



**"Data, Compute, and Ethics:
Navigating the Intersection
of AI and Humanity"**

Scott H. Hawley
Belmont University

CCCU Multi-Academic 2 Conference, March 1 2024

Outline:

- "What is AI?" ...is a tough question.

I'll offer a few views:

- AI is Data
- AI is Compute
- "AI is People"
...and about those people?

*SPOILER: I'm not going to tell you
"What to do about ChatGPT"*

...although...

*"Watermarking Makes Language Models Radioactive"
(Meta AI, Feb 24, 2024)*



History: The Rise of ML in AI, i.e.

Data+Compute

For most of its history, "AI" (McCarthy et al, 1956) was symbolic & rules-based, vs...

"Machine Learning" ([Arthur Samuel](#), 1959) was separate.

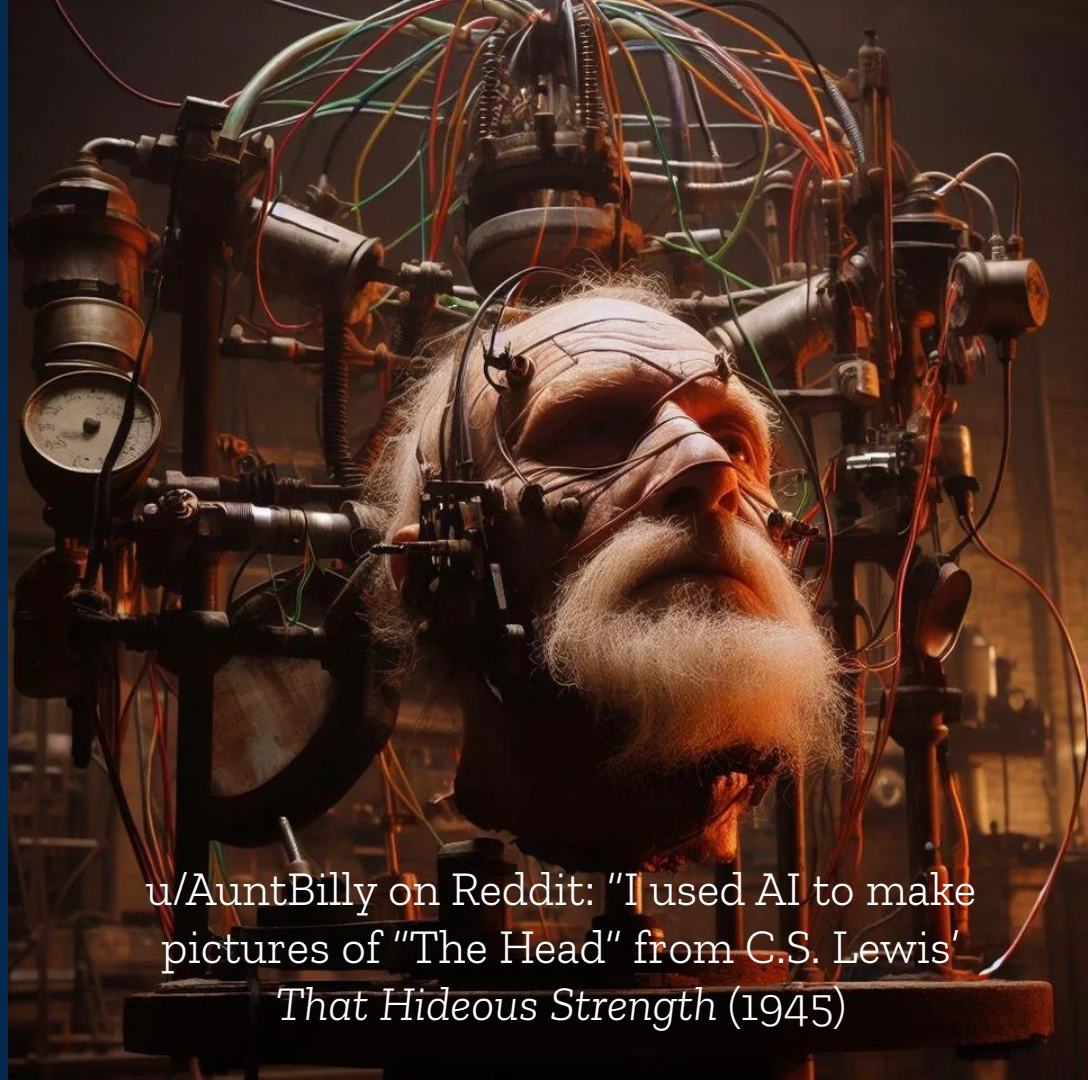
ML needs data:

- self-play (games)
- [datasets](#)



BTW:

Another early
view of "AI"



u/AuntBilly on Reddit: "I used AI to make
pictures of "The Head" from G.S. Lewis'
That Hideous Strength (1945)"

History: The Rise of ML in AI, i.e.

Data+Compute

Neural networks (one option for ML) were a cute and a bit useful, but were slow..

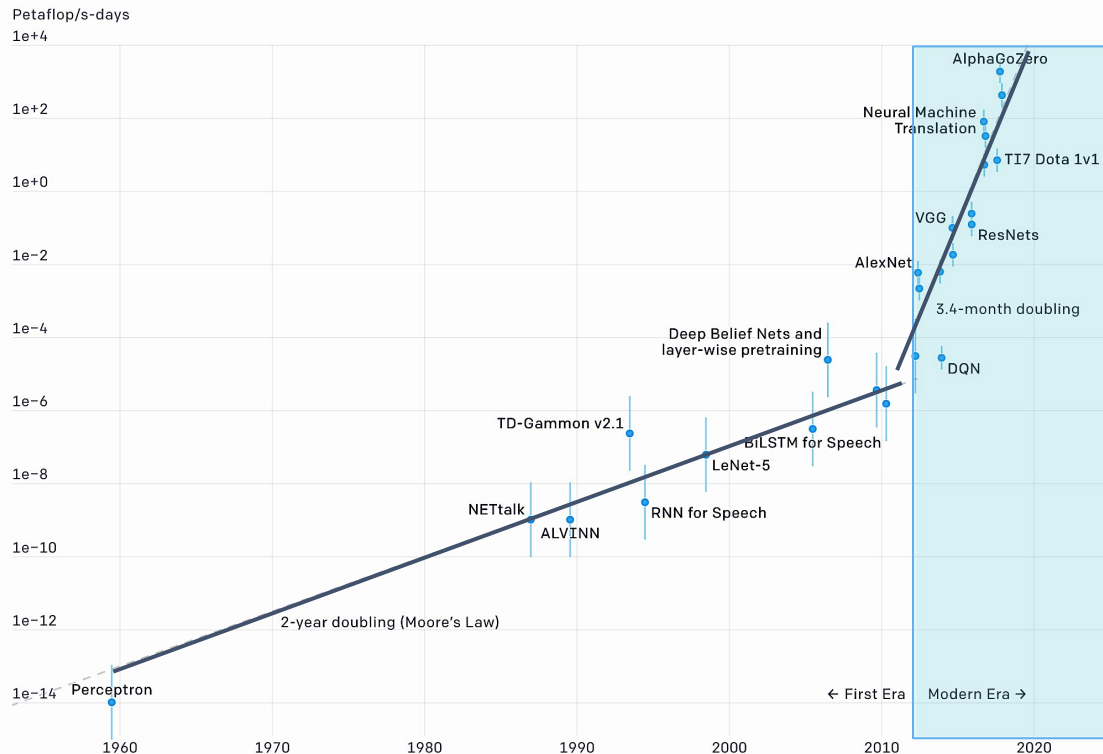
Then...

- the Internet provided the data
- GPUs provided the compute (Raina 2009, Ciresan 2010, Krizhevsky 2010,...)
- 2010s: Research(ers) from academia moved toward industry: access to data & compute



1993: Yann LeCun showing off Bell Labs' neural network-based character recognition

"Two Distinct Eras" (of Compute)



Moore's Law:
transistor density
doubles every 18
months.

AI training compute
doubles every ~3.5
months

Source:
OpenAI, 2018

Data: "Fair Use": Academia → Industry

- Grad students & profs left universities & carried their data policies with them
- Scraping the internet was/is standard, especially for education & non-commercial research
 - typically fair use for academic projects

Example ML datasets:

- [CelebA](#) (commercial photos)
- The Pile (included books)
- Anime Faces

OpenAI & Google are getting sued because they trained on The Pile

No accident that leaders in AI are "social" platforms (Meta, Google, ByteDance) b/c DATA!

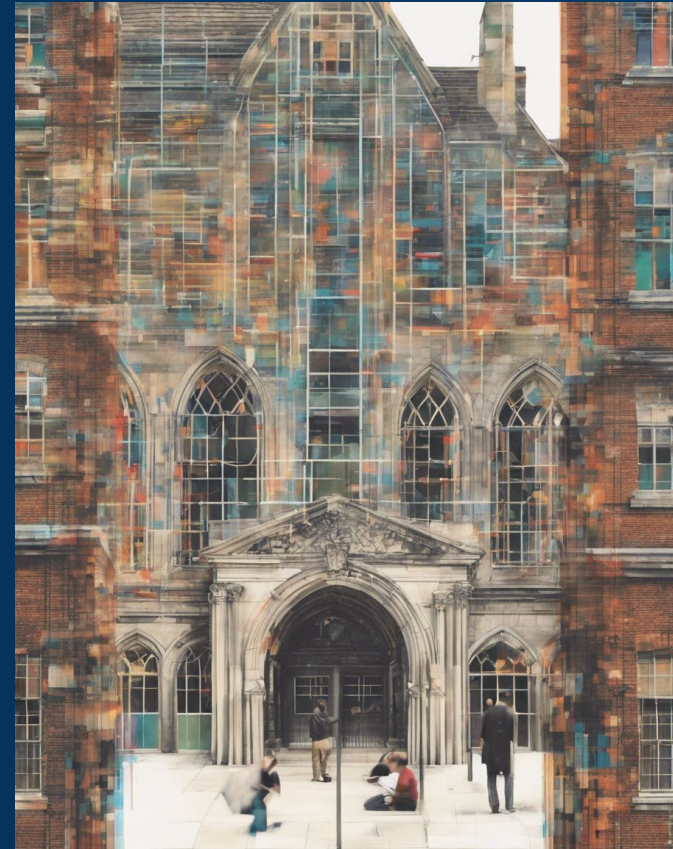


CelebA examples, all copyrighted material by professional photographers

Questions About Your (University's) Data

You already...

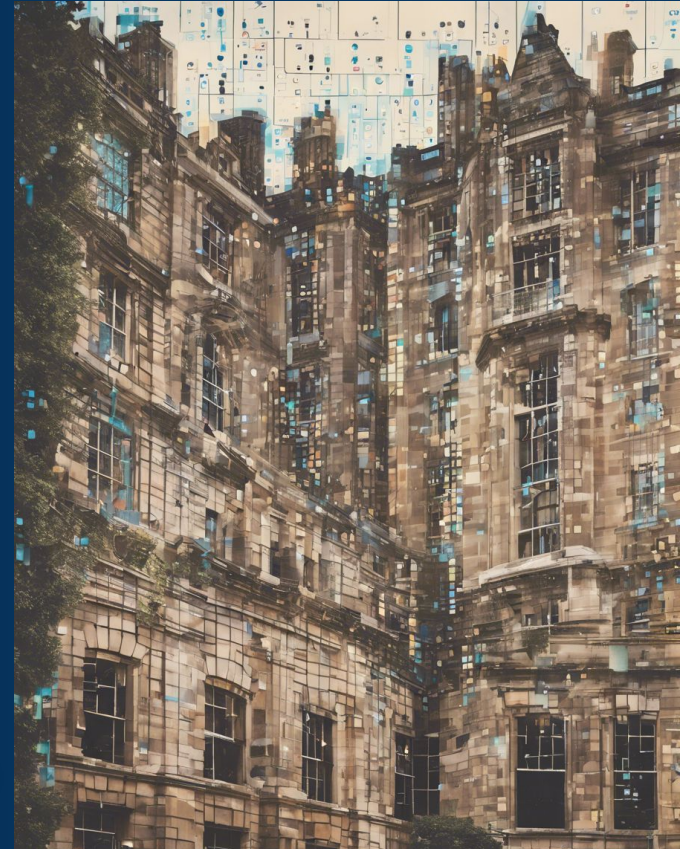
- use (aggregated) data, e.g.
 - ❖ for admissions / institutional planning
 - ❖ for assessment of students *and* faculty
- have policies protecting some kinds of data and/or data acquisition
 - ❖ student privacy
 - ❖ IRBs, e.g. for experiments on human subjects
- share data with 3rd parties
 - ❖ LMS (Blackboard, Canvas, etc, etc) / e-textbook providers
 - ❖ (what are *they* doing with that data?)



Questions About Your (University's) Data, pt.2

Now, what about...

- Internal: Faculty/students doing AI research, using
 - surveillance cams? (Duke debacle)
 - using "public(-ish)" data for questionable studies (Criminality, "Gaydar")
 - scanning student assignments/portfolios?
 - and instructor feedback?
 - My story: Using student audio production final projects
 - Who owns student/faculty data? (UK higher ed strike & recording lectures)
- When 3rd Parties come offering data licensing deals?
 - You are sitting on "gold mines" of text, visual, audio data – all *annotated* and *scored by experts!* Tech companies (will) want to train their AIs on your data!
 - Again: Who owns the data? Legal / Philosophical (Florida)



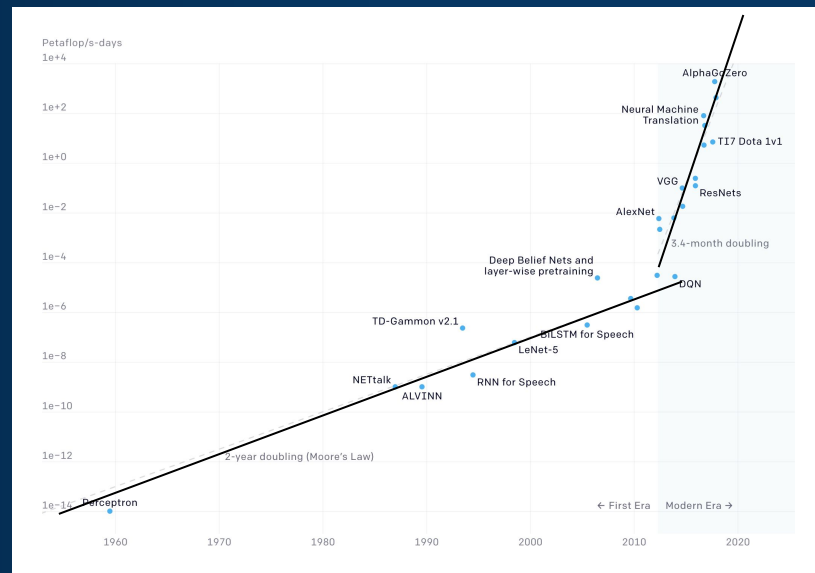
AI is Compute

CPUs are cheap and ubiquitous
But CPUs are slow compared to
GPUs

(e.g. 100x slower)

- inconvenient for inference
- impossible for training anything but tiny models

GPUs were and ARE hard to get, have gotten bigger. Gaming GPUs rarely cut it

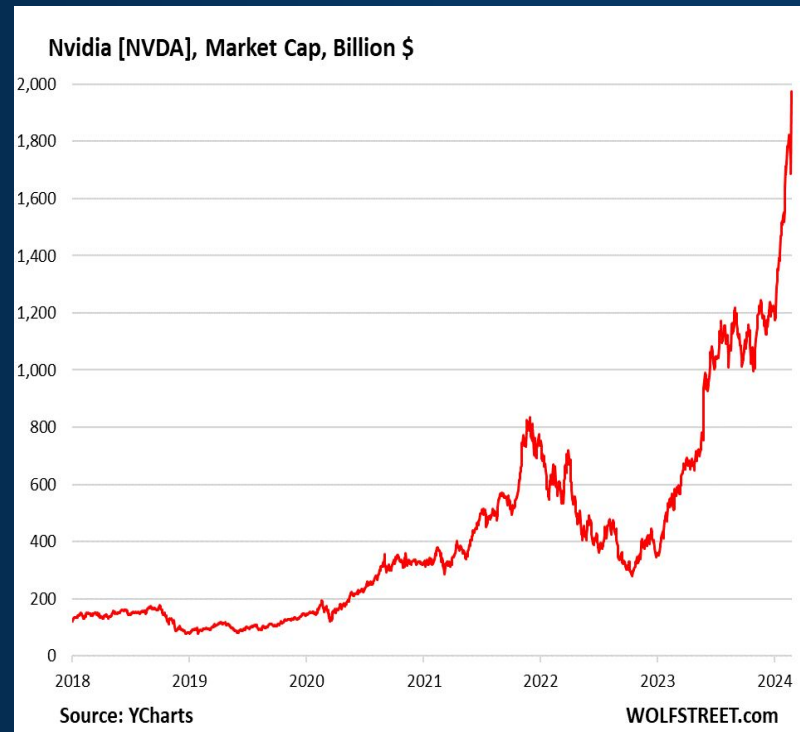


Same OpenAI graph as before

AI is Compute

NVIDIA makes GPUs.

Kind of like selling mining tools during the Gold Rush!



Compute: What's a small college to do?

Large schools get/build clusters, e.g. U. of Florida
LLMs typically* require > 24 GB of VRAM (inference)
"Specialized equipment," & they tend to sell out!

So you use a cloud service provider:

- you pay them money
 - you give them your DATA.
 - Colab, Replicate, Runway,
- ...My story: AWS SageMaker for Belmont

Focus on options with lower resource needs:

- Stable Diffusion (images), Whisper, LLaMA 7B
- fine-tuning e.g., LoRA training.



*But! New paper: "The Era of 1-bit LLMs", Feb 27, 2024 <https://arxiv.org/abs/2402.17764>

Small Uni, Compute + Data, Cont'd

Someone will tell you, "We need a chatbot for our university web page!"

- Who's going to power it?
- You probably won't make your own, so a company will lease it to you.
- Will the training data be "ethically sourced"?
- What happens to the data harvested from it (i.e. user prompts)?
- How will you handle the fallout when (not if) the chatbot says something embarrassing/horrific to someone?
 - "You don't need that class to graduate", "Tuition is free!", "We're a historically white university."
- What about when people ask you about the carbon footprint
 - of training?
 - of continued use?

AI is Made of People

Built by humans, with training data on humans, to separate humans from money, to shift power

- Training data is usually human-generated and/or about humans
- AI companies are companies,
 - AI ethics is business ethics
- Many problems in the world don't admit technical solutions
- AI systems merit regulation,
 - we need humans involved in all areas



How to make "more capable" humans? *(You do that already!)*

- Capable of *what*?
- College is not about learning how to regurgitate facts / produce artifacts
- College exists to provide a *transformative experience*.
 - How I retooled one class re. the rise ChatGPT...
- (Most college majors aren't for "job training")
- *The humanities* – philosophy, history, poli sci, literature, – are as necessary as ever for a functioning society
- "AI Ethics" folks moved from principles to *regulation & auditing*
- Christian university uniquely addresses human needs in a world decimated by AI-driven algorithms, can offer leadership toward solutions for problems "orthogonal" to technology

Thanks!

Contact me: [@drscotthawley](#) on socials, GitHub, etc

I blog tutorials at [drscotthawley.github.io](#), "Deep Learning & AI Ethics" class lessons on GitHub

Belmont's starting a new CS grad school focused on AI!

QR Code for these slides:

