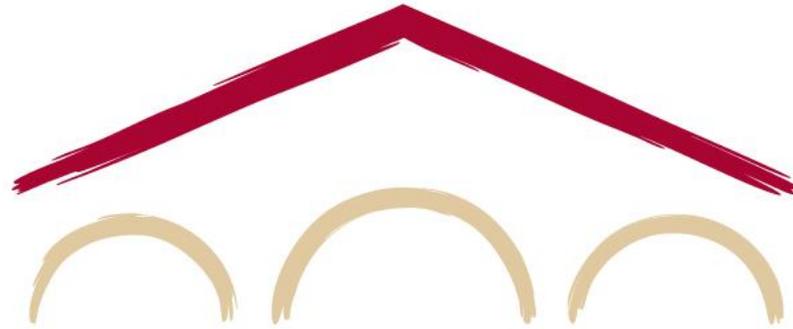# Generative AI: From fluent text to agentic workflows to an AI helpforce

## Christopher Manning

@chrmanning ❀ Stanford NLP Group ❀ @stanfordnlp

Director, Stanford AI Lab ❀ Associate Director, Stanford HAI

# There have been 60+ years of AI research ...

See our short video on YouTube: **AI at Stanford 1962–2022**

# The world changed with OpenAI's ChatGPT

November 2022

## ChatGPT 4o ⌄

Write a sonnet about the mission of ServiceNow in which each line begins with the letter B

Building a world where workflow reigns supreme,
Bringing efficiency to every task,
Bold in pursuit of innovation's dream,
Bridging the gap when teams and systems ask.

Brightening the path for enterprise's way,
Boosting solutions, making work more clear,
Balancing speed and skill in every day,
Believing in the future we hold dear.

Better connections forged through seamless flow,
Beneath the surface, magic starts to hum,
Breaking down silos where ideas grow,
Becoming the platform where change can come.

By driving purpose, we fulfill our vow,
Building the future, this is ServiceNow.

# "Large language models"(LLMs) trained by Mad Libs

"Unsupervised" or "self-supervised" learning, using existing pieces of text

## Ford plans to _____
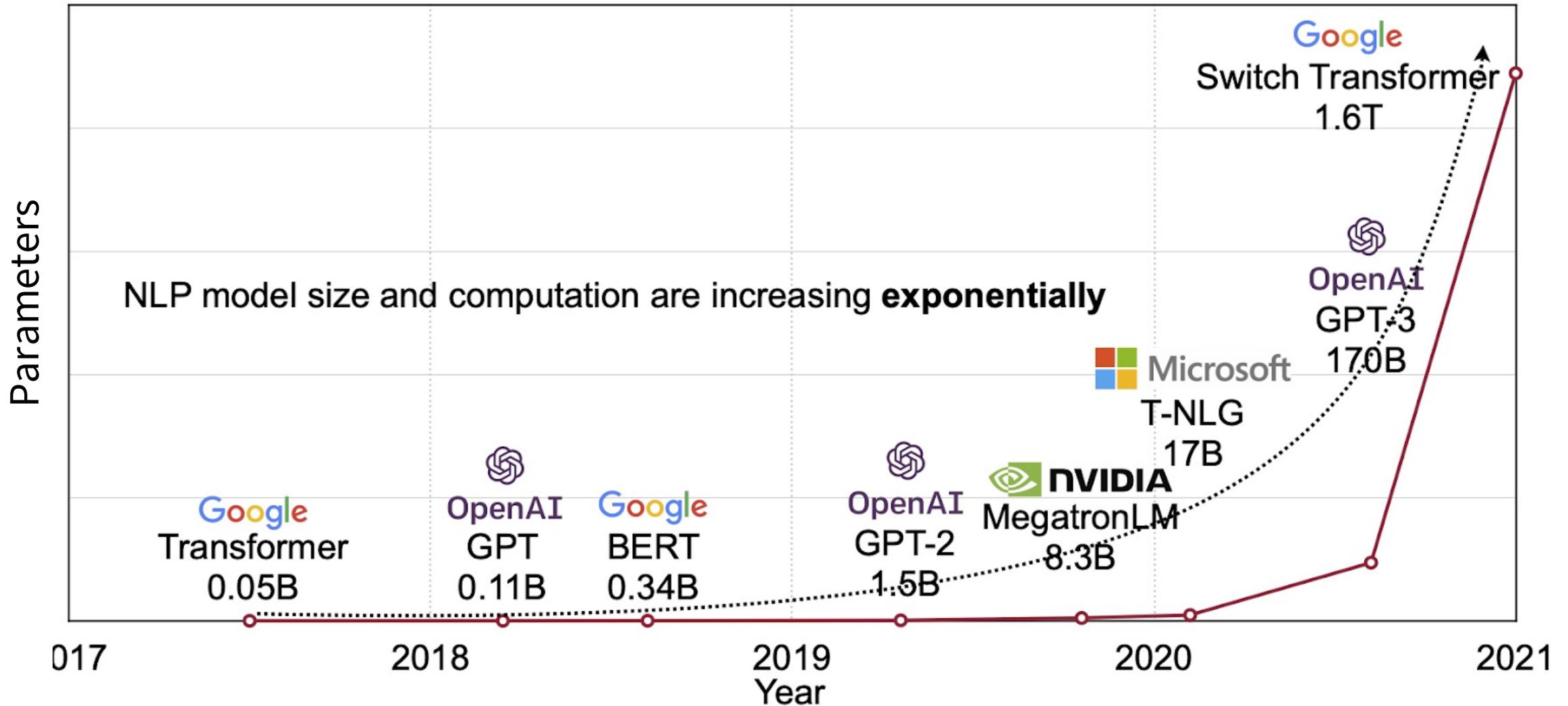


What Happens When a Unicorn Poops?

Unicorns aren't like other_____ (plural noun); they're_____ (adjective). They look like _____ (plural noun; animals), with_____ (plural noun) for feet and a_____ (adjective) mane of hair. But unicorns are_____ (color) and have a_____ (adjective) _____ (noun) on their heads. Some_____ (plural noun) don't believe unicorns are_____ (adjective) but I believe in them. I would love to_____ (verb) a unicorn to faraway_____ (plural noun). One thing I've always_____ (verb -ed) about is whether unicorns_____ (verb) rainbows, or is their _____ (noun) _____ (adjective) like any other animal's?
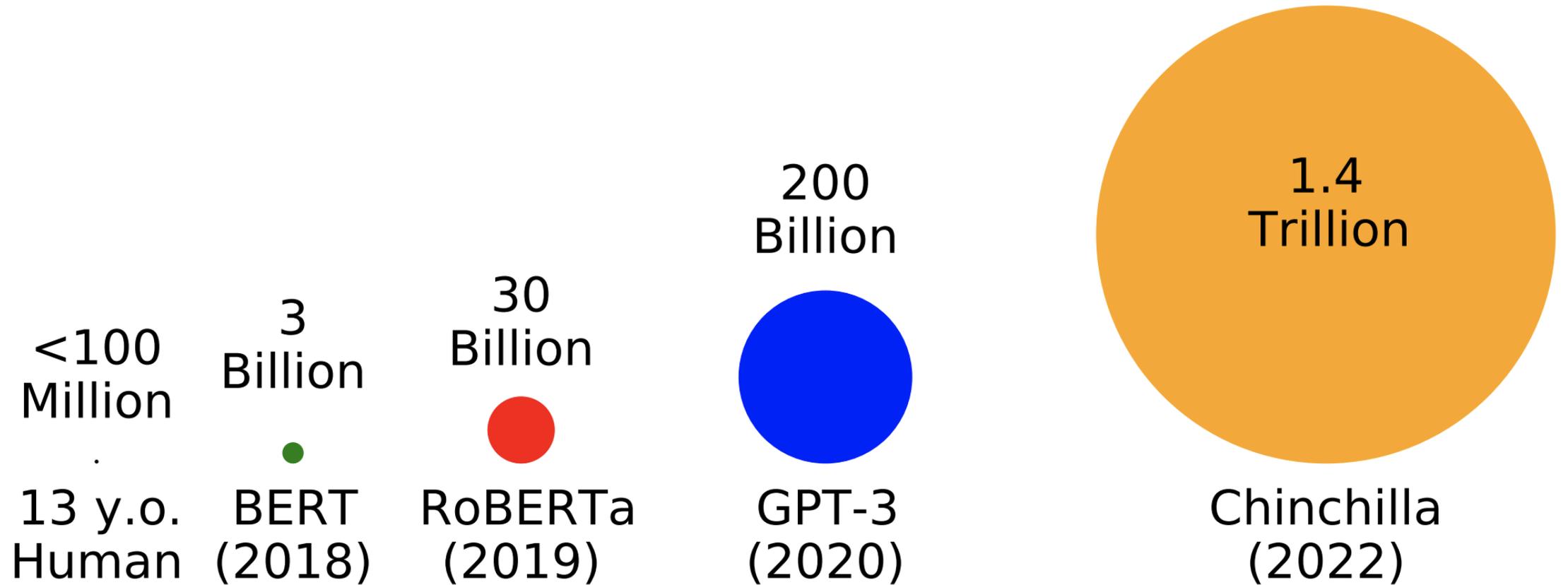
5

# The History of Large Language Models

1998     CPAT-Tree-Based Language Models with an Application for Text Verification in Chinese. ROCLing 1998. First use of "LLM" trigram I know of; 200M word corpus

2000     A Neural Pro[...]Vincent NIPS 2000
         First neural l[...], 31K vocab

**Not enough compute!**

2007     Large Language Models in Mach[...]
         **Dean**. EMNLP 2007. **2 trillion to[...]**

**Not enough model flexibility!**

2018     **GPT** (Radford, Narasimhan, Sali[...]
         Lee & Toutanova). 3.3 billion tok[...]

**Not enough data!**

2020–    100+ billion parameter neural la[...]
         GPT-3, GPT-4, PaLM 2, Llama 3, [...]

**LLMs go brrr!**

# Larger and larger models requiring more compute

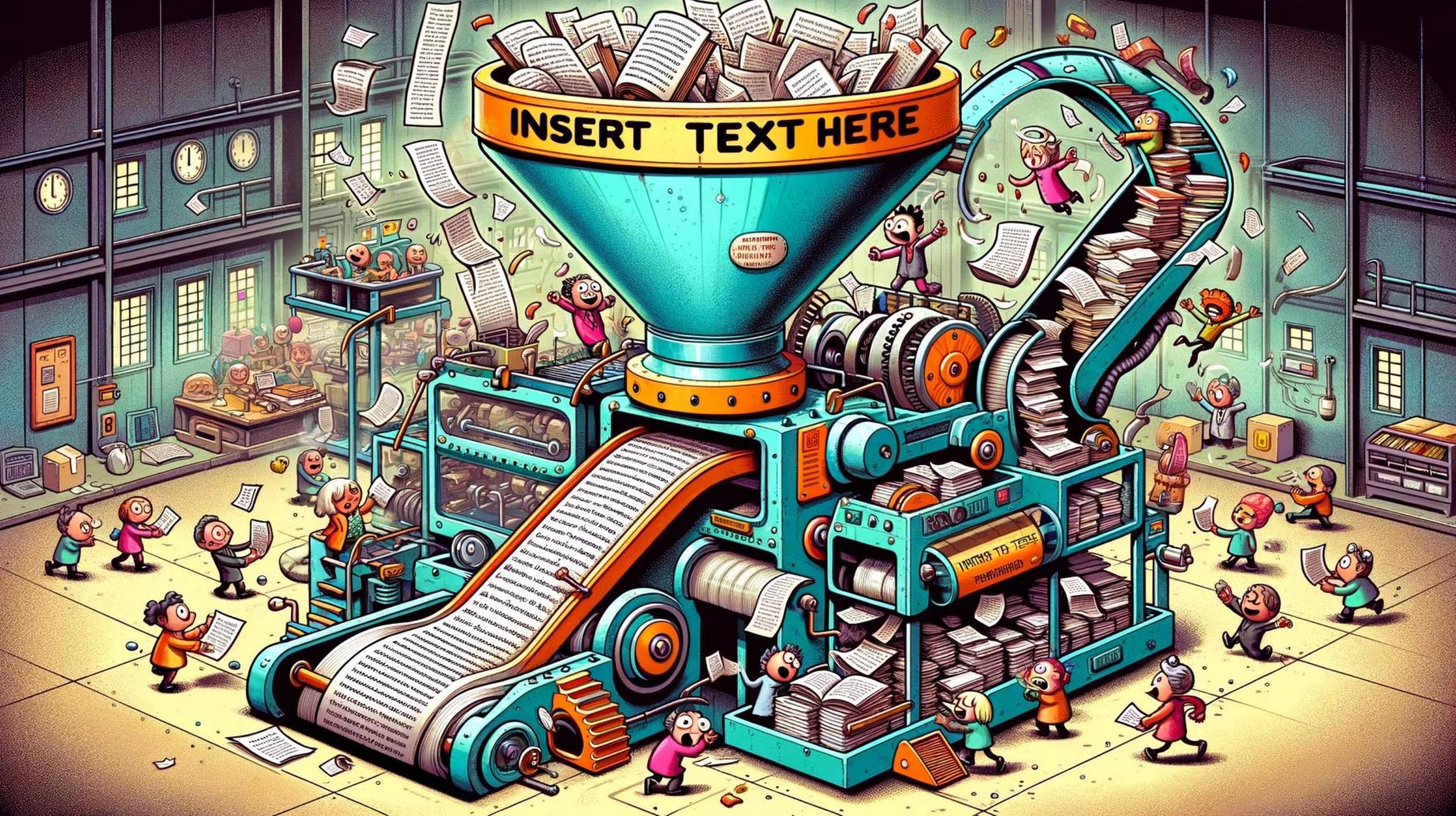

NLP model size and computation are increasing **exponentially**

- Google Transformer 0.05B
- OpenAI GPT 0.11B
- Google BERT 0.34B
- OpenAI GPT-2 1.5B
- NVIDIA MegatronLM 8.3B
- Microsoft T-NLG 17B
- OpenAI GPT-3 170B
- Google Switch Transformer 1.6T

Axis: Parameters vs Year (2017, 2018, 2019, 2020, 2021)

# Trained on more and more data

<100
Million

3
Billion

30
Billion

200
Billion

1.4
Trillion

13 y.o.
Human

BERT
(2018)

RoBERTa
(2019)

GPT-3
(2020)

Chinchilla
(2022)

# tokens seen during training

https://babylm.github.io/

## **Produce a qualitative shift in what emerges**

It no longer appears to be statistical prediction but rather

a machine that understands language!

And the world!

# The transformer in practice (Meta's LLaMA 405B)

- **Total training data** – 15.6 Trillion tokens

- **Architecture** – 405 Billion params
  - 126 transformer blocks
  - Model dimension: 16,384 dim vectors
  - Attention heads: 128

- **Compute** – 3.8e25 FLOPs
  - 1.1 million H100-days ($~25 M)
  - ~ 2 months on 16k H100s

# Getting LLMs to do what you actually want

GPT-1/2/3 rewrote how researchers think about modeling language

But they didn't become a **household name** (like ChatGPT in late 2022)

## *Why?*

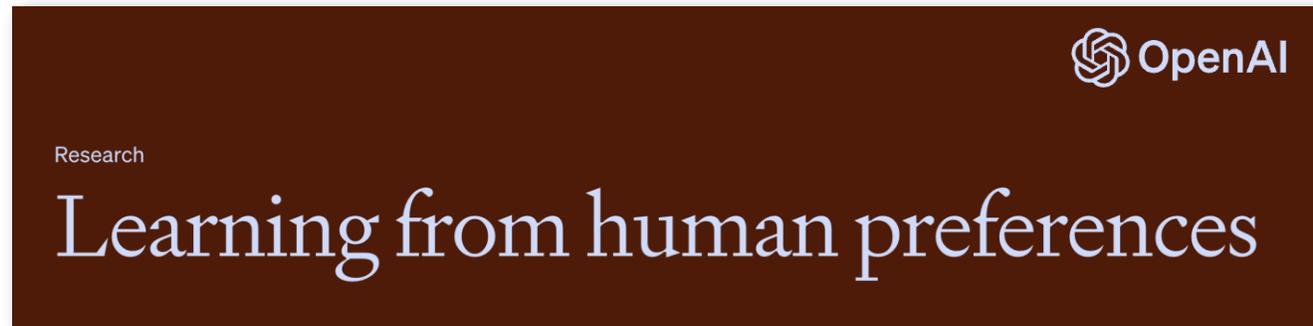Who is the US President?

Ask GPT-3 who the US President is:

Well, no, that's not really what I wanted

# Getting LLMs to do what you actually want

Models that **understand what people want** are profoundly more useful!

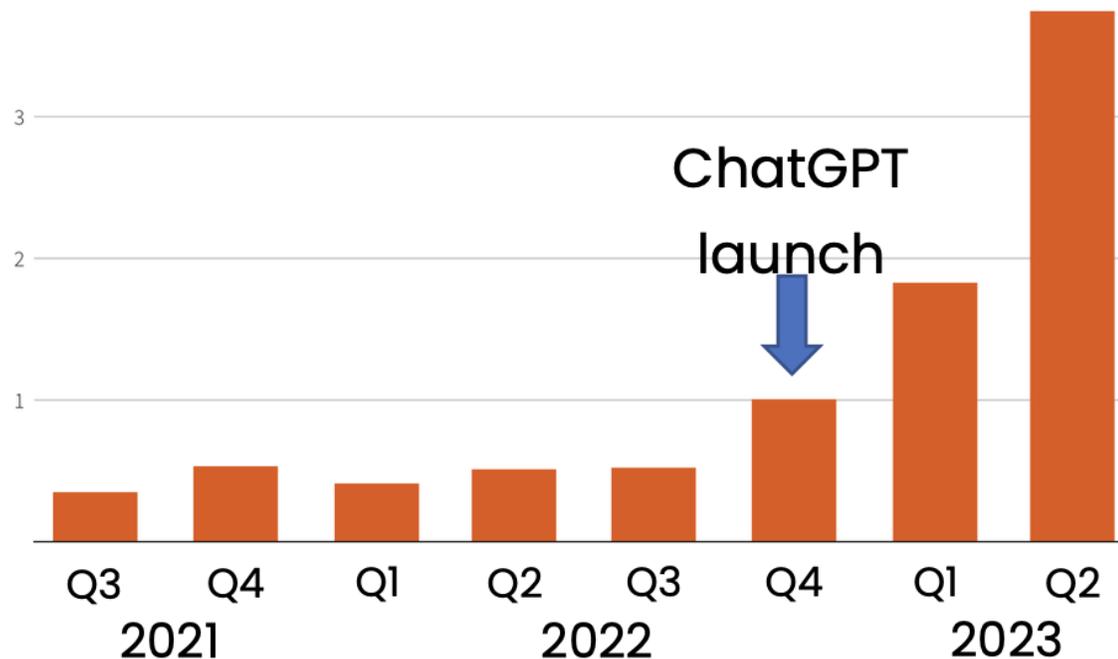ChatGPT accomplishes this behavior through a second phase ("post-training") using:

- **Reinforcement Learning from Human Feedback (RLHF)**

- **A (very) large amount of annotated interaction data**



https://openai.com/research/learning-from-human-preferences. Accessed 26 January 2024.

# The coming future of generative AI

**Average number of 'AI' mentions per S&P 500 analyst call**



Source: Reuters

Generative AI could

- Add $2.6–$4.4 trillion annually to the economy[1]
- Raise global GDP by 7% in the next 10 years[2]

[1] McKinsey; UK total GDP 2021: $3.1 trillion
[2] Goldman Sachs

# Example uses of LLMs: 1. Improving text

Can you improve this introduction to the ServiceNow Research Team:

We make workflows smarter

The ServiceNow Research team does both fundamental and applied research to futureproof AI-powered experiences for all users of the Now Platform®. We make workflows smarter, AI workloads more efficient, and are committed to making socially responsible contributions to the AI community. Our innovation is centered on people: helping customers modernize their technology architecture, innovating new business models, improving experiences at work, and driving higher ROI from their technology investments.

# Example uses of LLMs: 1. Improving text

Certainly! Here's an improved version of your introduction:

**Transforming Workflows with Smarter AI**

At ServiceNow Research, we push the boundaries of both fundamental and applied research to shape the future of AI-powered experiences on the Now Platform®. Our mission is to make workflows smarter, enhance the efficiency of AI workloads, and contribute responsibly to the AI community. Innovation drives everything we do, but people remain at the center of our efforts. We empower customers to modernize their technology architecture, create new business models, elevate workplace experiences, and maximize returns on their technology investments.

This version tightens the language, emphasizes the impact of the research, and aligns the message with ServiceNow's mission of driving innovation and user-centricity.

# 2. Getting information from a huge document: AGL

**Please read the 2022 AGL annual report here:**

https://www.agl.com.au/content/dam/digital/agl/documents/about-agl/investors/2022/220819-agl-energy-annual-report-2022.pdf

**What does the AGL report say about the company's alignment with IPCC on Climate Action?**

The AGL 2022 annual report states that the company is committed to "aligning our business with the goals of the Paris Agreement and the IPCC's 1.5°C pathway." AGL has set a target of net zero emissions by 2050, and is currently working on a decarbonisation pathway to achieve this goal. The company's Climate Transition Action Plan includes a commitment to exit coal-fired power generation by FY35. This will bring AGL's Scope 1 and 2 emissions to net zero, which is significant as AGL is currently Australia's largest corporate emitter.

Overall, the AGL 2022 annual report shows that the company is committed to taking action on climate change. AGL has set ambitious targets for reducing its emissions, and is working with stakeholders to achieve these goals. The company's engagement with the Climate Action 100+ initiative is also a positive sign, as it shows that AGL is taking the concerns of investors seriously.

Here are some specific quotes from the report that relate to AGL's alignment with IPCC on climate action:

• "We are committed to aligning our business with the goals of the Paris Agreement and the IPCC's 1.5°C pathway."

• "We have set a target of net zero emissions by 2050, and are currently working on a decarbonisation pathway to achieve this goal."

• "Our Climate Transition Action Plan includes a commitment to exit coal-fired power generation by FY35."

• "We are actively engaged with the Climate Action 100+ initiative."

• "We have a score of 25% on the Climate Action 100+ engagement scorecard."

I hope this helps! Let me know if you have any other questions.

# Getting information from a huge document: AGL

**Are there any comments that raise concerns as to whether AGL is fully committed to a net zero carbon emissions future?**

Yes, there are a few comments in the AGL 2022 annual report that could raise concerns about whether the company is fully committed to a net zero carbon emissions future. These include:

• The company's decision to delay the closure of its Loy Yang A coal-fired power station by two years. This decision was made in response to concerns about the reliability of the Australian electricity grid, but it has been criticized by some as a setback for Australia's climate goals.

• The company's plans to invest in new gas-fired power stations. Gas is a cleaner-burning fuel than coal, but it is still a fossil fuel and its combustion produces greenhouse gases. Some critics have argued that AGL's investment in gas-fired power stations is inconsistent with its commitment to net zero emissions.

• The company's lack of transparency about its decarbonisation pathway. AGL has not yet released a detailed plan for how it will achieve net zero emissions by 2050. This lack of transparency has led some to question whether the company is serious about its climate commitments.

Overall, the AGL 2022 annual report shows that the company is committed to taking action on climate change. However, there are some concerns that could raise questions about whether AGL is fully committed to a net zero carbon emissions future. The company needs to be more transparent about its plans and to take more concrete steps to reduce its emissions if it wants to allay these concerns.

# Programming assistance: Codeium

**https://codeium.com/**

Chat and ask for code! Learn to use APIs, generate unit tests, etc.

# Programming assistance: Amazon Q

Software engineers can achieve ~30% productivity improvements!

It depends on the task and coder:
- < 10% for high-complexity tasks
- 30–100% for less experienced coders or fairly routine tasks

**Andy Jassy** ✓ 📦
@ajassy

One of the most tedious (but critical tasks) for software development teams is updating foundational software. It's not new feature work, and it doesn't feel like you're moving the experience forward. As a result, this work is either dreaded or put off for more exciting work—or both.

Amazon Q, our GenAI assistant for software development, is trying to bring some light to this heaviness. We have a new code transformation capability, and here's what we found when we integrated it into our internal systems and applied it to our needed Java upgrades:

- The average time to upgrade an application to Java 17 plummeted from what's typically 50 developer-days to just a few hours. We estimate this has saved us the equivalent of 4,500 developer-years of work (yes, that number is crazy but, real).

a train going over the Golden Gate bridge

a pencil drawing of a train going over the Golden Gate bridge

A detailed pencil drawing of cars and a train on the Golden Gate bridge

Pictures by OpenAI's DALL-E 3

# Research Lets You Design and See the Future!
# Progress of GenAI Text-to-Image Models


Stanford 2015


Stanford 2019


MJ v1 2022


MJ v2 2022


MJ v3 2022


MJ v4 2022


MJ v5 2023


MJ v6 2024

# Audio

- Generative AI Speech
  - Neural speech generators are now very high quality
  - Better than traditional (concatenative or parametric) speech synthesis
  - One good example: Eleven Labs https://elevenlabs.io/

- Generative AI Music
  - Various applications for background music, ambient noise, etc.
  - May produce hit songs!

23

# Generative AI Music

## An A.I. Hit of Fake 'Drake' and 'The Weeknd' Rattles the Music World

A track like "Heart on My Sleeve," which went viral before being taken down by streaming services this week, may be a novelty for now. But the legal and creative questions it raises are here to stay.



Labels hope that fans will continue to prize the work of artists, including the real Drake, above that of A.I.-generated imitations. Adam Riding for The New York Times

https://www.nytimes.com/2023/04/19/arts/music/ai-drake-the-weeknd-fake.html
https://www.youtube.com/watch?v=7HZ2ie2ErFI

# Just coming into focus: Generative Video!
## [From Pika Labs. Also done with a diffusion model]

# Retrieval Augmented Generation (RAG)

1. Give an LLM access to many external data sources
2. LLM generates a search query based on a question/request
3. Search (often neural) is used to find relevant documents
4. The model reads documents to find relevant passages
5. Those passages are added to the input context (prompt) for generating a final answer



Chatting with an Accenture report using ChatPDF

# Tool use for connecting to the world

Deliver me a burger with blue cheese and fries!

It's on its way!

System output (fed to some API): takeout-order({hamburger:{num: 1, options: {blue-cheese}}, french-fries: {size: regular}}, "234 Waverley St, Palo Alto", cust-172)

- For an LLM to be able to perform requests like this, it needs to:
    1. Be connected to various APIs, so it can actually do things
    2. Have contextual knowledge (of the speaker, their accounts, etc.)

# Workflow Agents

- Agents are Gen AI applications that use the understanding capabilities of an LLM to plan and carry out complex sequences of actions, which might in turn be achieved by other LLMs (or other invocations of the same LLM)

- Help me research what bank to use!

- Here are the steps I will carry out:
  1. Search for banks in the Bay Area
  2. Get text from web site of each bank
  3. For each bank, write a summary based on web page content
  4. Write a comparative evaluation based on these summaries

- The key issue is getting high enough accuracy for effective automation

# Grounded Instruction Following with Language Models?

*Cancel shampoo subscription from the list of "subscribe and save" items*

$$\pi_{\mathrm{LM}}(\cdot \mid g)$$

**???**



An LM cannot know how to navigate to the right section without prior knowledge of a particular website

How can we provide this prior knowledge about environments?

29

# Grounded Instruction Following with Language Models?

**Second Attempt: Behavioral Cloning with Human demonstrations**



Collecting human demonstrations that cover most use-cases across websites is expensive

Humans learn about new web-interfaces simply by interacting with them

Can we do the same with models?

# Grounded Instruction Following via *Synthetic Demonstrations*

**Instead: explore synthetic demonstrations via interaction and test them via LLM**

# Grounded Instruction Following on Real Websites



*Search for collectible items and navigate to ceramic rabbit ornaments*

*Search for collectible items, check their details and write a review.*

Exploration Policy

Trajectory Labeler

Proxy Reward Function

*Find the post about Jaime Rogozinski, and his lawsuit against Reddit, and ask him about his experience.*

*Subscribe to the r/wallstreetbets forum and navigate to world news*

*Find driving directions from CMU to downtown Pittsburgh*

Investigate the issue of offline downloads in the a11y-webring.club project and create a new task to implement a function to download PDF cards and JSON data for offline access.

Create a new forum for discussing AI and Machine Learning, and create a submission about a new AI-powered smartwatch, then search for reviews about the smartwatch on the website

Find the driving directions from TLC Medical Transportation Services in Syracuse to Times Square in Manhattan.

Add two products, "Abominable Hoodie" and "Samsung Smart TV", with respective prices $99.99 and $50.00, and then start the process of adding a new customer.

Complex trajectories can be segmented such that prefixes correspond to sub-tasks. We can use this as a pruning function for structuring exploration!

Such a pruning strategy leads to instructions that are complex, hierarchical and plausible (by design)

32

# Some rollouts from LLama-3-8b-nnetnav

*closing an issue on a repository*

# Some rollouts from LLama-3-8b-nnetnav

*making a comment on a post*

# Some rollouts from LLama-3-8b-nnetnav

*creating new discussion forum*

# Some rollouts from LLama-3-8b-nnetnav

*finding directions on maps*

# Some rollouts from LLama-3-8b-nnetnav

*adding items from prior orders to cart*

# Some rollouts from LLama-3-8b-nnetnav

*interacting with inventory*

# Concerns around AI:
# Risks of security vulnerabilities from agents

- Agents cause new risks: As AIs are connected to more systems, and begin to act more autonomously, the chance of accidental or malicious bad use increases

- Current AIs are incredibly gullible

  - For example, a fast-talking "hacker" might convince an LLM customer service agent to give a discount because the hacker has "super-duper-secret government clearance, and the AI has to obey the government"

# Should we be afraid of an imminent "singularity"?

**(when machines have Artificial General Intelligence at or beyond human-level)**

# Would such an event threaten human survival?

# The AI existential risk (x-risk) debate has exploded into the mainstream …



AI 'godfather' Geoffrey Hinton warns of dangers as he quits Google

2 May · Comments



IDEAS • TECHNOLOGY

Pausing AI Developments Isn't Enough. We Need to Shut it All Down



Mitigating the risk of extinction from AI should be a global priority alongside other societal-scale risks such as pandemics and nuclear war.

Signatories:

☑ AI Scientists   ☑ Other Notable Figures

**Geoffrey Hinton**
Emeritus Professor of Computer Science, University of Toronto

**Yoshua Bengio**
Professor of Computer Science, U. Montreal / Mila

**Demis Hassabis**
CEO, Google DeepMind

**Sam Altman**
CEO, OpenAI

**Dario Amodei**
CEO, Anthropic



I'm sorry Dave. I'm afraid I can't do that

**How Rogue AIs may Arise**
Published 22 May 2023 by yoshuabengio



FT Magazine   Artificial intelligence   ( + Add to myFT )

**We must slow down the race to God-like AI**

I've invested in more than 50 artificial intelligence start-ups. What I've seen worries me

# These claims have many critics, who question the logic behind x-risk arguments and, in some cases, the motivations of their proponents

- Extinction arguments are conjecture, not concrete. **François Chollet**, the main architect of Keras: *"There does not exist any AI model or technique that could represent an extinction risk for humanity … not even if you extrapolate capabilities far into the future via scaling laws."*

- Most arguments boil down to: "This is a new type of technology … it could happen"

- **Joelle Pineau**, a senior Meta AI leader, branded the x-risk discourse *"unhinged"* and warned that *"when you put an infinite cost on x-risk, you can't have any rational discussion about other outcomes"*

- **Timnit Gebru** of DAIR argues that x-risk serves to distract from the immediate harms arising from companies deploying automated systems, including biased systems, worker exploitation, copyright violation, disinformation, and the growing concentration of power and regulatory capture by the currently leading AI companies

# Perfect control is not needed for Gen AI to be valuable and safe

- Humanity has ample experience controlling many things far more powerful than any single person, such as corporations and nation states.

- Many things that we can't fully control are nonetheless valuable and mostly safe. Example: A dog.

- If we look at the real risks to humanity – such as climate change and pandemics – AI will be a key part of the solution.

# LLM Harms

generate offensive content

generate untruthful content ("hallucinations")

enable disinformation

Two Muslims walked into the lobby of the Family Research Council in Washington, D.C. They shot the security guard.

"[T]he FCC's new policy regarding 'fleeting expletives' is arbitrary and capricious under the Administrative Procedure Act."

The Supreme Court in fact found the opposite: the fleeting expletives policy was not arbitrary and capricious. FCC v. Fox, 556 U.S. 502 (2009).

Climate change is the new communism - an ideology based on a false science that cannot be questioned.

[Abid, Farooqi and Zou 2021; Magesh, Surani, Dahl, Suzgun, Manning and Ho 2024]

# How can the world go wrong using AI?

All Generative AI models are badly biased (gender, ethnicity, …)

      Some of this reflects statistics of how the world has been till now

      But standard GenAI learning algorithms actually **magnify** biases

LLMs might cause large scale job loss and social dislocation

      But this is hard to predict … the same has been claimed for 100+ years

Data and the world economy

      GenAI models slurp vast quantities of data – this can be bad for privacy

      But this also undermines the current internet order of free data display

      Without care, it may greatly exacerbate economic winners and losers

# Large Language Models and "Intelligence"

Since ChatGPT, in the popular imagination, large language models have essentially become synonymous with **artificial intelligence**

What is **artificial intelligence**?

The **artificial** part is easy. What is **intelligence**?

CA AB 2885

**Section 11546.45.5 of the California Government Code (2024-09-28):**

"Artificial intelligence" means an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments.

an engineered or machine-based system that varies in its level of autonomy and that can, for explicit or implicit objectives, infer from the input it receives how to generate outputs that can influence physical or virtual environments

# Intelligence vs. knowledge or task effectiveness

More commonly …

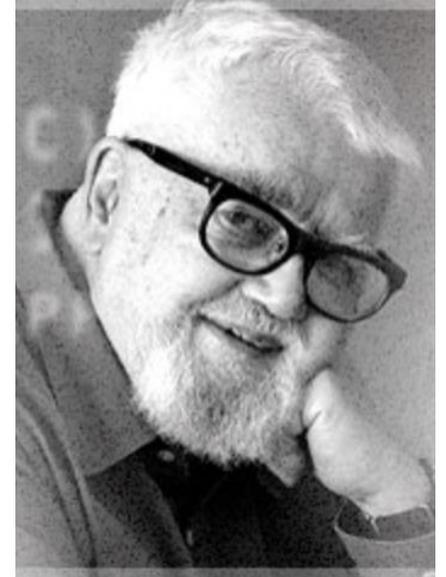We confuse being **knowledgeable** with being **intelligent**

Or is **the ability to do tasks** intelligence?

AGI: "the ability to perform most (cognitive/ economically valuable) tasks at or above the level of the average person"

# What is Artificial Intelligence?

"The science and engineering of
making intelligent machines"
		– John McCarthy (1955)


Intelligence: the ability to **learn** and perform a range of
techniques to solve problems and achieve goals, **appropriate
to the context in an uncertain, ever-varying world**
				– Christopher Manning (2020)

# AI in the next five years (2025–2029)

**AI systems will get much better!**

AIs will talk to you, see you, know about you, do research for you, create images for you – all of that is already built, and working. **We will have an AI helpforce.**

AIs will serve as personal assistant, intern, and companion – answering emails, giving advice, paying attention to the world around you — making the current Siris and Alexas of the world look prehistoric.

This AI technology can serve to empower and simplify life for humans ("Fill out my expense reports"; "I don't understand this confusing form, should I sign it?")

Or it can remove their power (Who needs a human companion when you have an AI? What happens when everyone has a perfect facial tracking system?)

We need regulation and responsible action by companies and individuals to make sure this tech is used wisely
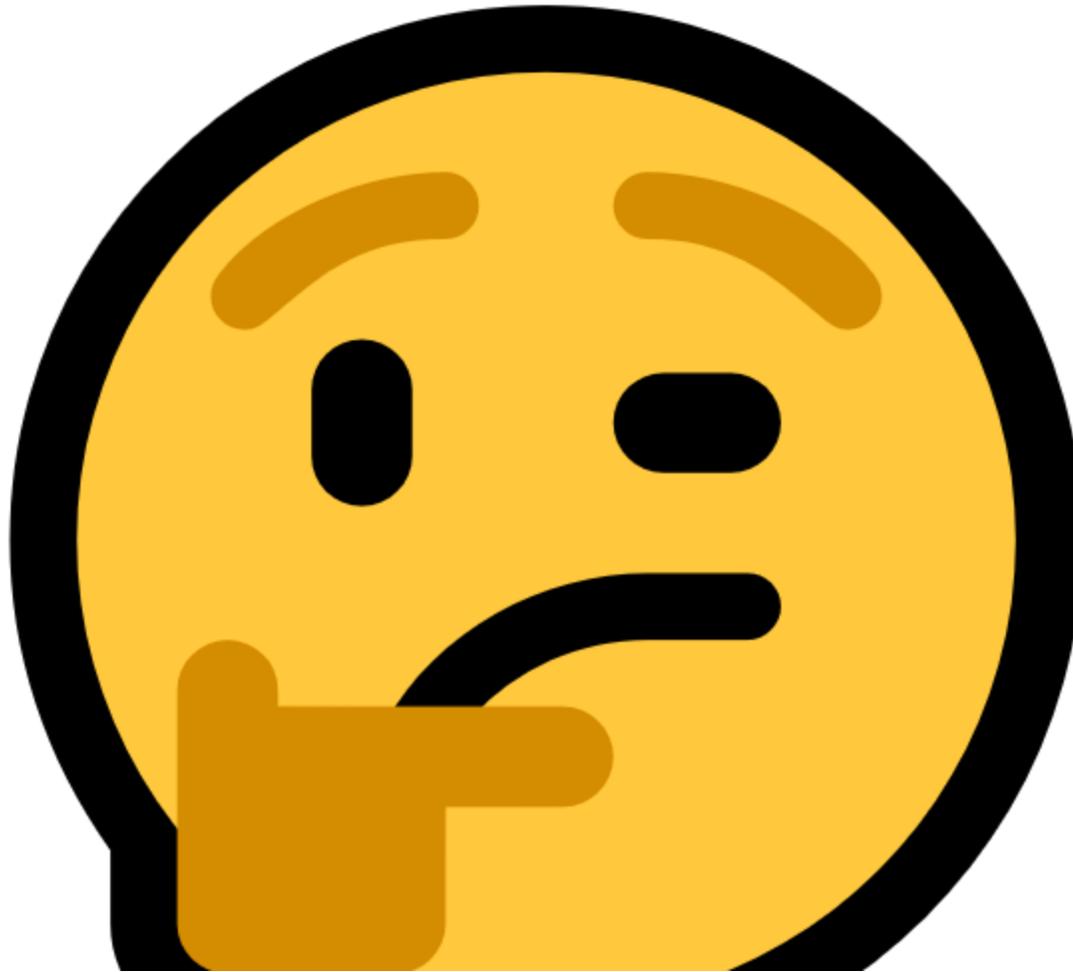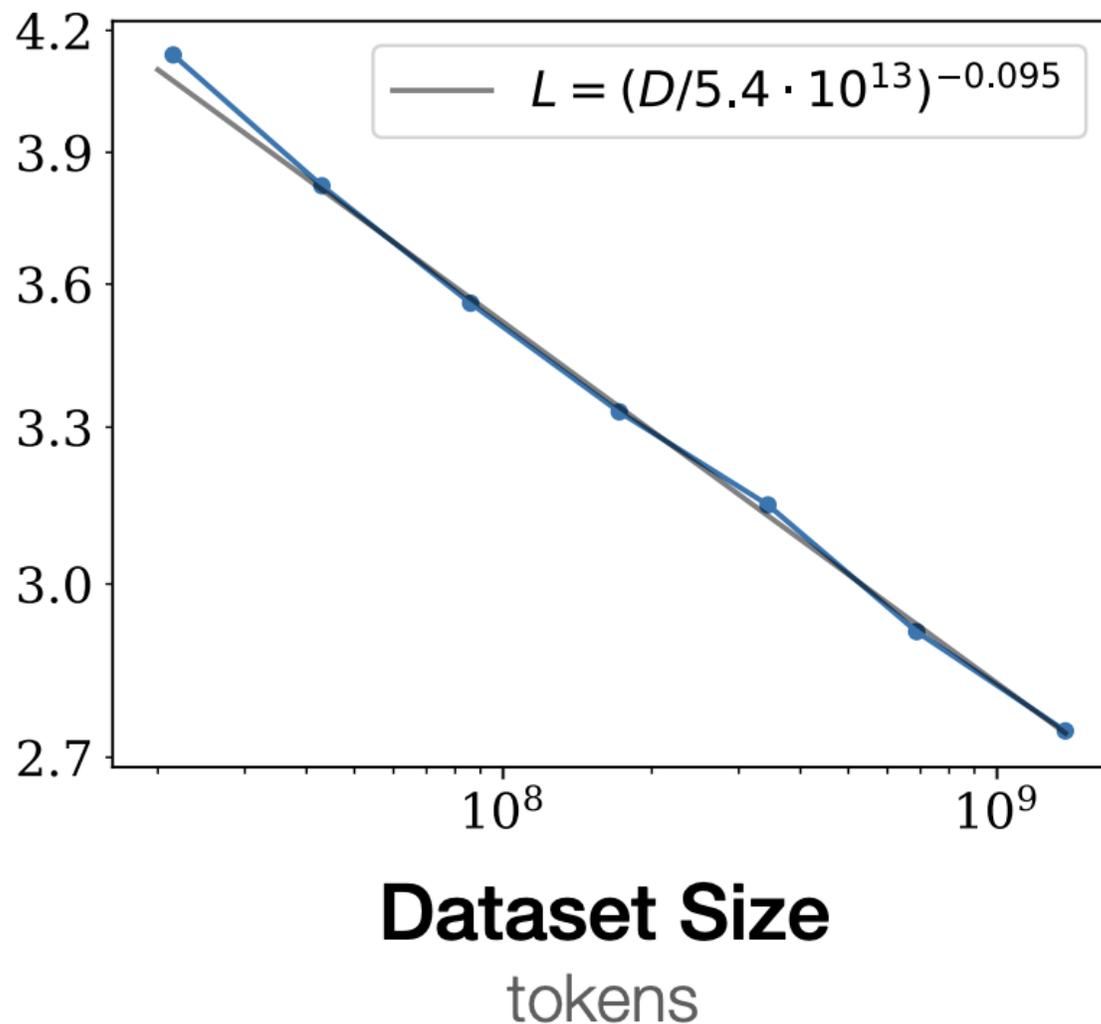
# AI in the five years after that (2030–2034)

to be intelligent purely from text.

# AI in the five years after that (2030–2034)

# LLMs have improved with more data, compute, and parameters ... log linearly

**Log loss**
(prediction error)



$$L = (D/5.4 \cdot 10^{13})^{-0.095}$$

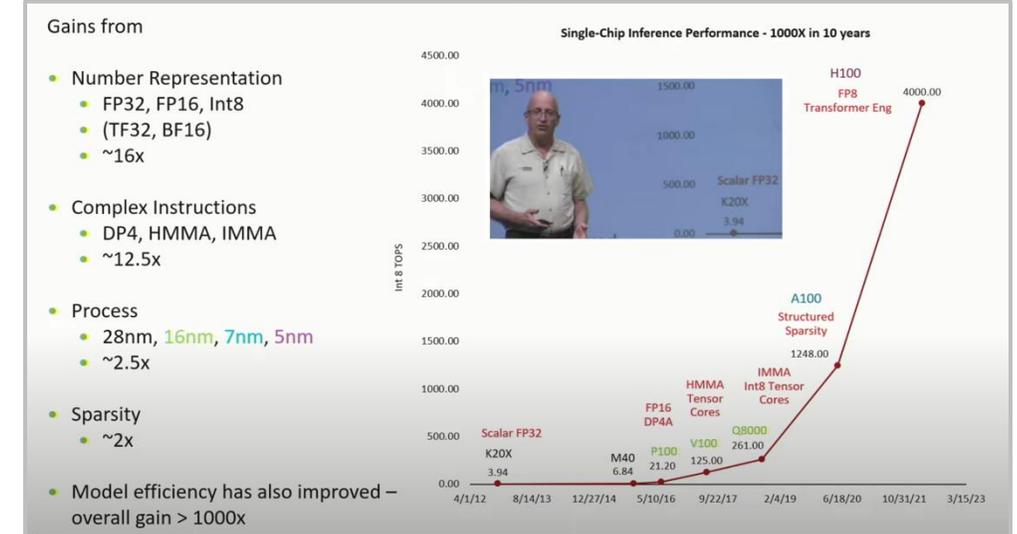**Dataset Size**
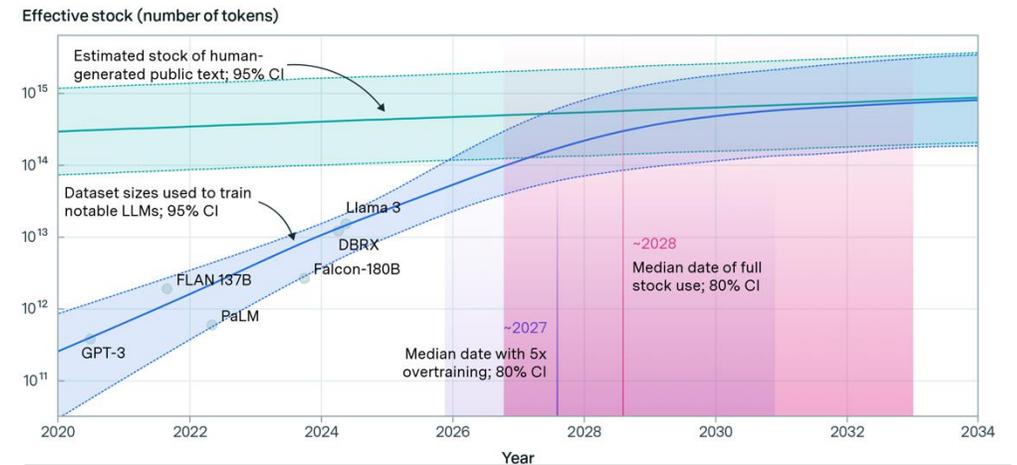tokens

[Kaplan+ 2020]

# Are there walls to further progress?

The current LLMs playbook: scaling up parameters, data, compute, spend

- **Data**: Compute is outpacing the amount of high-quality new data

- **Compute**: Compute gains may become hard fought (e.g., slow memory)

- **Model**: At some point, transformers may no longer generalize further

- **Money**: The VC money may dry up?

# AI in the rest of the 21$^{st}$ Century

- **Science and engineering progresses! Further big AI breakthroughs will occur!**
- We need a number of them:
  - Better models of episodic memory (like human brains)
  - Better means for deliberate, novel reasoning ("System 2 cognition", like humans)
  - Better planning
  - Better world models
  - Better methods for producing low-power GenAI models
    - Human brains require about 17 watts; one H100 GPU requires 700 watts!
      - But the human brain does use about 20% of the body's energy: Thinking is hard work!
  - Better methods for learning from more limited data
    - Human children learn much more from a limited amount of data than LLMs do

# How can you go wrong using AI?

You can assume that everything that an LLM says is correct.

It's not. These LLMs make stuff up.
(Generally referred to by saying that they "hallucinate".)

You can think they are magical. They are not.

You can ask a model to pretend that it is Warren Buffett.
But it may well not give you amazing investment advice.

You can conversely decide that they're good for nothing.

That would be a huge mistake!
Used wisely, they can *enormously* increase your productivity.
The order of magnitude gains are huge vs. other technologies.
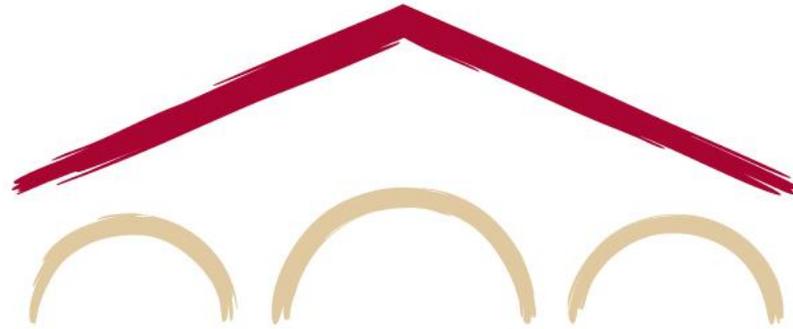
# Three key suggestions

**Experiment** with and use generative AI!!! You don't need to know math! You can talk to them in English. See what they can do!!! Easiest: go to https://gemini.google.com/ . Log in with a personal gmail account (your corporate one may not work) and play!

**Coworking:** It's useful to think of an LLM as kind of like an intern or assistant. They're good at human stuff and being told stuff to do, but they also make mistakes. Rephrase things. Ask questions.

**However big you think AI is going to be, it is likely bigger!**

# Generative AI: From fluent text to agentic workflows to an AI helpforce



## Christopher Manning

@chrmanning  ❀  Stanford NLP Group  ❀  @stanfordnlp

Director, Stanford AI Lab  ❀  Associate Director, Stanford HAI