# AI Policy

Miles Brundage

# AGENDA

## 01 BACKGROUND

- Policy in general
- My road to AI and AI policy
- Some things I worked on at OpenAI

## 02 AI POLICY IN GENERAL

- Key concepts
- Key tensions

## 03 SOME HOT TAKES

- [ redacted – you'll have to wait until later ]

## 04 WHERE YOU FIT IN

- Policy-related research/engineering opportunities
- Your voice in companies and public discussions

# 01

# Background

# "Policy"

"X policy" basically just means "the decisions that society makes about X, and how they are and should be made."

AI policy/governance is the theory and practice of governmental and non-governmental decision-making about AI.

# "Policy"

Healthcare policy: how insurance is regulated, how drug approvals work, etc.

Energy policy: how utilities are regulated, how research and development is encouraged, etc.

We're still figuring out exactly what AI policy should involve.

# "Policy"

Regulation is a part of it, but not all of it.

E.g. the CHIPS Act (subsidizing the American semiconductor industry) is an example of AI policy.

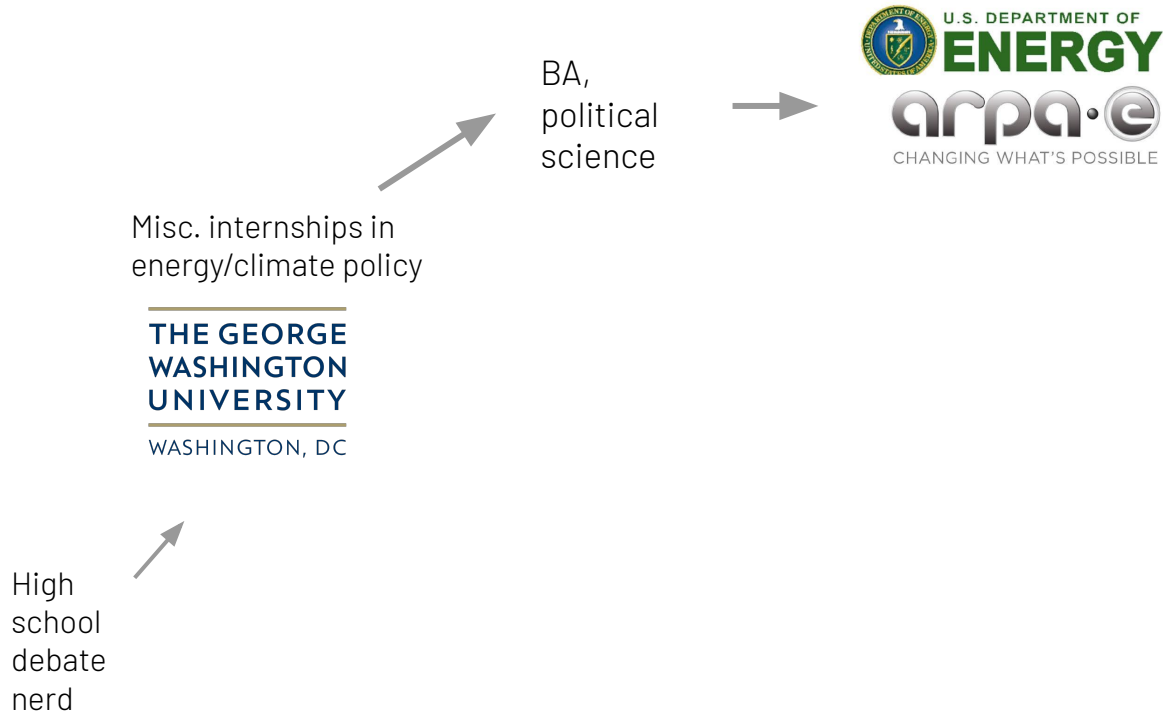It's also not all about the government, either.

# "Policy"

- A lot of AI policy is informed by and related to safety, broadly defined, including alignment, reliability, etc.
- But there's more to AI policy than safety, and I won't say much about safety on its own here
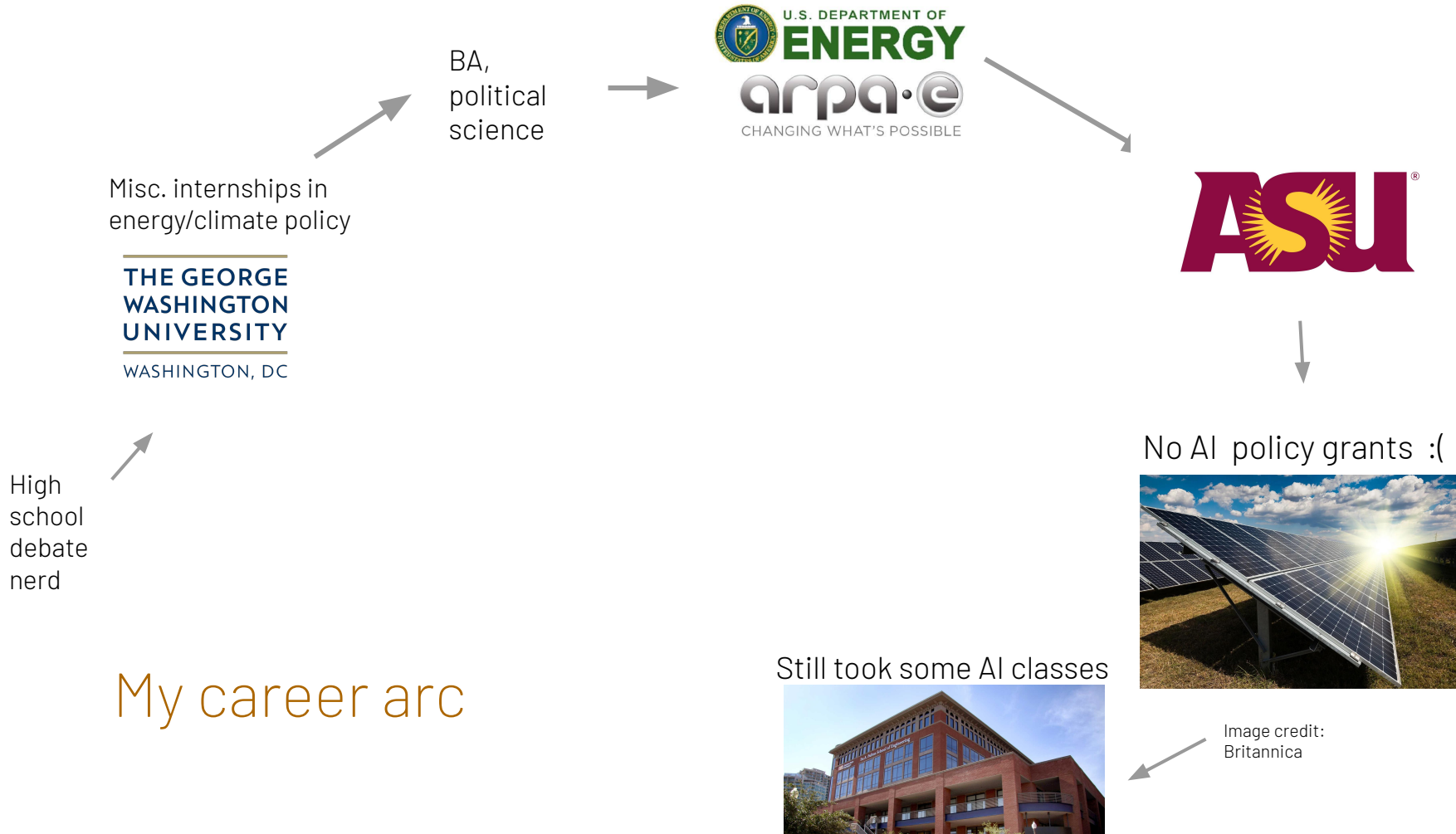
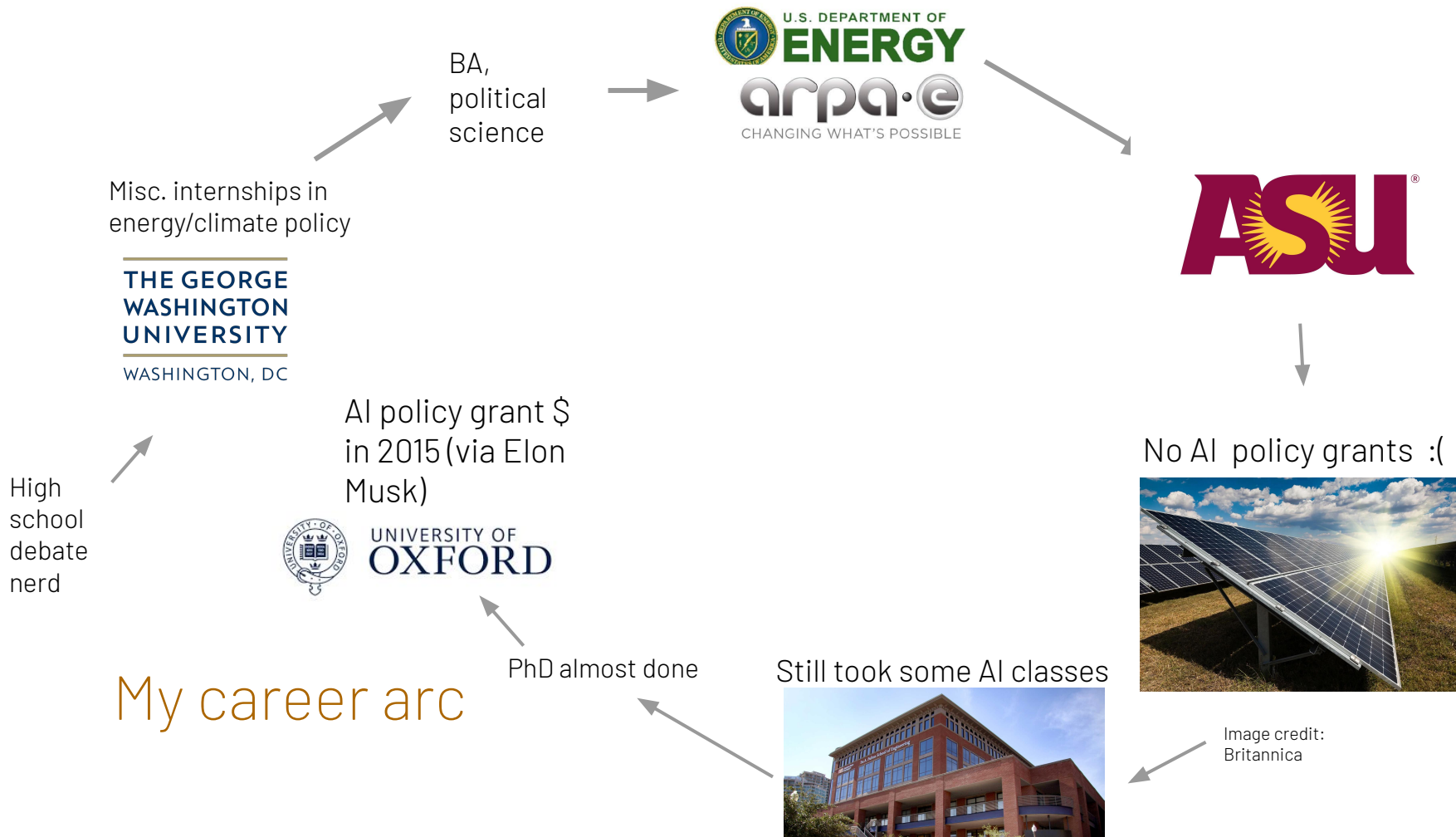# AI policy examples

Illustratively, AI policy involves:
- Company decision-making about the right use case policies for APIs and first party products
- Company decisions around deployment of technologies
- Regulations like the EU AI Act
- Industry norms/best practices
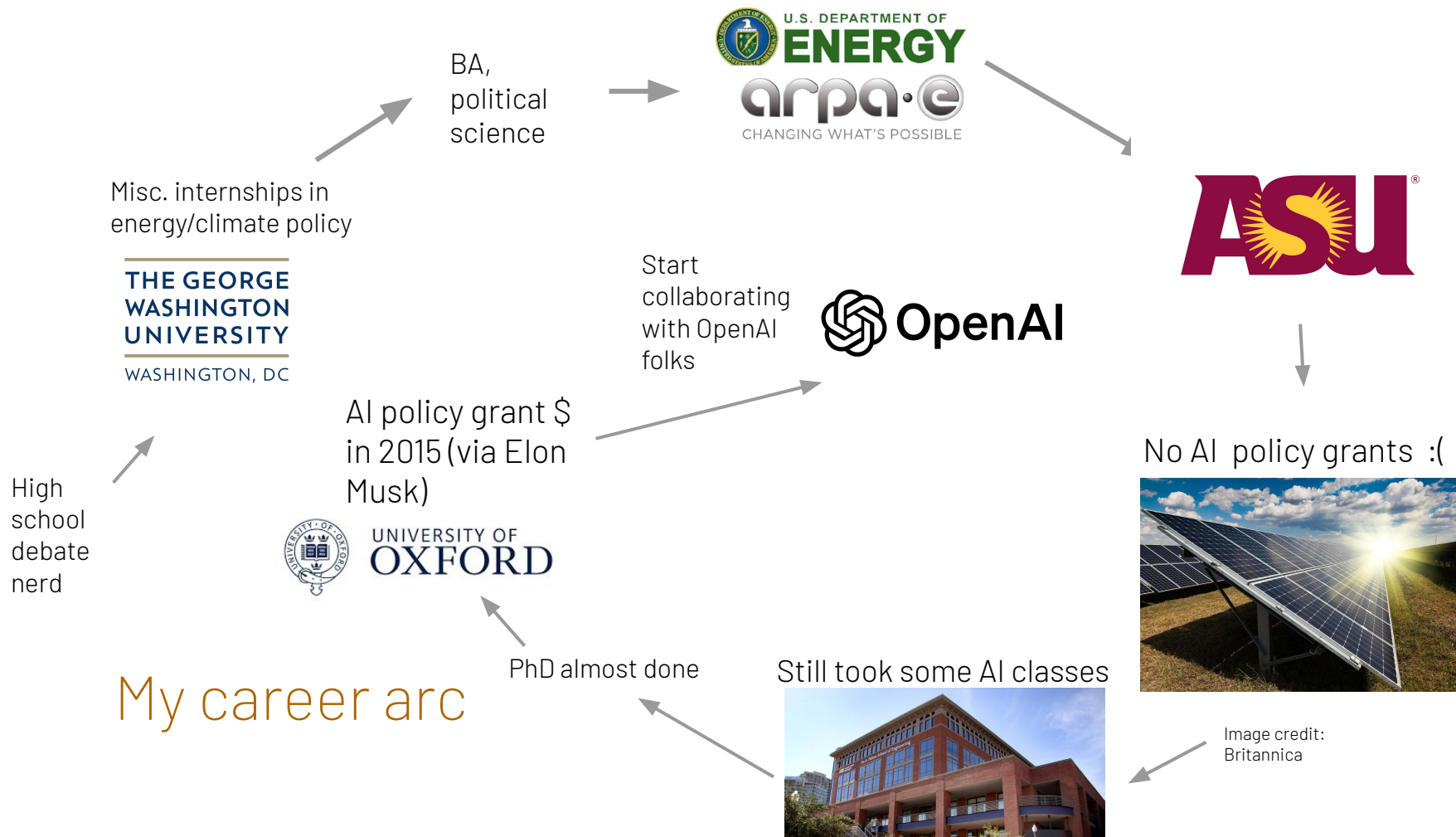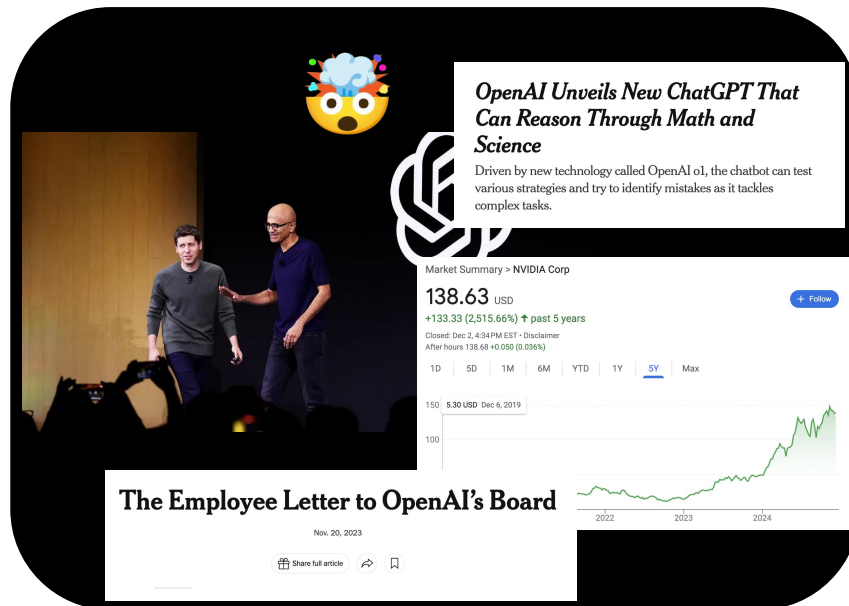- Academics' decisions about open sourcing and documenting models
- Etc.

# My career arc







IEEE Spectrum, 2008

BA,
political
science

Misc. internships in
energy/climate policy

THE GEORGE
WASHINGTON
UNIVERSITY
WASHINGTON, DC

High
school
debate
nerd

# My career arc

My career arc

High school debate nerd

Misc. internships in energy/climate policy

THE GEORGE WASHINGTON UNIVERSITY
WASHINGTON, DC

BA, political science

U.S. DEPARTMENT OF ENERGY
arpa·e
CHANGING WHAT'S POSSIBLE

ASU

No AI policy grants :(

Still took some AI classes

Image credit: Britannica

My career arc

High school debate nerd

Misc. internships in energy/climate policy

THE GEORGE WASHINGTON UNIVERSITY
WASHINGTON, DC

BA, political science

U.S. DEPARTMENT OF ENERGY
arpa·e
CHANGING WHAT'S POSSIBLE

ASU

No AI policy grants :(

Image credit: Britannica

Still took some AI classes

PhD almost done

AI policy grant $ in 2015 (via Elon Musk)

UNIVERSITY OF OXFORD

High school debate nerd

Misc. internships in energy/climate policy

THE GEORGE WASHINGTON UNIVERSITY
WASHINGTON, DC

BA, political science

U.S. DEPARTMENT OF ENERGY
arpa·e
CHANGING WHAT'S POSSIBLE

ASU

AI policy grant $ in 2015 (via Elon Musk)

UNIVERSITY OF OXFORD

Start collaborating with OpenAI folks

OpenAI

No AI policy grants :(

My career arc

PhD almost done

Still took some AI classes

Image credit: Britannica

# My time at OpenAI in short



Not really - GPT-2→3, Dota, etc. were quite exciting

# A few things I worked on

Analysis of and policies for, e.g.:

- GPT-2, 3, 3.5, **4**, 4v, 4o
- **Codex**
- DALL-E 2, 3
- CLIP
- **o1**

Helping shape/scale red teaming, economic impact analysis, etc.

More general research on, e.g., agents, compute, frontier AI regulation, etc.

# 02

## AI Policy in General

# Key concepts in AI policy

I'll try to be mostly* uncontroversial in this section, compared to the next one

*there is basically nothing that's totally uncontroversial in AI policy

# AI as a general-purpose technology





Images: IEA, Britannica

# AI as a general-purpose technology

→ will have impacts across all sectors

→ "AI policy" will become, or interact with, "everything policy" by default (how to avoid overreach?)

→ can be differentially sped up/slow down in some aspects, though such interventions will degrade by default

# AI is a fast-moving technology



Exam results (ordered by GPT-3.5 performance)

Estimated percentile lower bound (among test takers)

OpenAI

# AI is a fast-moving technology



AI performance on a set of Ph.D.-level science questions

# AI is a fast-moving technology

It also features semi-regular paradigm shifts:

- RL in videogames/simulations →

- Large-scale unsupervised pretraining + a bit of SL supervised and reinforcement learning at the end →

- "Actual RL" on language models

# AI development is (in part) a collective action problem



Prisoners' dilemma

© 2010 Encyclopædia Britannica, Inc.

Not nec. this exact dilemma

[Submitted on 10 Jul 2019]

**The Role of Cooperation in Responsible AI Development**

Amanda Askell, Miles Brundage, Gillian Hadfield

In this paper, we argue that competitive pressures could incentivize AI companies to underinvest in ensuring their systems are safe, secure, and have a positive social impact. Ensuring that AI systems are developed responsibly may therefore require preventing and solving collective action problems between companies. We note that there are several key factors that improve the prospects for cooperation in collective action problems. We use this to identify strategies to improve the prospects for industry cooperation on the responsible development of AI.

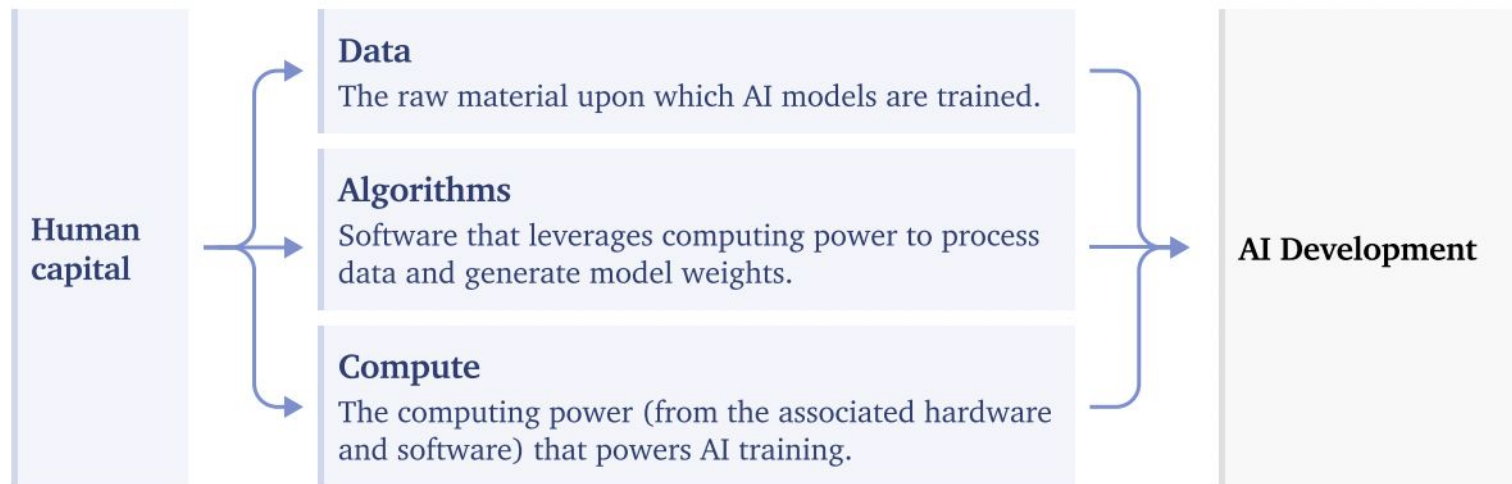**Frontier Model Forum: Advancing frontier AI safety**

The Frontier Model Forum draws on the technical and operational expertise of its member companies to benefit the entire AI ecosystem, advancing AI safety research and supporting efforts to develop AI applications to meet society's most-pressing needs.

SEPTEMBER 12, 2023

FACT SHEET: Biden-Harris Administration Secures Voluntary Commitments from Eight Additional Artificial Intelligence Companies to Manage the Risks Posed by AI

BRIEFING ROOM    >    STATEMENTS AND RELEASES

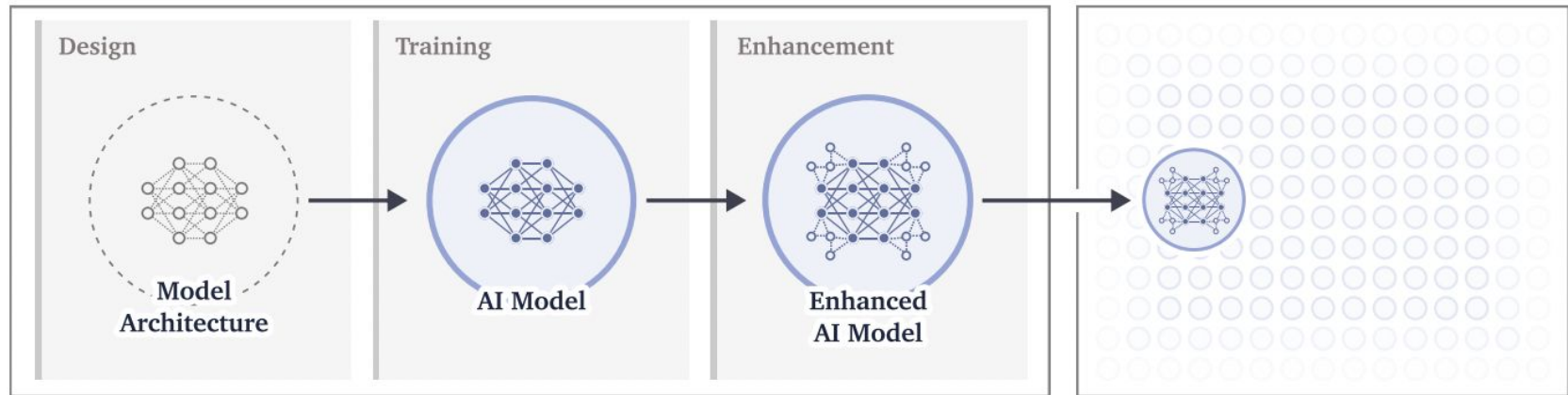# Different inputs to AI capabilities are easier/harder to govern than others



Sastry et al., 2024

# Different phases of AI development and deployment raise different policy questions

## Simplified AI Lifecycle

**Development**

Design — Model Architecture

Training — AI Model

Enhancement — Enhanced AI Model

**Deployment**

Copyright?  Bias?  Misuse?  Economic impact?

Sastry et al., 2024

# Different phases of AI development and deployment raise different policy questions

**Development & Deployment Lifecycle**

**Initial Development**
Problem identification, goal setting
Initial impact assessment
Data sourcing, curation, filtration

**Alignment**
Instruction generation
Fine-tuning
Alignment evaluations

**Evaluation & Iterative Development**
Model evaluations
Revised impact assessment, hazard analysis
Red teaming, user testing

**Downstream Assessment**
Retrospective reviews
Retrospective impact assessment
Platform-level risk measurement

**Deployment & Ongoing Evaluation**
Private betas
Use case pilots
Misuse detection & response

OpenAI

# AI is still behind many other areas of policy, analytically

Many AI forecasts/opinions etc. are way off, or unfalsifiable/"not even wrong"

Many seemingly big topics are basically brand new (e.g. test-time compute)



Jacob Steinhardt, IPCC

# It's not all about the (base) model

Platforms

Overview    Documentation    API reference

Use cases

Systems

ChatGPT PLUS                    DALL·E

casetext

duolingo

Models

Whisper                         GPT-4

# We're transitioning from self-regulation to "real" regulation

OCTOBER 30, 2023

FACT SHEET: President Biden Issues Executive Order on Safe, Secure, and Trustworthy Artificial Intelligence

BRIEFING ROOM ▸ STATEMENTS AND RELEASES

## Artificial Intelligence Act: MEPs adopt landmark law

Press Releases PLENARY SESSION IMCO LIBE 13-03-2024 - 12:25

# But there will be gaps for the foreseeable future, where the behavior of people within companies makes a difference

(Will return to this at the end)

# Tensions between different objectives

*Spoiler: I don't think we're on, or even near the Pareto frontier for all of these!*

*Often the tradeoffs are overstated, and a great thing about AI is that it can help "grow the pie" through automated labor.*

*The point is just that there are sometimes actually tradeoffs.*

# Tension: preventing risks vs. unlocking benefits

Where the rubber meets the road:
- The general pace of progress
- Use case policies / fine-tuning approaches
- General permissiveness of regulation
- Open source

# Tension: competing with other companies/countries economically and militarily while collaborating on shared safety challenges



## US hits China's chip industry with new export controls

Parting measures by Biden administration aim to slow Beijing's development of AI with military applications

The new measures will hit chip manufacturers including Semiconductor Manufacturing International Corporation and Chinese companies that produce chipmaking tools © Reuters

**Demetri Sevastopulo** in Washington  2 HOURS AGO

70



## Biden, Xi agree that humans, not AI, should control nuclear arms

By **Jarrett Renshaw** and **Trevor Hunnicutt**

November 16, 2024 5:04 PM PST · Updated 16 days ago

# Tension: preventing risks vs. concentrating power

- Open source has a lot of potential for misuse…
  - …but also decentralizes "deep" access to AI capabilities
- Fine-tuning of closed models can reduce direct misuse risk…
  - …but also represents an imposition of a certain set of values

# Tension: addressing existing vs. anticipated issues

- There are synergies, but also tradeoffs in policymaker attention, compute, researcher/engineer time, etc.
- E.g. bias/hallucination/"mundane" misuse vs. catastrophic misuse and accidents

# 03

## Some ~~Hot Spicy Personal~~ Hot Takes

# 03

Some Hot Takes

# Ranking methodology

I'm giving a talk on AI policy. Please rank the following "takes" from the talk in descending order of "hotness"/"spiciness":

# Ranking methodology

Saving one of them, which may or may not be the spiciest, for the section on what you can do

# You should watch/read more sci-fi

Battlestar Galactica
Person of Interest
Pantheon
Travelers
Transcendence*
WALL-E
Terminator 2

Westworld
Interstellar
The Diamond Age
Robopocalypse
The Player of Games
Altered Carbon

# You should watch/read more sci-fi

...but not too much, for the reason I'll give at the end...

# The economic impacts will be huge soon

Already, some gig workers are being negatively affected (e.g. copywriters, illustrators)

Others are gaining a lot of productivity, with unclear job consequences

Soon more interactive roles will be impacted, especially ones that are already outsourced/done remotely

# Huge != no one has jobs, but…

Double digit annual
economic growth within
five years

Hundreds of millions
displaced from their
previous jobs

**INNOVATIONS**

## ChatGPT took their jobs. Now they walk dogs and fix air conditioners.

Technology used to automate dirty and repetitive jobs. Now, artificial intelligence chatbots are coming after high-paid ones.

By Pranshu Verma and Gerrit De Vynck
June 2, 2023 at 6:00 a.m. EDT

Washington Post

# We need to talk about post-work futures now

# Investing in society's resilience to AI's impacts is ~a free lunch

Not literally free, but –

If all countries in the world had strong social safety nets (to cushion people from harms associated with unemployment), robust investment in cybersecurity, ubiquitous physical defenses against biological threats (e.g. far-UVC)...

We could distribute AI capabilities more widely/evenly and have more transparency – right now, these are sometimes in tension with safety/security

# Access to AI is unusually equal by the standards of previous technologies

- Piggybacking on the internet → fast distribution
  - Though the digital divide is still a huge issue
- For part of this year, the best free model and the best paid model were the same model (other than rate limits)
- The point here is not to excuse cases where companies fall short! Just to calibrate.

# Access to AI is unusually equal by the standards of previous technologies



**GPT-4 3-shot accuracy on MMLU across languages**

| Language | Accuracy |
| --- | --- |
| Random guessing | 25.0% |
| Chinchilla-English | 67.0% |
| PaLM-English | 69.3% |
| GPT-3.5-English | 70.1% |
| GPT-4 English | 85.5% |
| Italian | 84.1% |
| Afrikaans | 84.1% |
| Spanish | 84.0% |
| German | 83.7% |
| French | 83.6% |
| Indonesian | 83.1% |
| Russian | 82.7% |
| Polish | 82.1% |
| Ukranian | 81.9% |
| Greek | 81.4% |
| Latvian | 80.9% |
| Mandarin | 80.1% |
| Arabic | 80.0% |
| Turkish | 80.0% |
| Japanese | 79.9% |
| Swahili | 78.5% |
| Welsh | 77.5% |
| Korean | 77.0% |
| Icelandic | 76.5% |
| Bengali | 73.2% |
| Urdu | 72.6% |
| Nepali | 72.2% |
| Thai | 71.8% |
| Punjabi | 71.4% |
| Marathi | 66.7% |
| Telugu | 62.0% |

Legend: Random, Chinchilla, PaLM, gpt-3.5, gpt-4

Accuracy →

# ...but in the worst case, it could concentrate power as never before...

- More compute → better answers → by default, you should expect AI services to bifurcate dramatically between free and paid
- AI can also automate surveillance, censorship, political messaging, etc. at scale with precision

# ...but it could still be a massive force for equality

...because the world of today is a very weak baseline on this score...

...and because there may be diminishing returns on cognition in many aspects of life.

# ...but it could still be a massive force for equality



Life expectancy, 2021
The period life expectancy at birth, in a given year.

Data source: UN WPP (2022); HMD (2023); Zijdeman et al. (2015); Riley (2005) – Learn more about this data
OurWorldInData.org/life-expectancy | CC BY

# ...but it could still be a massive force for equality



Our World in Data

# ...but it could still be a massive force for equality



Test score in mathematics

Poorest households in Brazil:
Annual household income: $630
Test scores of the children: 375

Richest households in Brazil:
Annual household income: $11,630
Test scores of the children: 477

Brazil

Annual household income (adjusted for price differences between countries) on a logarithmic axis

Our World in Data

# ...but it could still be a massive force for equality

Consider perhaps the most basic measure of a functioning school: that there are teachers in the school teaching classes. On any given day, nearly a quarter of teachers in rural India simply do not show up. And when they do turn up, they're often not teaching. A World Bank report found that even when Kenyan teachers were present, they were absent from their classrooms 42% of the time.

Low quantity of education

Low quality of education

Even if we ignore these constraints, developing-country schools struggle with ineffective curricula and overly prescriptive pedagogy. National curriculums rarely meet students where they are, and few students are at "grade level," but teachers are still instructed to teach as if they are. Instruction consists largely of memorization. Rather than foster critical thinking, teachers effectively train students' ability to repeat back what the teacher wants to hear. And perhaps worst of all, students are often taught in a language they don't even speak.

Gilbert, 2024

# …but it could still be a massive force for equality

## Lifespan

Lifespan uses GPT-4 to radically improve health literacy and patient outcomes.



With over half of Americans reading at or below the 6th grade level, Dr. Ali and Dr. Mirza proposed using GPT-4 to simplify surgical consent forms from a college reading level to a middle school reading level. To mitigate the risk of bias and hallucination, Lifespan leadership created a system where GPT-4 would do a first pass, and then legal and medical reviewers would check the output.

# There should be much more serious consideration of a "CERN for AI" scenario



"The CERN approach" ~ = pooling resources to build and operate centralized infrastructure in a transparent way, as a global scientific community, for civilian rather than military purposes

# There should be much more serious consideration of a "CERN for AI" scenario

 x 20     vs.     1x 

# There should be much more serious consideration of a "CERN for AI" scenario



I feel thin, sort of stretched, like butter scraped over too much bread.

Lord of the Rings: Fellowship of the Ring

**The bread and butter problem in AI policy**

There is too little safety and security "butter" spread over too much AI development/deployment "bread."

**MILES BRUNDAGE**
NOV 05, 2024

# There should be much more serious consideration of a "CERN for AI" scenario

- Ensure that the very most capable models are developed extremely securely and safely:
  - Combine the world's talent on security, then safety, then capabilities – in that order, otherwise you just speed up development and everyone steals it and fine-tunes it in dangerous ways
- When models are derisked, distribute and deploy them widely

# There should be much more serious consideration of a "CERN for AI" scenario

- It's not obvious that this is the right thing to do but it deserves serious debate and being fleshed out
- Key question: how could this be designed such that there is *distributed* control over these *centralized* capabilities (e.g. multiple parties can stop a dangerous training/inference run).
  - This is partly a technical and partly a political question.

# AI sentience will also be a huge issue

The costs of error in either direction are huge

ROBERT LONG, JEFF SEBO · OCTOBER 30, 2024

## New report: Taking AI Welfare Seriously

Our new report argues that there is a realistic possibility of consciousness and/or robust agency — and thus moral significance — in near-future AI systems, and makes recommendations for AI companies. (Joint output with the NYU Center for Mind, Ethics, and Policy.)

# 04

## Where You Fit In

Even if you don't care about AI policy, AI policy cares about you.

Public (government) and private (corporate/non-profit etc.) decisions will affect your career in various ways, as well as your life as a citizen more generally.

See also my blog post: "FAQs and General Advice on AI Policy Careers"

# We're running out of time

## Scoring Humanity's Progress on AI Governance

Miles Brundage
8 min read · May 28, 2023

| Category | 2022 | 2023 | 2024 |
|----------|------|------|------|
| Shared Understanding of the Challenge | D+ | B- | B |
| Technical Tooling | D+ | C- | C |
| Regulatory Infrastructure | D+ | C+ | C |
| Legitimacy | D- | D+ | C- |
| Societal Resilience | F | F | D- |
| Differential Technological Development | F | D+ | D |

# We're running out of time

| Category | 2022 | 2023 | 2024 |
|---|---|---|---|
| Shared Understanding of the Challenge | D+ | B- | B |
| Technical Tooling | D+ | C- | C |
| Regulatory Infrastructure | D+ | C+ | C |
| Legitimacy | D- | D+ | C- |
| Societal Resilience | F | F | D- |
| Differential Technological Development | F | D+ | D |

2025     2026     2027...

?

# We're running out of time

| Category | 2022 | 2023 | 2024 |
|---|---|---|---|
| Shared Understanding of the Challenge | D+ | B- | B |
| Technical Tooling | D+ | C- | C |
| Regulatory Infrastructure | D+ | C+ | C |
| Legitimacy | D- | D+ | C- |
| Societal Resilience | F | F | D- |
| Differential Technological Development | F | D+ | D |

**Your career may not (have to) be as long as you thought**

Not sure now is the best time to start a PhD...

...except, perhaps, if you're OK multitasking

# The AI policy Pareto frontier

# The AI policy Pareto frontier

Maximizing benefits and minimizing risks...

...privacy/copyright protection vs. knowledge about the world...

...raising the ceiling of capabilities/raising the floor of access

# Ways you can help

Pushing out the frontier with technical innovation

Helping organizations and the world make informed decisions about how to get to the frontier and pick a point on it

# Pushing out the frontier with technical innovation

# Pushing out the frontier with technical innovation



**Access**

- Privacy-Preserving Third-Party Access to Datasets
- Preservation of Evaluation Data Integrity

- Provision of Compute Resources

- Facilitation of Third-Party Access to Models

- Access to Downstream User Logs and Data

# Getting to the frontier and moving around on it: your voice matters!

- Companies have, and may always have, some latitude in how they interpret voluntary commitments and regulations.
  - This is good in some respects by fostering innovation, but also can enable recklessness
- Companies will generally want you to be happy, which gives you power.
  - You won't always get your way (company behavior isn't *just* a function of employee views – there are also competitive pressures, etc.), but every bit of thoughtfulness contributes to company culture

# Getting to the frontier and moving around on it: your voice matters!

- Public and private debates matter: it can be helpful to push back when someone is mischaracterizing their opponents, dismissing legitimate concerns, etc.
  - Though don't spend all your time on it – there's a lot of real work to do, as well

Terminator 2: Judgment Day

# Acknowledgments

# Thank you!

@miles_brundage on Twitter
Miles Brundage on Google Scholar
milesbrundage.substack.com