

# LOGISTICS

Join the Discord!

**<https://fsdl.me/llmhc-discord>**

Slides will be shared.

Talks will be recorded  
and posted online.

WiFi: llmhc-2023

password: SolidGoldMagikarp



# Tensor Tier Sponsors



# Matrix Tier Sponsors



# Vector Tier Sponsors



# Compute Credit Donors



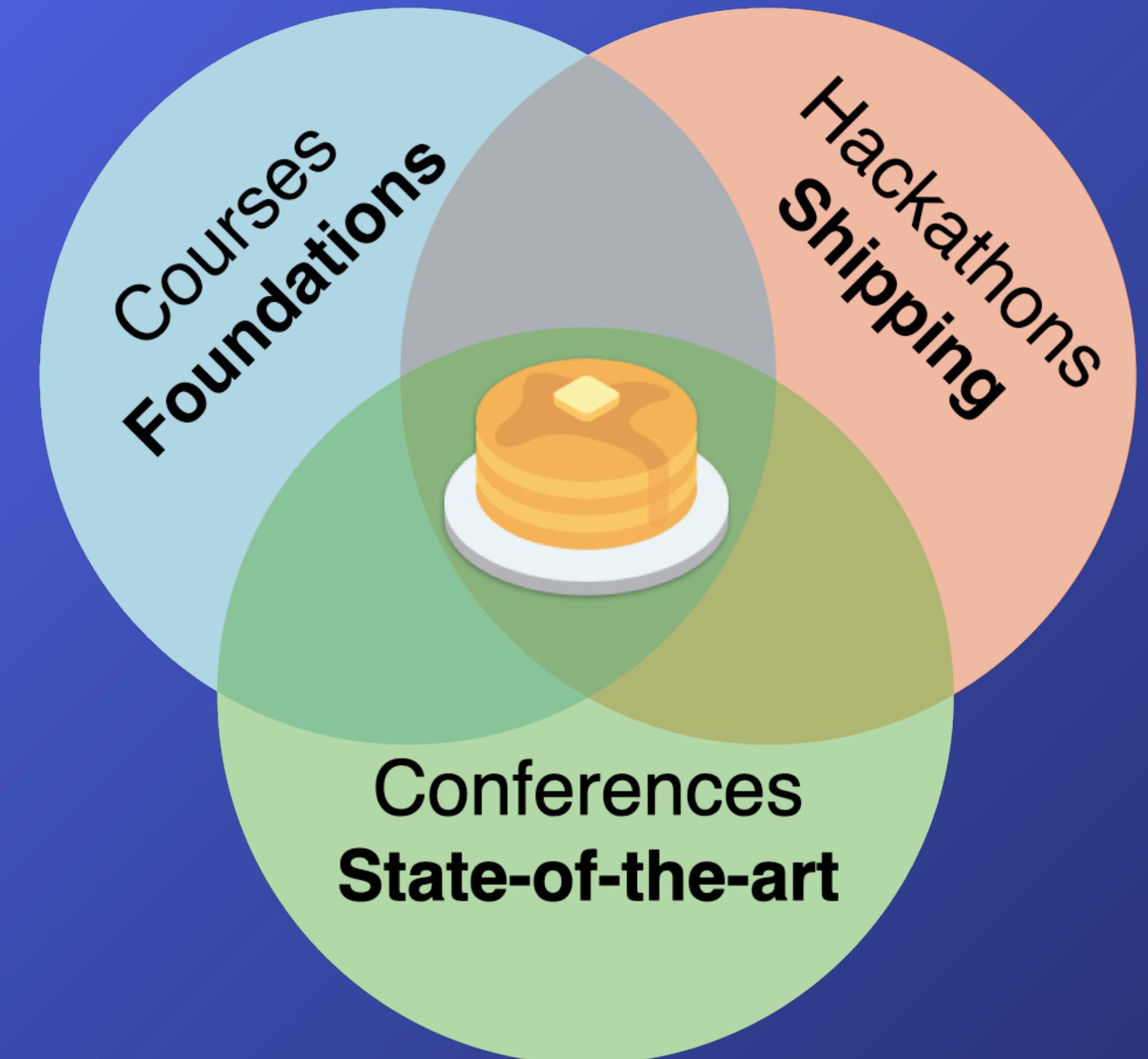
# LLM Bootcamp 2023

## Launch an LLM App in an Hour

Charles Frye

APRIL 21, 2023





# Instructor Team



@charles\_irl



@sergey\_karayev



@josh\_tobin\_

Berkeley PhD alums with years of experience  
building and teaching about ML-powered apps.



# Agenda

00

## INTRODUCTION

Why are we here?

01

## PROTOTYPING & ITERATION

Move fast  
and write prompts.

02

## DEPLOYING AN MVP

No one can  
stop you.

03

## NEXT STEPS

Why are there 10  
more hours of talks?

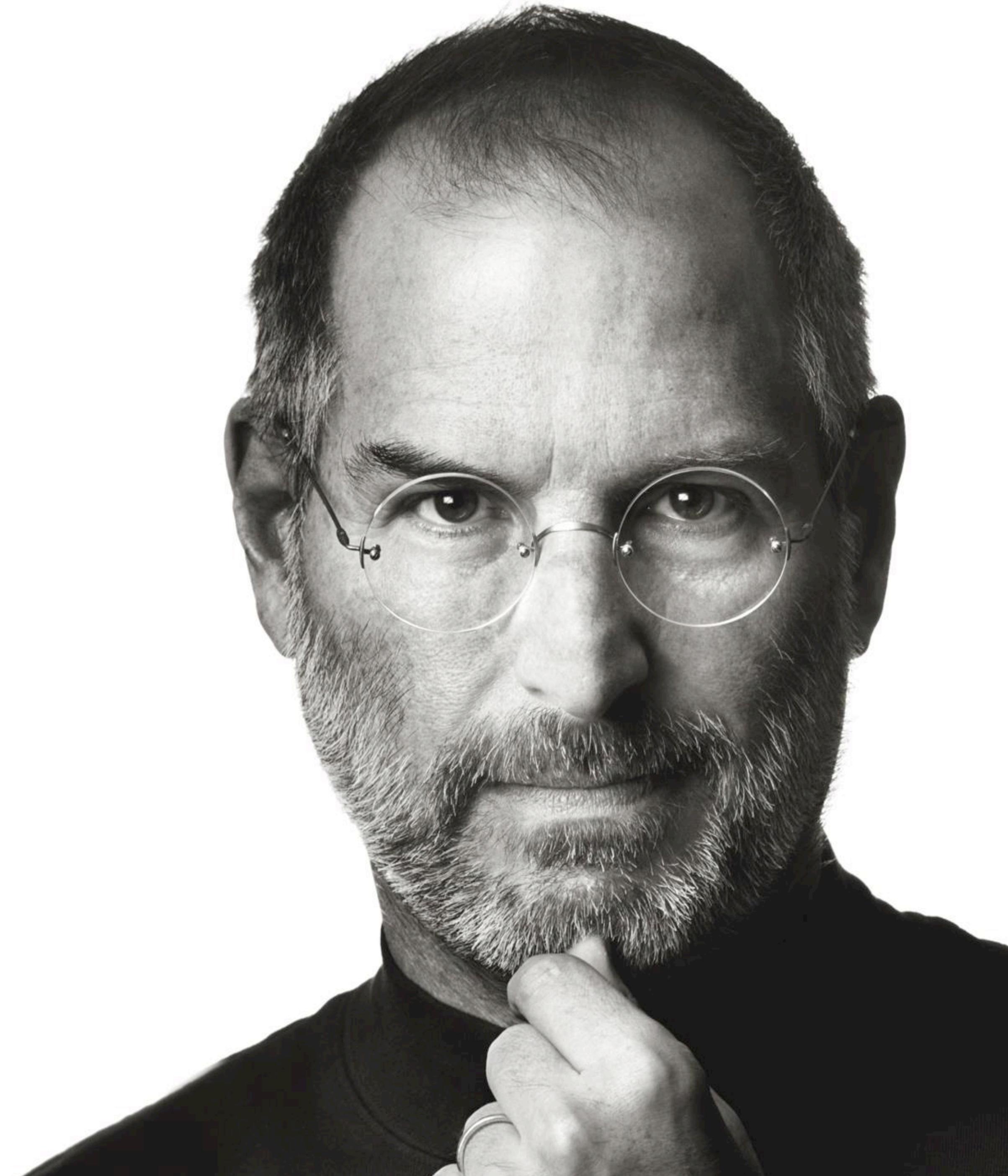
00

# Introduction



# Imagine a world...

- Disturbingly simple and heuristic computer programs are used to **mimic human cognitive processes**.
- Humans are forming emotional **attachments to chatbot therapists**.
- Computers can **play chess**, **write mathematical proofs**, and even **pass some exams**.
- You can even get **investment advice from an artificial intelligence!**



# It happened in 1965.

By using for heuristic programs tasks where there are good measures of human performance, much has been learned about the general level of complexity and sophistication that is required if programs are to operate at levels of effectiveness comparable to human levels. We know now that such programs need not be enormous in size or terribly intricate provided that they incorporate powerful heuristics for searching the problem space selectively. Thus, a recent chess program of modest size matches human abilities in discovering deep checkmating combinations.<sup>4</sup> Programs for discovering proofs for theorems in geometry and for finding indefinite and definite integrals are powerful enough to pass standard school examinations in these subjects.<sup>5</sup> A program of very simple structure, having access to a large body of data, simulates closely the behavior of an investment trust officer selecting portfolios for trust funds.<sup>6</sup>

## HEURISTIC PROBLEM SOLVING BY COMPUTER\*

HERBERT A. SIMON AND ALLEN NEWELL  
Carnegie Institute of Technology  
Pittsburgh, Pennsylvania



# What's different this time?

One tool can do it all:  
a large language model.

# Language Models model language.

Imagine a paragraph appearing word by word.

If you learned to guess the next word,  
you'd get pretty smart.

Especially if you could predict things like:

$$2 + 2 = 4$$

```
def square(x): return x * x
```

The chemical formula for glucose is C<sub>6</sub>H<sub>12</sub>O<sub>6</sub>

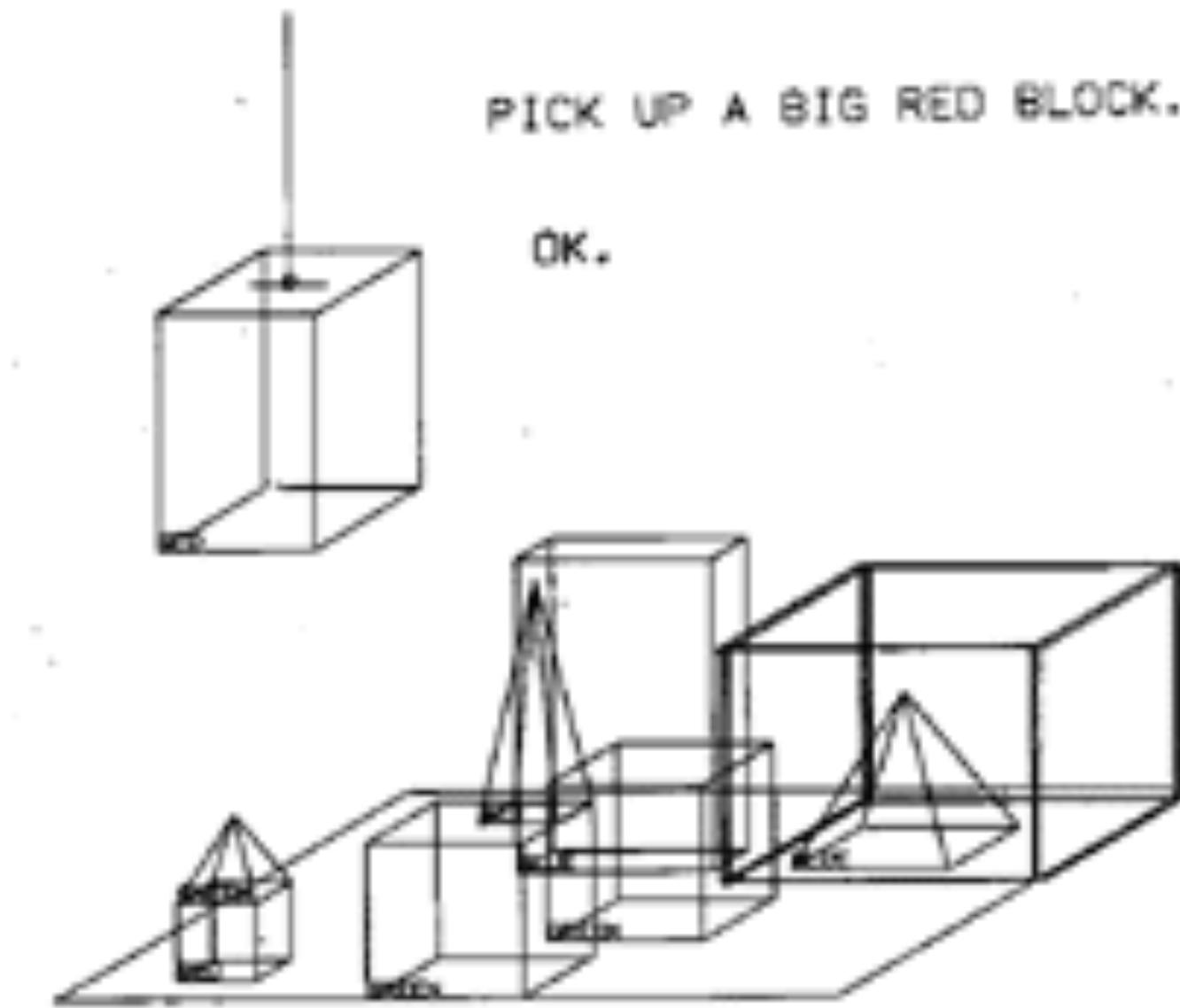
*Large Language Models are really good at it.*

# LLMs unlock LUIs: Language User Interfaces.

```
=====
EEEEEEEEE L      IIIIIIII ZZZZZZZZ      AAA
E          L      I          Z      A   A
E          L      I          Z      A   A
EEEEEE    L      I          Z      A   A
E          L      I          Z      AAAAAAAA
E          L      I          Z      A   A
EEEEEEEEE LLLLLLLL IIIIIIII ZZZZZZZZ      A   A

=====
ELIZA > (PLEASE TYPE IN ALL CAPS) WHAT'S YOUR NAME DEAR ?
AMIT
ELIZA > HI, I'M ELIZA. WHAT DO YOU WANT TO TALK ABOUT ?
AMIT > CAN I TALK ABOUT MY PROBLEM ?
ELIZA > SURE... !
AMIT > DUE TO ON GOING EXAMS IN COLLEGE, I HAVE A LOT OF STRESS.
ELIZA > PLEASE RELAX A LITTLE. SLEEP WELL.
AMIT > THANKS FOR YOUR ADVICE.
ELIZA > NO MENTION.
AMIT > BYE
ELIZA > BYE AND KEEP IN TOUCH...
```

# ELIZA, 1966



# SHRDLU, 1970

• Home • About • Help • Corporate Services

**PERSONAL Jeeves is here!**

**May I Suggest:**

- PERSONAL JEEVES
- ANSWER POINT
- MONEY
- TRAVEL
- HEALTH
- COMPUTERS
- ENTERTAINMENT
- HOME & FAMILY
- SHOPPING

**ASK JEEVES FOR Kids!**

**TAXING QUESTIONS?**  
Jeeves Tax Guide to the rescue! ▶

**Ask Jeeves<sup>SM</sup>** Ask.com

**Just Curious,  
Jeeves**

**Have a Question?**  
Just type it in and click **Ask!**

Ask!

**Most Recent Questions About Business:**

**Where can I find a currency converter?**

Ask!

**What are people asking RIGHT NOW?** ▶

• Make Jeeves Your Homepage • Ask Jeeves U.K. • Advertise On Jeeves • Become an Affiliate • Investor Relations

# Ask Jeeves, 1996

# Language modeling is considered AI-Complete\*.

Computational Complexity [19]. According to the Encyclopedia of Artificial Intelligence [40] published in 1992 the following problems are all believed to be AI-Complete and so will constitute primary targets for our effort of proving formal AI-Completeness on them [40]:

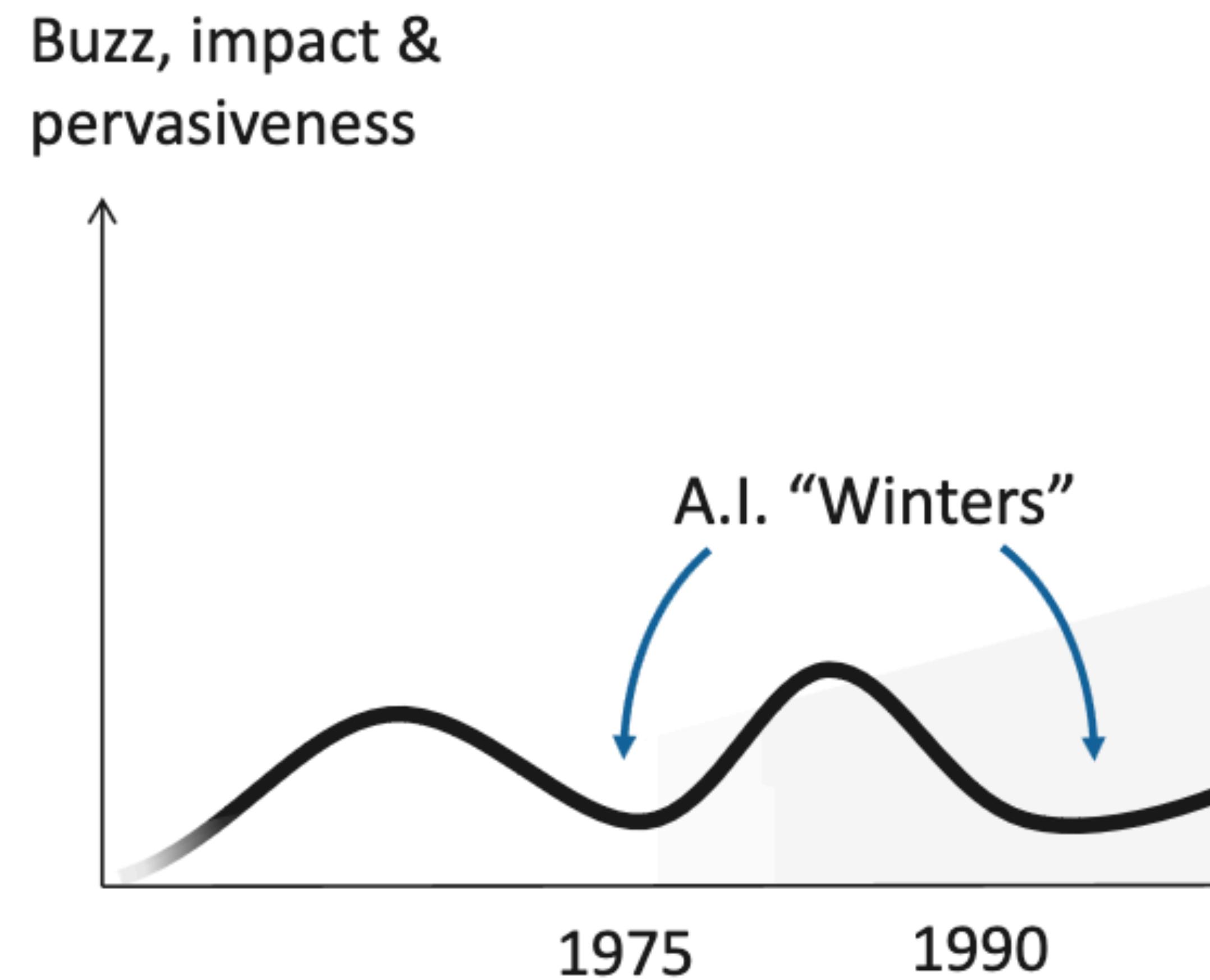
- **Natural Language Understanding**

\*this concept is slipperier than NP-Complete etc.



# What happened last time?

# Over-selling and under-delivering led to an “AI Winter”.



[https://www.cambridgewireless.co.uk/media/uploads/resources/AI%20Group/AIMobility-11.05.17-Cambridge\\_Consultants-Monty\\_Barlow.pdf](https://www.cambridgewireless.co.uk/media/uploads/resources/AI%20Group/AIMobility-11.05.17-Cambridge_Consultants-Monty_Barlow.pdf)

# We failed to deliver products to match high expectations.

“high hopes that are very far from having been realized”

“[N]otorious disappointments in the area of machine translation, where enormous sums have been spent with very little result”

“Suggestions from recent research ... augur badly for the future availability of machine-translation programs versatile enough to be commercially valuable.”





To avoid AI Winter, we need to  
build products that people value.



# Good news: a lot more than just research!

**Spaces**  
Discover amazing ML apps made by the community!  
Create new Space or learn more about Spaces.

Spaces of the week 🔥

- Running on A10G Baize 7B 243 project-baize 10 days ago ai-force
- Running on A10G VideoCrafter 120 VideoCrafter about 21 hours ago abhishek

**microsoft / DeepSpeed**  
DeepSpeed is a deep learning optimization library that makes distributed training and inference easy, efficient, and effective.

**Nat Friedman ✅ @natfriedman**

Finally, there's tinkering.

**Nat Friedman ✅ @natfriedman · Sep 30, 2022**  
Still not enough tinkering happening, for whatever reason  
[Show this thread](#)

3:54 PM · Apr 12, 2023 · 116.4K Views

<https://twitter.com/natfriedman/status/1646285680941862912>

# Bad news: the gap from demo to product is big.



<https://www.youtube.com/watch?v=fmVWLr0X1Sk>

# Like, REALLY big.

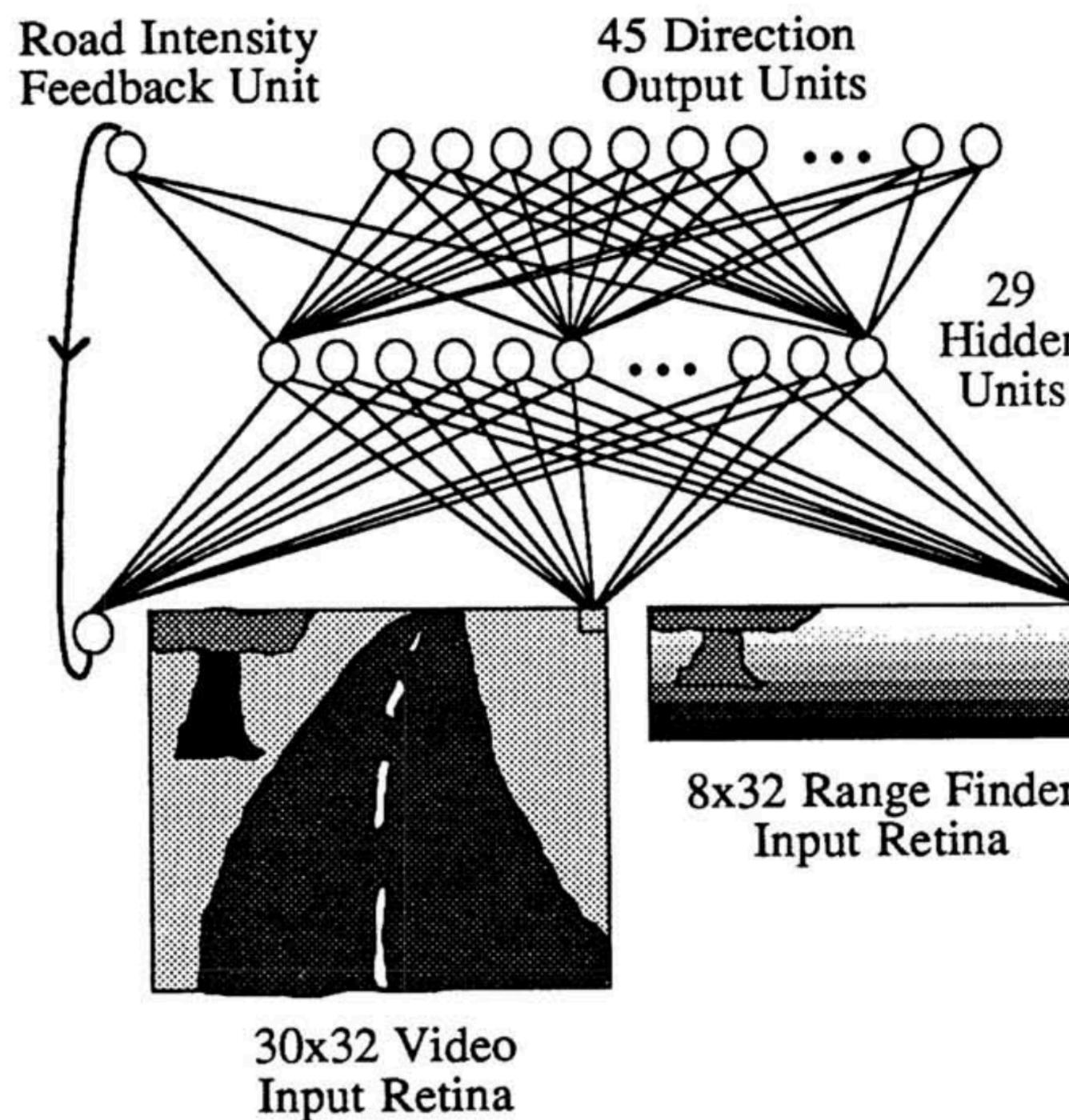


Figure 1: ALVINN Architecture

## ALVINN: AN AUTONOMOUS LAND VEHICLE IN A NEURAL NETWORK

Dean A. Pomerleau  
Computer Science Department  
Carnegie Mellon University  
Pittsburgh, PA 15213



Figure 3: NAVLAB, the CMU autonomous navigation test vehicle.

# NeurIPS, 1988

# Good news: we're building products (finally!)

2 minute read · February 2, 2023 7:33 AM PST · Last Updated 2 months ago

## ChatGPT sets record for fastest-growing user base - analyst note

By Krystal Hu



There's a new way to make video and podcasts. **A good way.**

Descript is the simple, powerful, and fun way to edit.



**GitHub**  
Copilot

<https://www.reuters.com/technology/chatgpt-sets-record-fastest-growing-user-base-analyst-note-2023-02-01/>



# So let's build an app!

01

# Prototyping & Iteration: The First Half Hour





# Takeaways

- “Foundation models” have unblocked a lot of applications
- Prototype with a high-capability hosted model in chat first
- Tinker with prompts and build on open source frameworks

02

## Deploying an MVP: The Second Half Hour



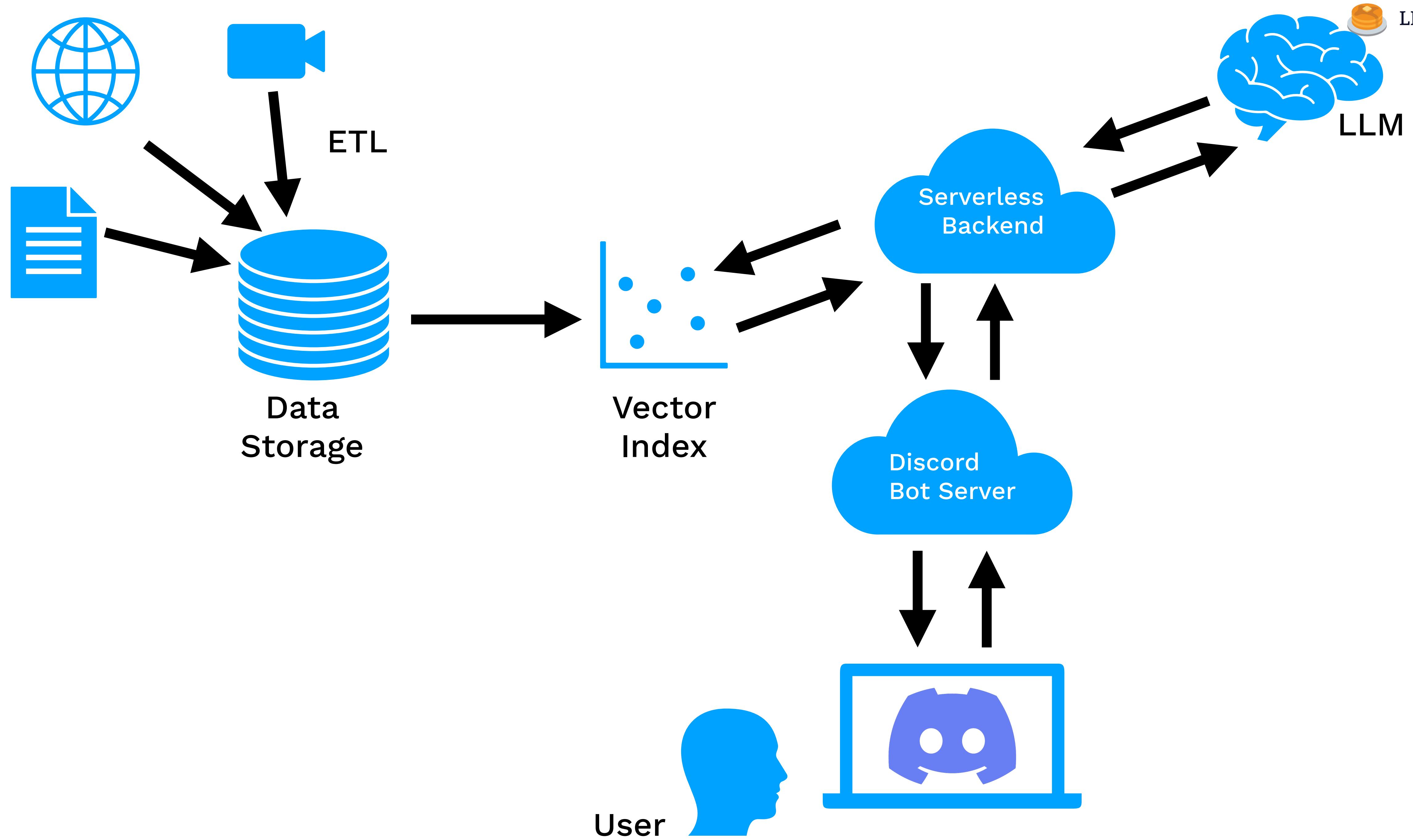


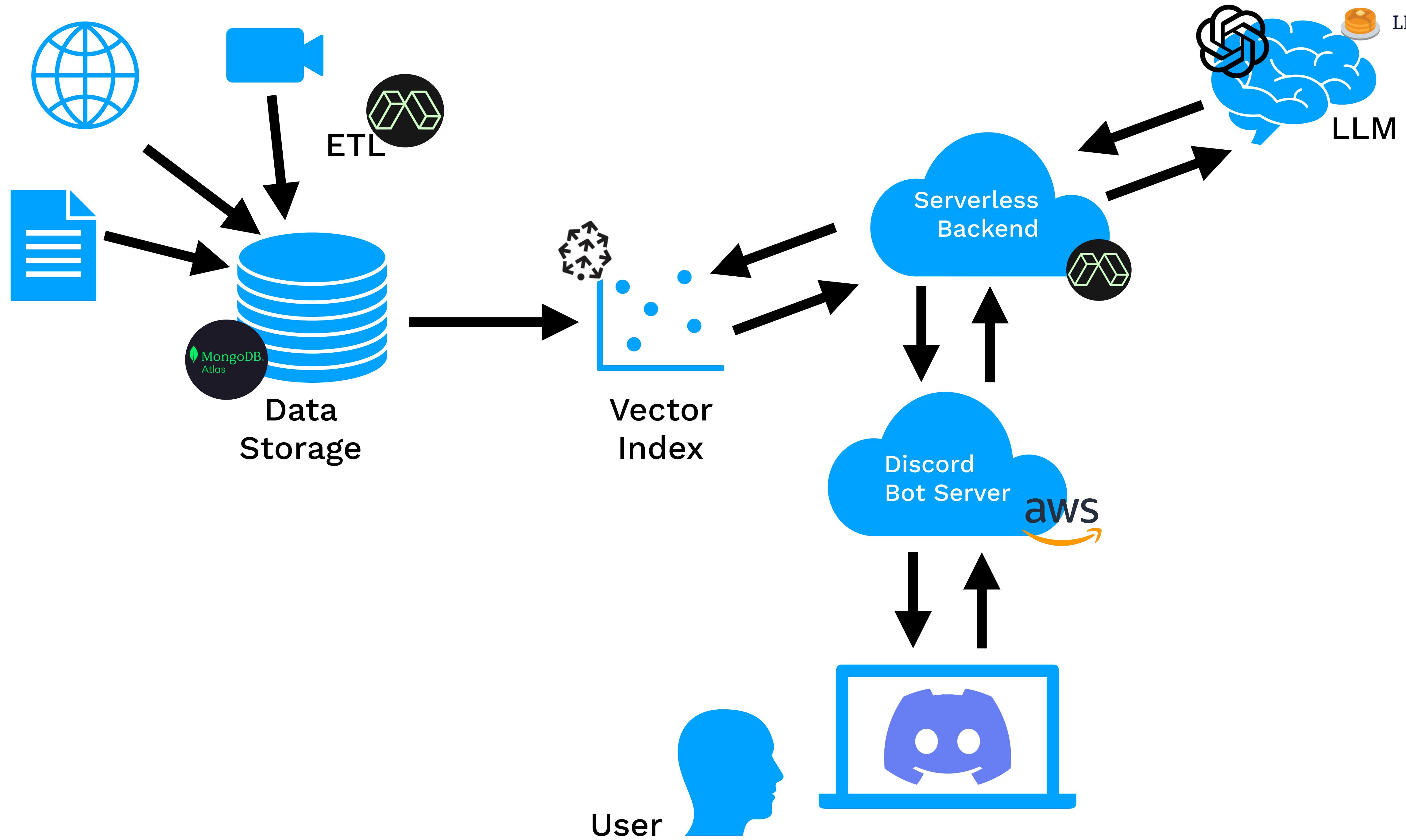
# Takeaways

- Cloud tooling makes it easy to get started
  - But make sure to limit spend!
- Find a simple UI and start getting feedback from users ASAP



Let's build a bot  
to do sourced Q&A  
on the FSDL Discord!







# Try it out!

#ask-fsdl-bot  
in the Discord.

We'll analyze its behavior  
later today!

03

## Next Steps: The Rest of Time





Doing it is easy,  
knowing what you're doing is hard.



# This morning: foundations.

## Schedule

Friday (April 21)	
9 am	Registration & Breakfast
10 am	 <b>Launch an LLM App in 1 Hour</b> <ul style="list-style-type: none"><li>• We'll go from idea to user-ready website in one hour</li></ul>
11 am	 <b>LLM Foundations</b> <ul style="list-style-type: none"><li>• Speed-run core ideas in ML</li><li>• Learn the core concepts behind Transformer architecture</li><li>• Get a tour of notable LLMs, their training data, and abilities</li></ul>
12 pm	Lunch and networking

<https://fullstackdeeplearning.com/llm-bootcamp/#schedule>

# This afternoon: core LM skills.

1 pm

## ✨ Learn to Spell: Prompt Engineering and Other Magic

- Prompt techniques for language modeling, instruction, and agent simulation
- Arcana and curiosities from the world of prompt wizardry
- Heuristics for where LLMs can and cannot be effective

---

2 pm

## 🔨 Augmented Language Models

- Moving beyond pure language modeling with tools and retrieval
- Choosing between vector stores/indices (e.g. FAISS, Milvus, Pinecone, Weaviate, Vespa) for long-term memory

---

3 pm

Coffee and networking

# Late afternoon: repo and invited talks.

3:45  
pm

## askFSDL Walkthrough

- Detailed breakdown of a well-documented sample project demonstrating use of LLM APIs and frameworks, traditional and vector databases, and user feedback ingestion

---

4:30  
pm

## Invited Talk from **Richard Socher:** CEO and co-founder of You.com

---

5:15 pm

## A Fireside Chat with **Peter Welinder:** VP of Product at OpenAI



# Tomorrow morning: user experience.

9 am	Breakfast
9:30 - 10:45 am	 <b>Demo Garden</b> <ul style="list-style-type: none"><li>• Present your cool project</li><li>• Review other cool projects</li></ul>
11 am	 <b>UX for LUIs</b> <ul style="list-style-type: none"><li>• Principles of successful UX design</li><li>• Emerging Language User Interface design patterns</li><li>• Case studies on the best AI-powered apps today</li></ul>
12 pm	Lunch and networking



# Tomorrow afternoon: LLMOps and the future.

1 pm

## 🏎️ LLMOps: Deployment and Learning in Production

- How vendors like OpenAI, Anthropic, and Cohere compare with each other and with open-weights options like T5, Pythia, and LLaMA
  - How to manage prompts and programs, monitor models, trace chains, and record feedback
- 

2 pm

## 🔮 What's Next?

- What are the limits of scale, large and small?
  - Has multimodality unlocked general purpose robotics?
  - Is AGI already here?
  - Can we make it safe?
- 

3 pm

Coffee and networking



# The end: invited talks and a panel.

---

3 pm	Coffee and networking
4 pm	🎤 Invited Talk from <b>Harrison Chase:</b> Creator of LangChain
4:45 pm	🎤 Invited Talk from <b>Reza Shabani:</b> Training LLMs at repl.it
5:30 pm	Discussion Panel: <b>Building a Defensible Business</b>

---



LLMBC 2023



# Thanks!



@charles\_irl



@full\_stack\_dl

/imagine green parrot perched on a laptop, coder, #developer, pancake logo, pixar style, big round eyes, dorky glasses, 3d animation, Blender, 4k render, OpenGL - shaders, ray tracing