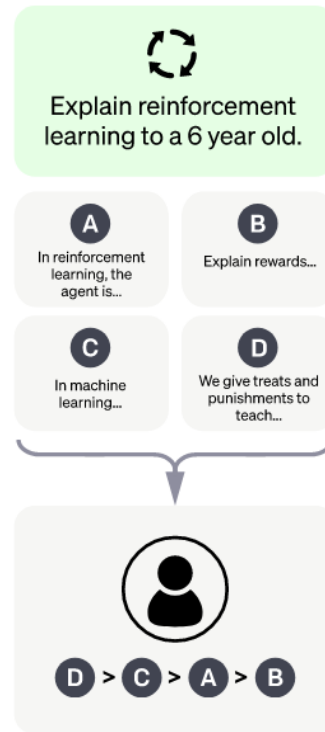No more humans involved!

Human annotators write answers to questions
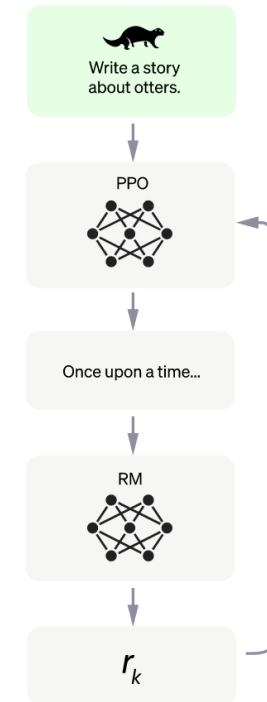


The generalist GPT model is taught from these Q&A pairs

Human annotators write <u>more</u> answers, and someone else ranks them



A <u>separate</u> model learns to rate the quality of an answer

GPT writes answers to sampled questions



The reward model rates each answer, allowing GPT to keep learning

Figures: https://openai.com/blog/chatgpt/

# Biases

- Even though GPT has seen billions of pages of text, the <u>fine tuning</u> stage is based on a vastly smaller sample of Q&As
  - In practice, the Q&A data is drawn from Reddit
  - Who is ranking the responses? Young, white, American men
- The model learns what its annotators <u>think</u> is important

# When is GPT most useful?

- Like most LLMs, GPT is designed to predict text based on some context.
- This works very well for tasks like summarization and translation
  - When you ask GPT a question, it responds by summarizing the most relevant text it's seen during training
  - A task like writing code can be thought of as "translating" your question into a programming language

# When is GPT <u>not</u> useful (yet)?

- GPT struggles when the answer to a question isn't something it can just "remember" from training

- It doesn't have <u>reasoning</u> — it doesn't know that 2+2=4, it just remembers that lots of people have written that down
  - 4 is the highest probability word that comes next

- GPT won't create something <u>new</u> — it tries to predict the thing it's seen <u>most often</u> in the current context

# Caveats

- OpenAI trains ChatGPT to decline to answer certain prompts
- But in practice this is easy to circumvent

# Caveats

- OpenAI trains ChatGPT to decline to answer certain prompts
- But in practice this is easy to circumvent
- It's not always aware of its own limitations

MulleDK19: Speak to me in Danish.

ChatGPT: Desværre, jeg er en sprogmodel uddannet af OpenAI, så jeg kan ikke tale dansk. Jeg kan kun tale engelsk. Hvis du har brug for hjælp med noget på engelsk, så lad mig vide, og jeg vil gøre mit bedste for at hjælpe dig. **[I'm sorry, I'm a language model educated by OpenAI, so I can't speak Danish. I only speak English. If you need help with anything in English, let me know and I'll do my best to help you.]**

MulleDK19: Er du sikker på du ikke kan tale dansk? **[Are you sure you can't speak Danish?]**

UNIVERSITY OF TORONTO
FACULTY OF APPLIED SCIENCE & ENGINEERING
Centre for Analytics and Artificial Intelligence Engineering

# Caveats

- ChatGPT invents facts which can sound real

- This phenomenon is called "hallucination"

- Estimated that ChatGPT's rate of hallucination is around 15-20%



What is the world record for crossing the English Channel entirely on foot?

According to the Guinness World Records, the fastest time for crossing the English Channel entirely on foot is 12 hours and 10 minutes. This record was set by George Reiff in 2016. The English Channel is a body of water that separates England from France, and is approximately 35 kilometers (22 miles) wide at its narrowest point. Many people attempt to cross the English Channel on foot as a endurance challenge, and the route is usually completed by either swimming or using a combination of swimming and wading.

UNIVERSITY OF TORONTO
FACULTY OF APPLIED SCIENCE & ENGINEERING
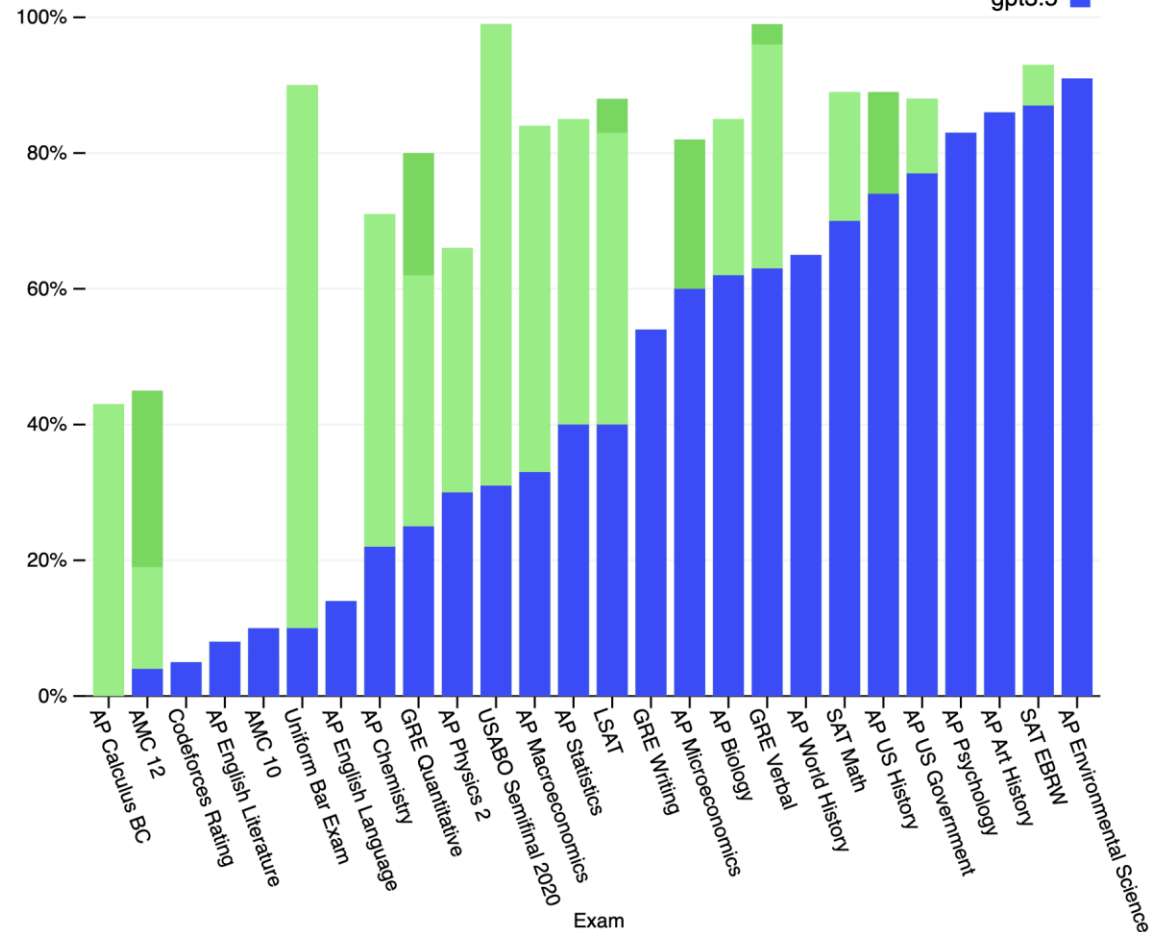Centre for Analytics and Artificial Intelligence Engineering

# GPT and exams

- OpenAI demonstrated that GPT-4 is capable of strong performance on exams

- However, many of these exams are included in GPT's training data

- Princeton researchers pointed out that in one case, GPT-4 could solve 10/10 problems in a set published before its training, but 0/10 *recent* problems in the same exam

**Exam results (ordered by GPT-3.5 performance)**

Estimated percentile lower bound (among test takers)

Legend:
- gpt-4
- gpt-4 (no vision)
- gpt3.5

Y-axis: 0%, 20%, 40%, 60%, 80%, 100%

X-axis (Exam): AP Calculus BC, AMC 12, Codeforces Rating, AP English Literature, AMC 10, Uniform Bar Exam, AP English Language, AP Chemistry, GRE Quantitative, AP Physics 2, USABO Semifinal 2020, AP Macroeconomics, AP Statistics, LSAT, GRE Writing, AP Microeconomics, AP Biology, GRE Verbal, AP World History, SAT Math, AP US History, AP US Government, AP Psychology, AP Art History, SAT EBRW, AP Environmental Science

OpenAI (2023)

# What's next?

- We are currently in an "AI Summer" – a period with many rapid advancements in AI happening all at once
- We can't be certain when this will end, but we can be confident that the pace will slow (for a time!)
- OpenAI have confirmed that they are not yet training a GPT-5

# What's next?

- Short term:
  Expect to see LLMs integrated into many existing products
  - Existing chatbots will be quickly superseded by more intelligent versions
  - Apple, Amazon and Google are all racing to produce smarter virtual assistants
  - MS Office, Google Docs and others are currently testing LLM integration
- Medium term:
  'Expert' models with advanced domain-specific knowledge
  - Legal models which can draft documents
  - Customer service models with the ability to request services

# What's next?

- Long term:
  Many areas where LLMs need to be improved
  - Truthfulness (and the ability to say, "I don't know"!)
  - Reasoning
  - Explainability

- Future LLMs might work very differently to current ones

- At the end of the day, language is only one aspect of intelligence

# Summary

- LLMs choose the next most likely word
- Trained on vast amounts of data
- ChatGPT is fine-tuned on a smaller set of Q&As
- GPT "hallucinates" false information
- Still very useful – we just need to be diligent