

# Getting Started with AI for .NET Devs with Semantic Kernel

A large QR code is positioned on the left side of the slide, enclosed in a white square with a soft glow. A thin white vertical line extends from the bottom center of the QR code down to the name 'Daniel Ward'.

Daniel Ward





# What this talk isn't

---

- In-depth
- Philosophy
- Morals/ethics



# Who am I?



LEAN  
TECHNIQUES

<https://leantechniques.com>



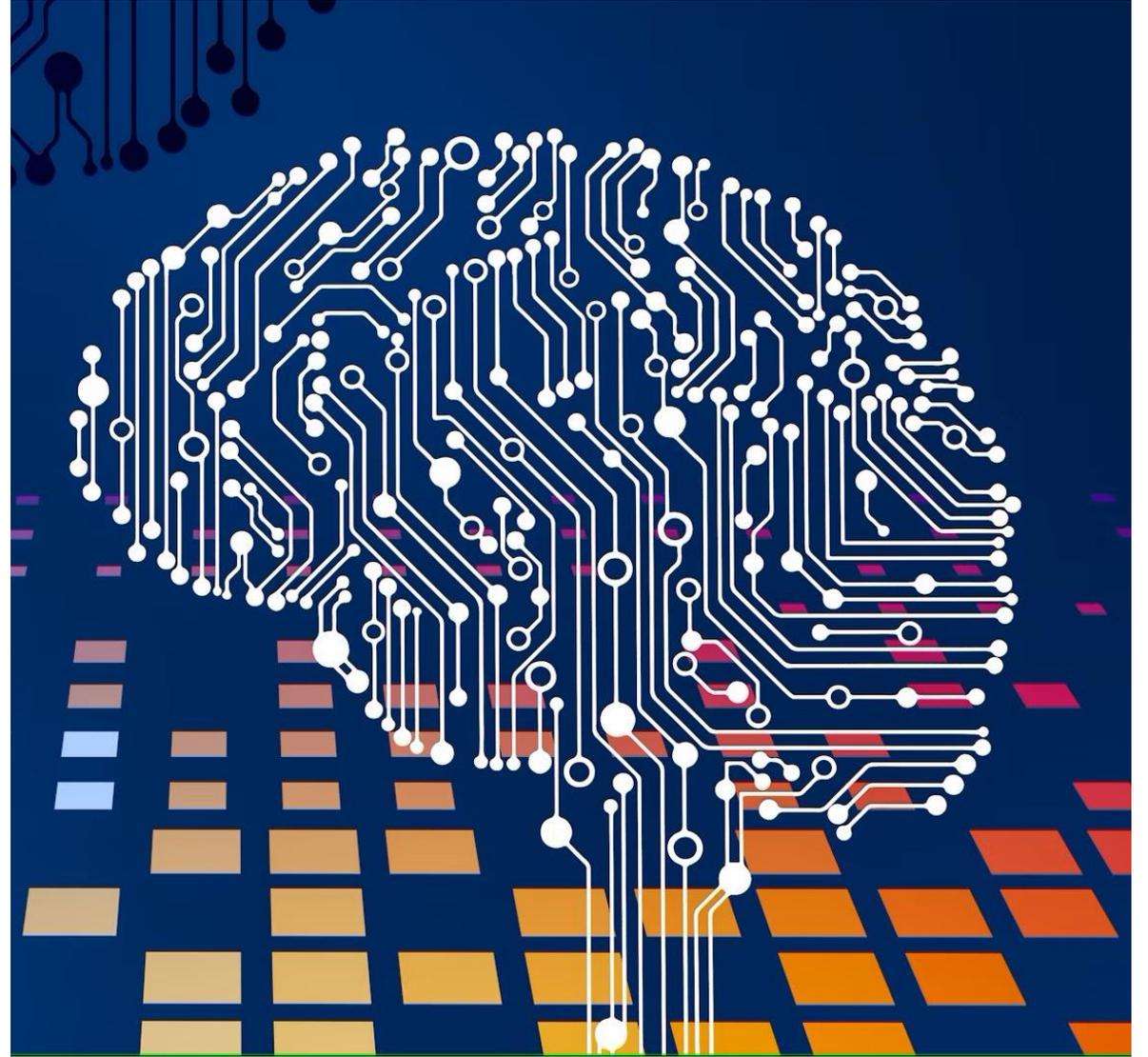
- Software developer, consultant
- Microsoft .NET MVP
- Co-organizer of the San Antonio/Austin .NET User Group
-  daninacan.com
-  @danielwarddev
-  daniel-ward-dev
-  danielwarddev.bsky.social

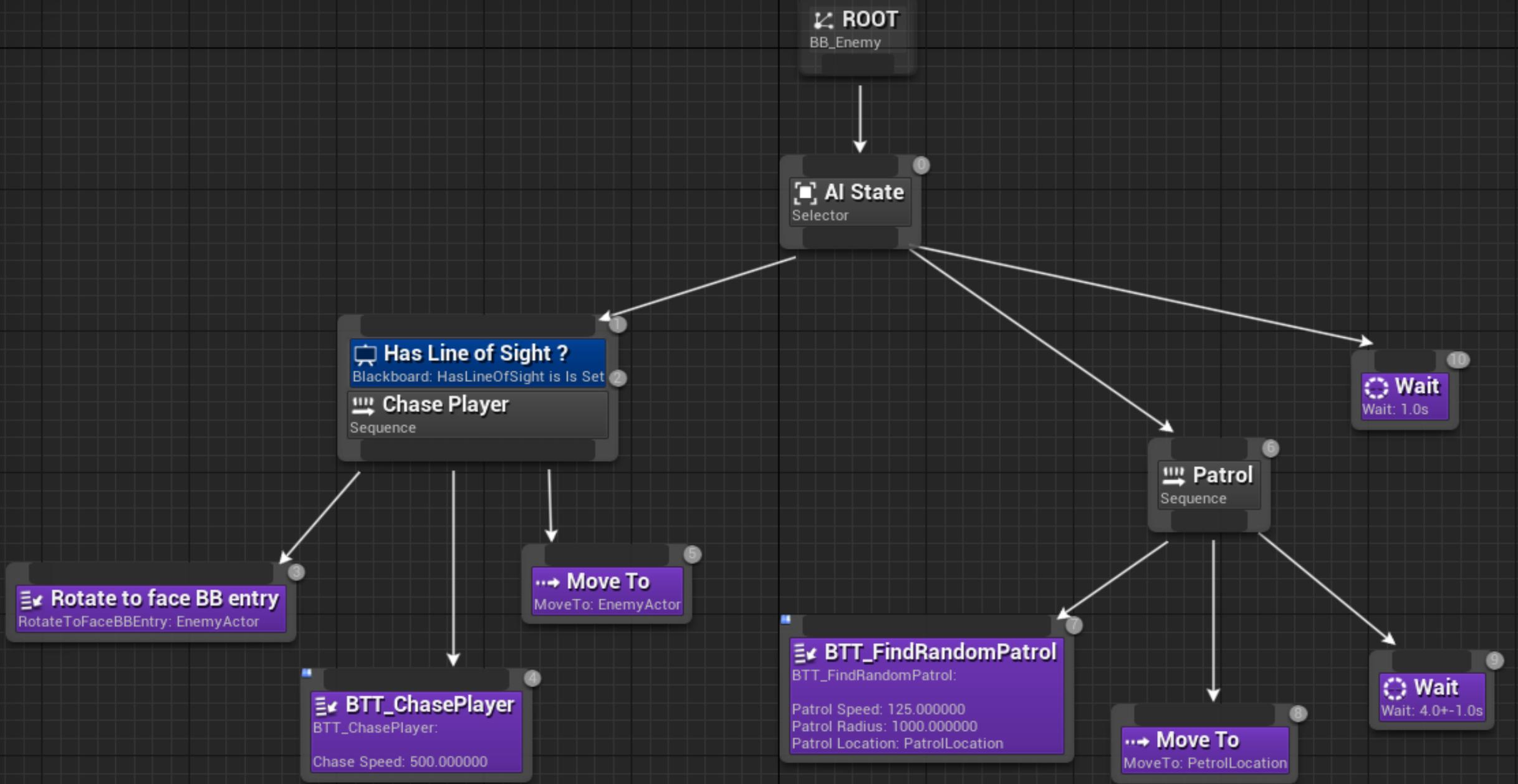


# What is AI, anyway?

---

- Simulating human intelligence
- Not new

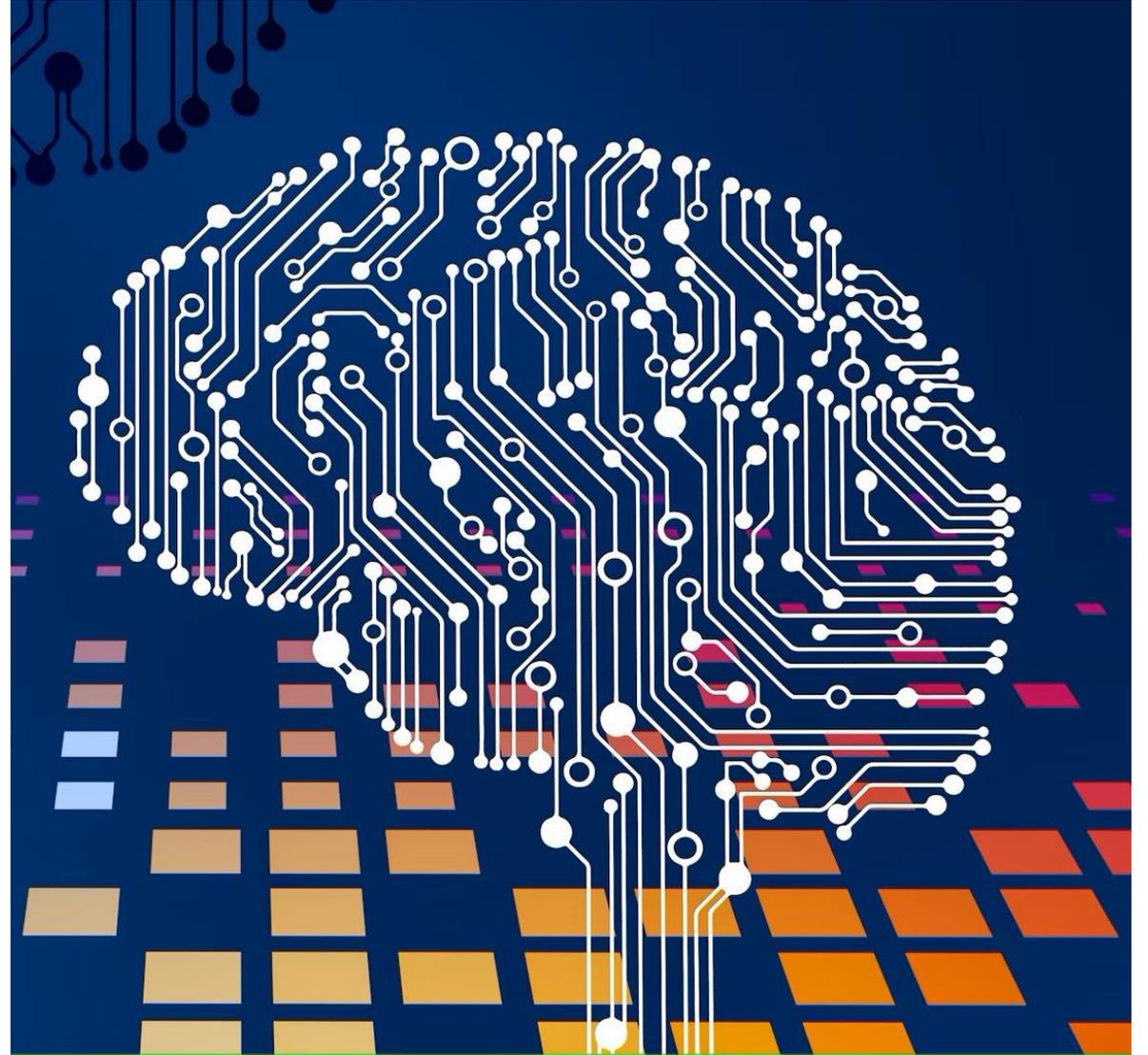




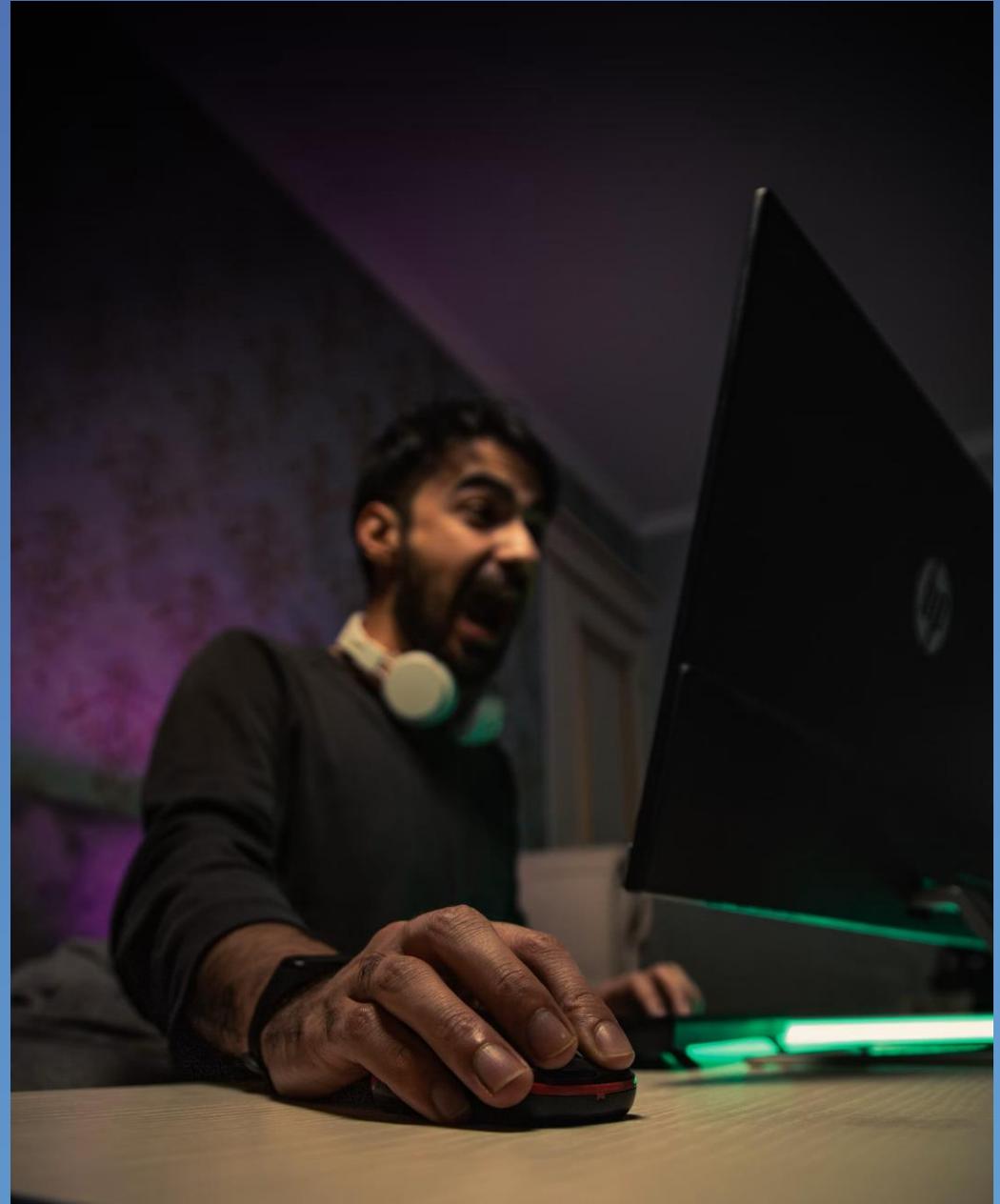
# What is AI, anyway?

---

- Simulating human intelligence
- Not new
- ★ 1950: Turing Test
- 1960s: learning checkers; algebra word problems; English; neural networks
- 2010s – machine learning (ML), deep learning
- 2010s, 20s – transformers, large language models (LLMs), generative AI
  - Copilot, ChatGPT, etc.



# AI vs ML vs deep learning vs generative AI

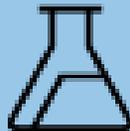


## Artificial Intelligence



Any technique that enables computers to mimic human intelligence. It includes *machine learning*

## Machine Learning



A subset of AI that includes techniques that enable machines to improve at tasks with experience. It includes *deep learning*

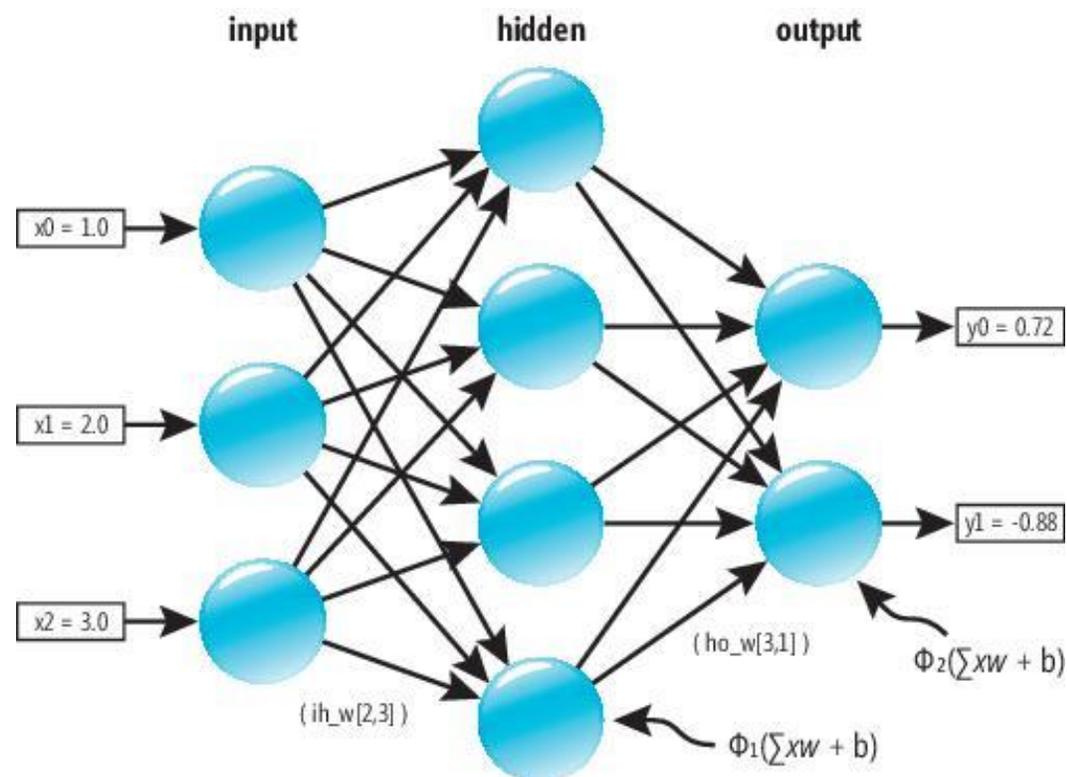
## Deep Learning



A subset of machine learning based on neural networks that permit a machine to train itself to perform a task.

# Neural network

- Loosely based on biological neurons/synapses
- Central part of modern ML
- 3 layers: input, hidden, output
- Each node is a math function
- The more layers, the deeper it is

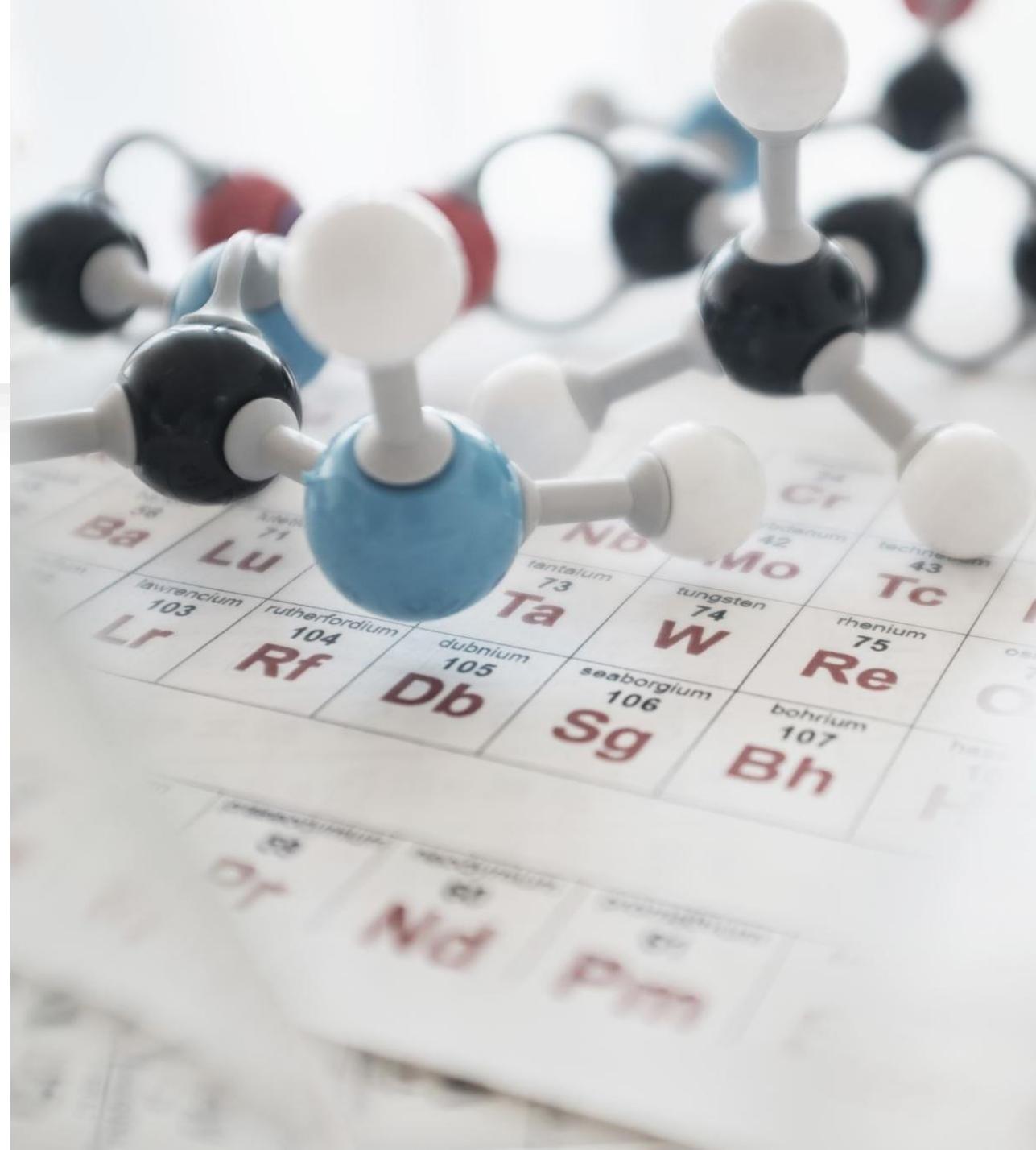


# Machine Learning (ML) vs deep learning

- Machine learning
  - Subset of AI
  - Specific kind of AI that allows machines to improve over time in their tasks with experience
  - Instead of hardcoding rules, train models to recognize patterns
- Deep learning
  - Subset of ML
  - Uses neural networks
  - Difference is in scale
  - Lots of hidden layers, lots of processing

# How is it possible AI creates new stuff?

- BBC. (2024, March 20). *NHS AI test spots tiny cancers missed by doctors.* <https://www.bbc.com/news/technology-68607059>
- MIT Technology Review (2023, December 14). *Google DeepMind used a large language model to solve an unsolved math problem.* <https://www.technologyreview.com/2023/12/14/1085318/google-deepmind-large-language-model-solve-unsolvable-math-problem-cap-set/>



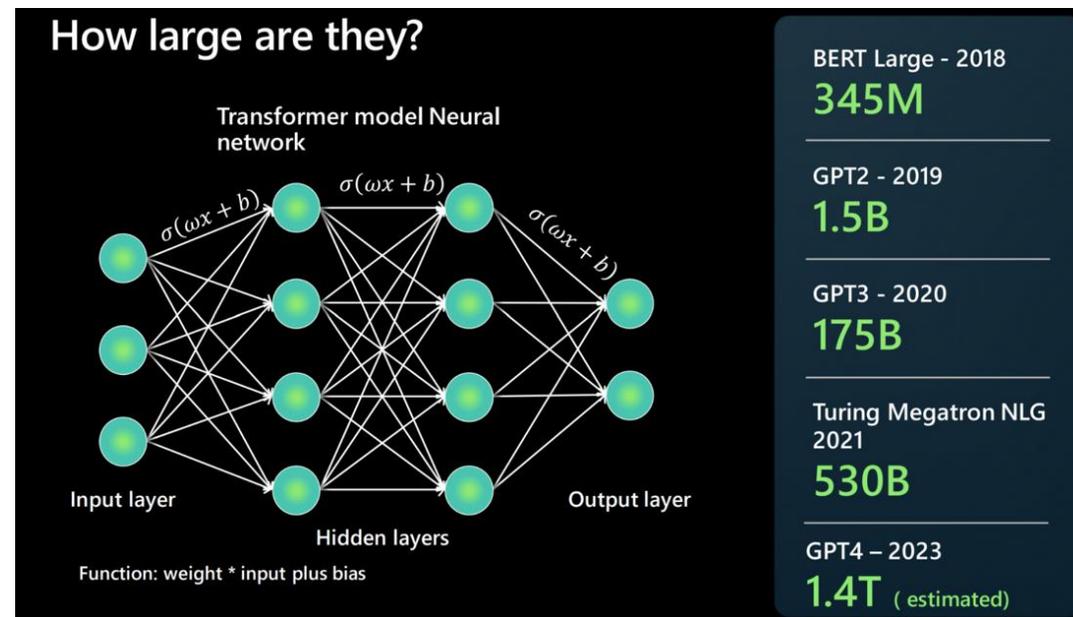
# Generative AI

---

- Subset of AI
- Uses techniques like deep learning to generate new content
- Depends on a huge amount of pre-trained knowledge
- ChatGPT, GitHub Copilot, etc.

# Large Language Models (LLMs)

- Type of ML model designed for natural language processing tasks



<https://microsoft.github.io/Workshop-Interact-with-OpenAI-models/llms/>

# GPT

- Generative pre-trained transformer
- Uses LLMs
- Takes  $n$  tokens as input and produces one as output
  1. How does it understand what we ask it and respond?
  2. How does it refer back to its own knowledge?

# Tokens

- “Building blocks” of input
- Middle ground between characters and words
- Model keeps a complete list of all possible tokens

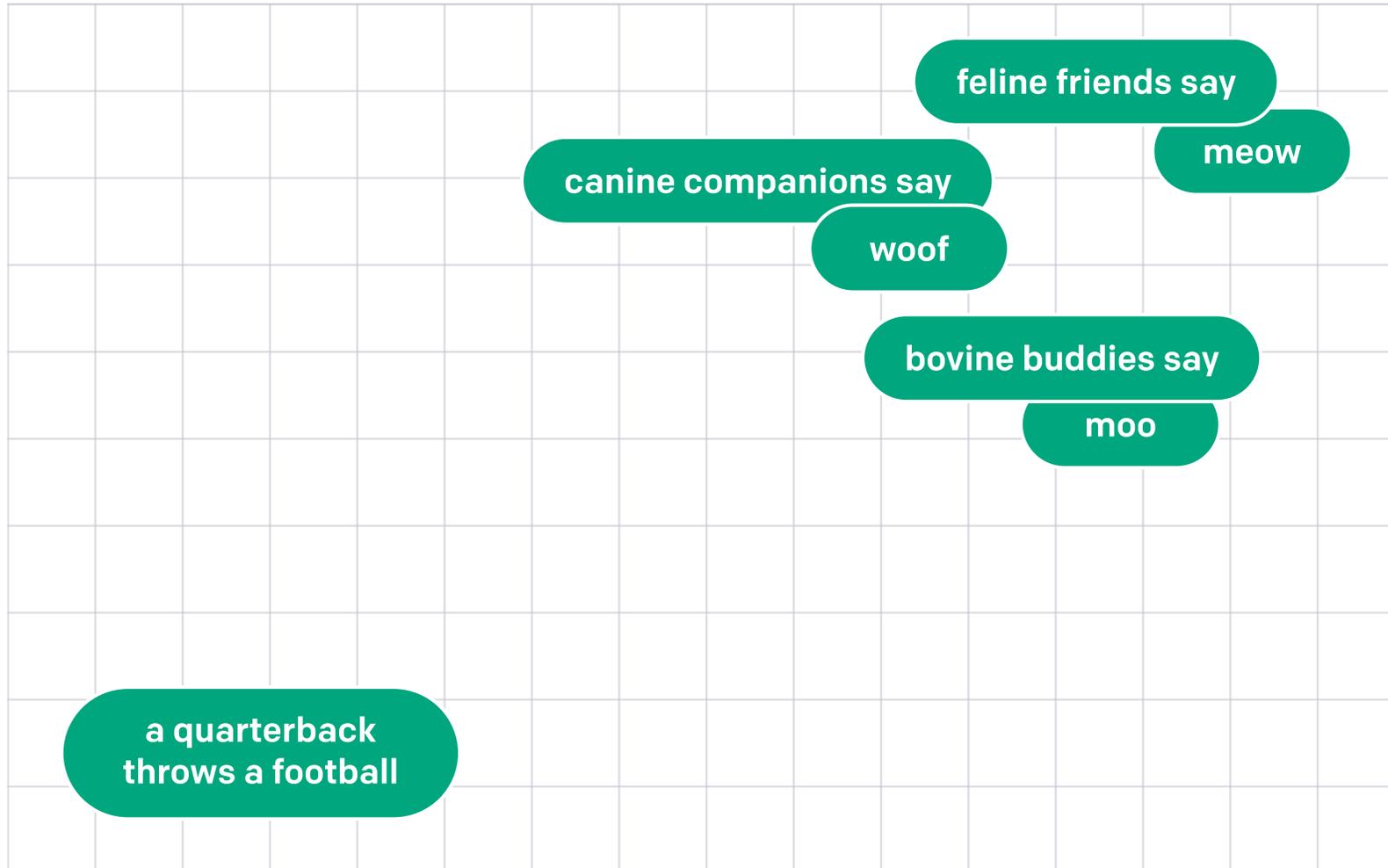
Tokens	Characters
9	51

The English language is sometimes incomprehensible.



# Vectors

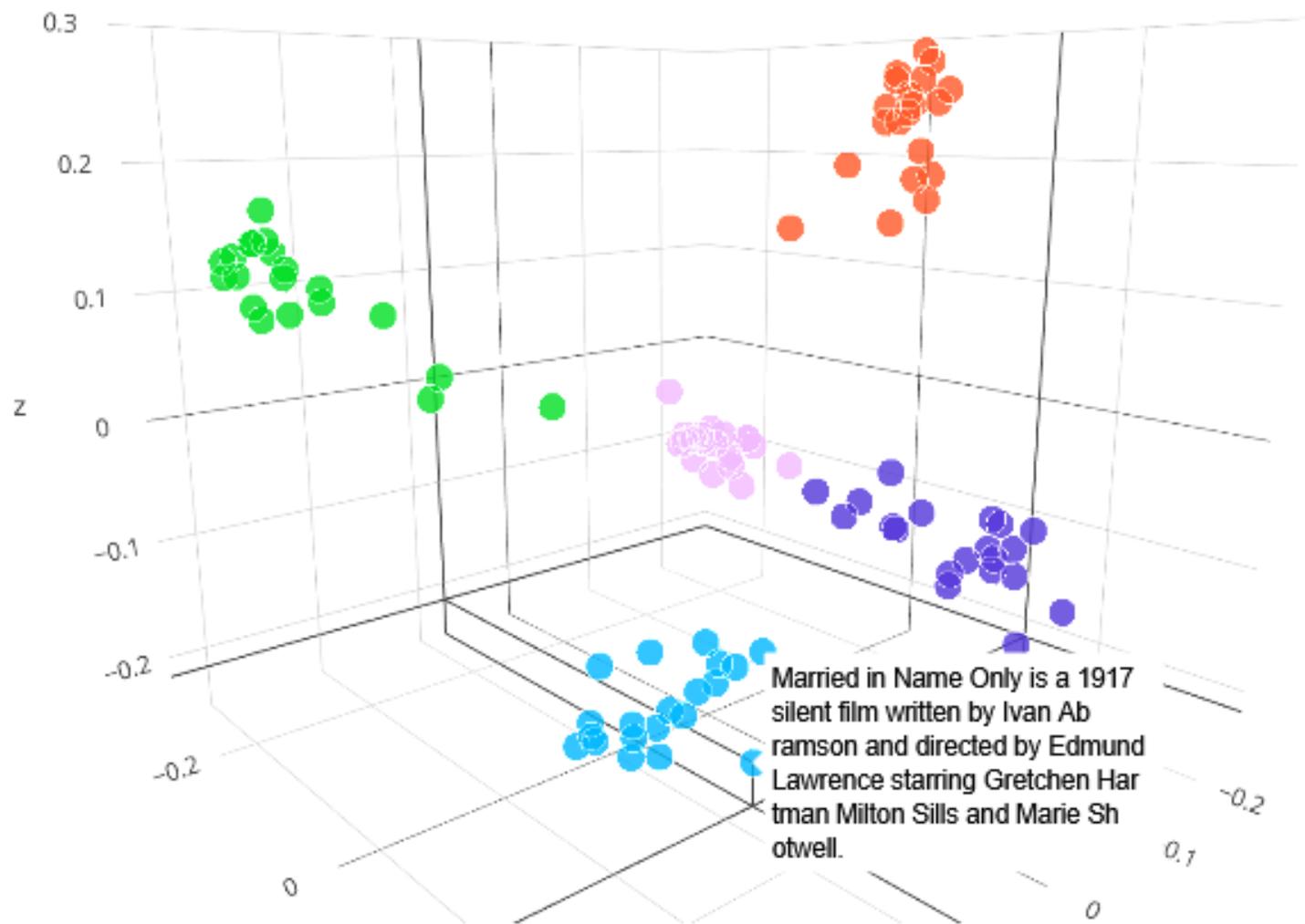
- Tokens are embedded into vectors
  - Mathematical representation of data
  - “king” -> [0.12, -0.98, 0.55, ..., 1.06]
  - Semantic similarity
  - “apple” and “orange” are closer than “apple” and “car”
  - GPT-3 had 12288 dimensions
  - Each dimension is a “feature”
-



<https://openai.com/index/introducing-text-and-code-embeddings/>

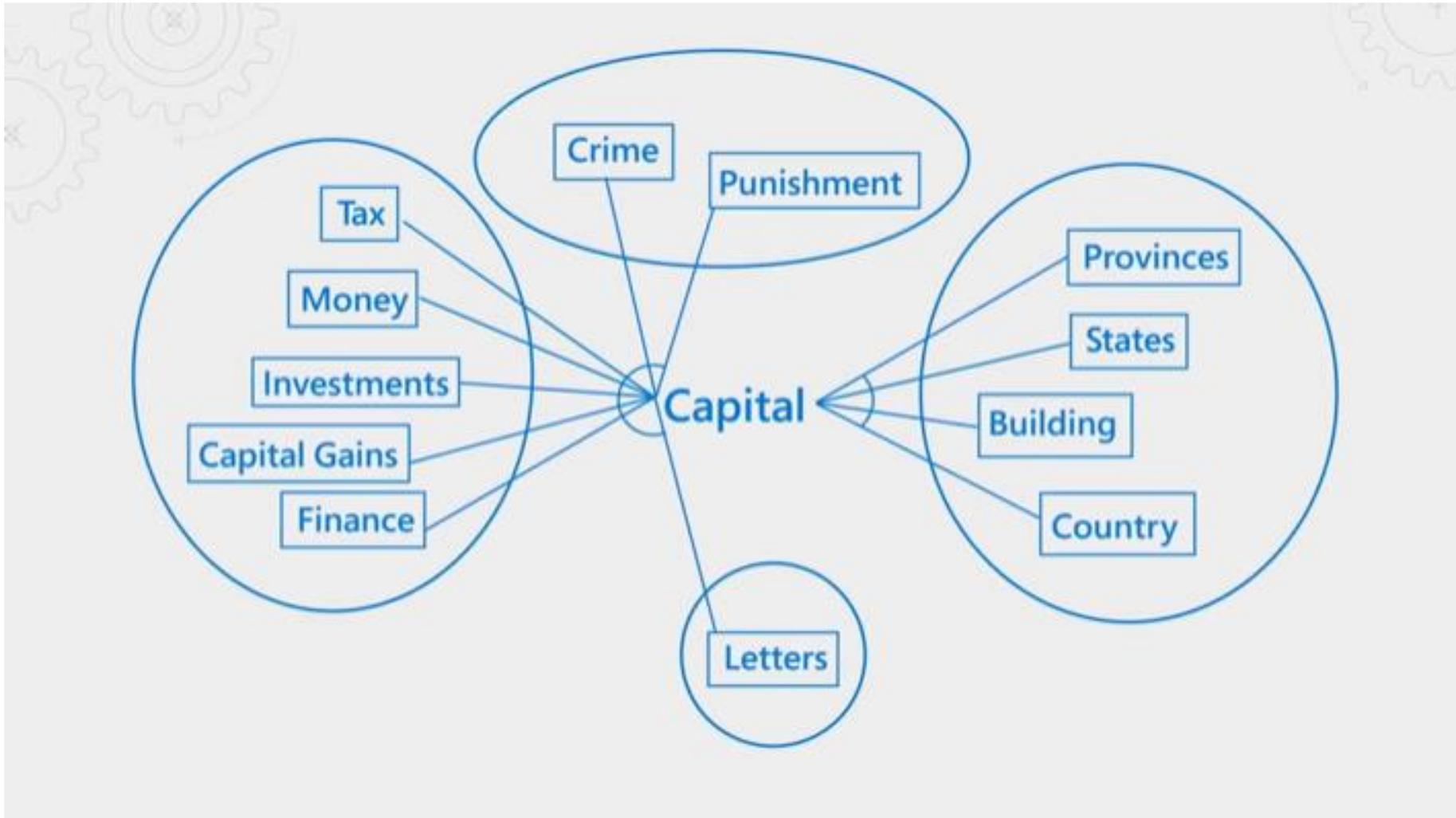
Drag to pan, scroll or pinch to zoom

● animal ● athlete ● film ● transportation ● village



# Attention

- Attention is All You Need, Google, 2017
- “Attention block” – does two primary things
  1. Allows vectors to modify each other based on all the other vectors
  2. Allows for parallel calculation
- “What is the capital of France?”



<https://learn.microsoft.com/en-us/azure/search/semantic-search-overview>

# Training

- 1. Pretraining. Give it a massive dataset, gets embedded, finds relationship. Get a **base model**
- 2. Fine-tuning. Manually done by humans. Swap out to a different dataset by humans. Get lots of people to write prompts, as well as an ideal response. Get an **assistant model**
- This makes up the model

# Example: Ollama 2:70b

- **Step 1 – pretraining**

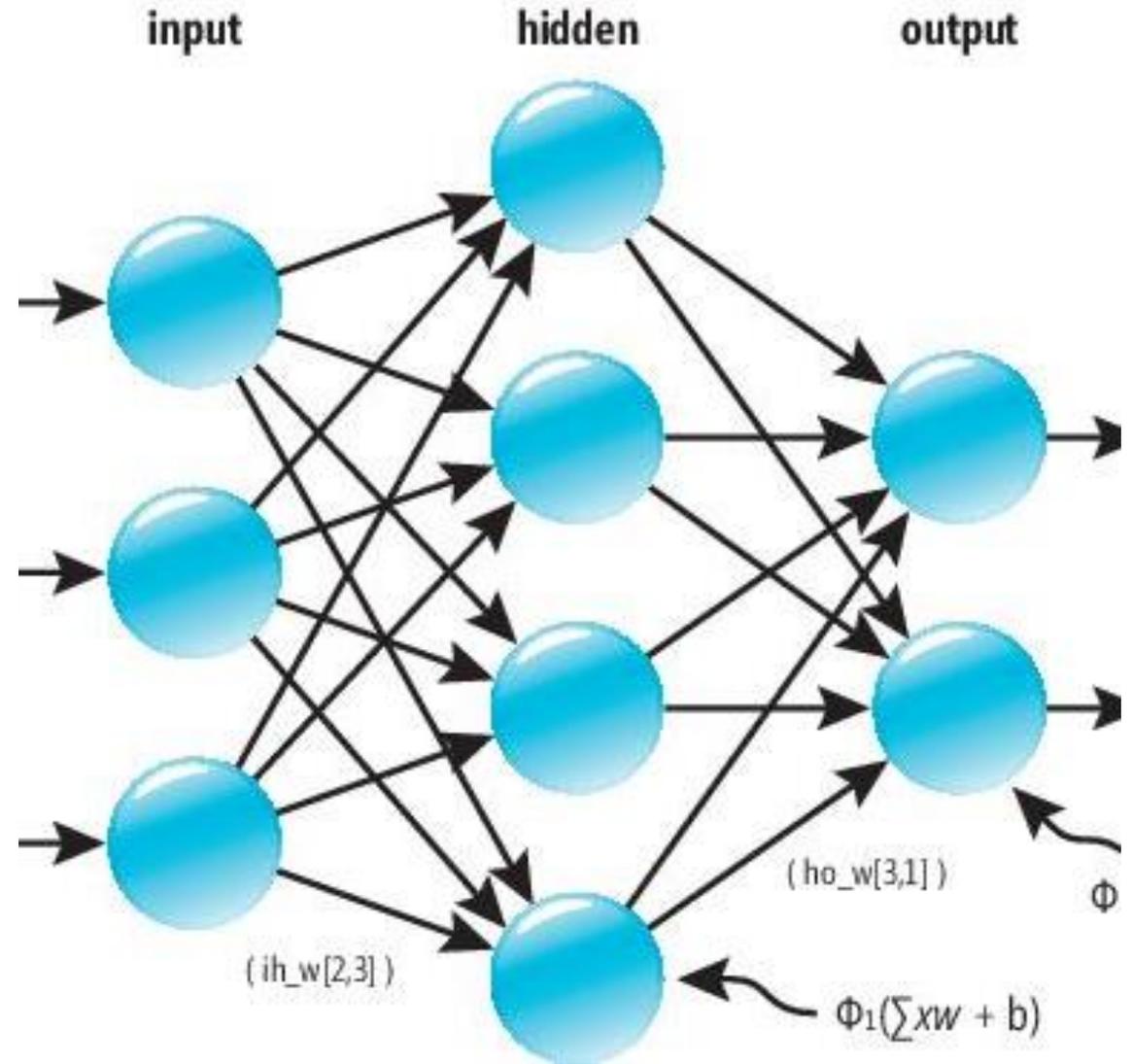
1. Get ~10TB of text from the internet
2. Use ~6,000 GPUs ~\$2M, and ~12 days to embed this into vectors in a neural network with rough semantic meaning

- **Step 2 – finetuning**

1. Write labeling instructions for humans
2. Hire lots of humans to write ~100k high quality idea Q&A
3. Use these to “test” the model and refine it

# Parameters

- GPT-4 (supposedly) has 1.7 trillion
- As the model trains, these change
- Temperature for randomness



# Example: ChatGPT

Brackets hold the flow,  
LINQ dances with sharp logic—  
Types bloom, strong and bold.

Write me a haiku about C#

Write me a haiku about C#

System: You are a helpful assistant willing to help users with a variety of tasks.

User: Write me a haiku about C#

System:

# Example: ChatGPT

- Br
- Brackets
- Brackets hold
- Brackets hold the
- etc.



```
Brackets hold the flow,  
LINQ dances with sharp logic—  
Types bloom, strong and bold.
```

flow	0.8
world	0.05
code	0.02
...	...
dog	0

# Retrieval Augmented Generation (RAG)

- Augment the response of an LLM with info retrieved from a data store
- Two main phases
  1. Retrieval. User prompts something, system retrieves info from an external knowledge store
  2. Generation. Retrieved info is used to augment user's prompt. AI model processes both user prompt and retrieved info to produce an enriched response

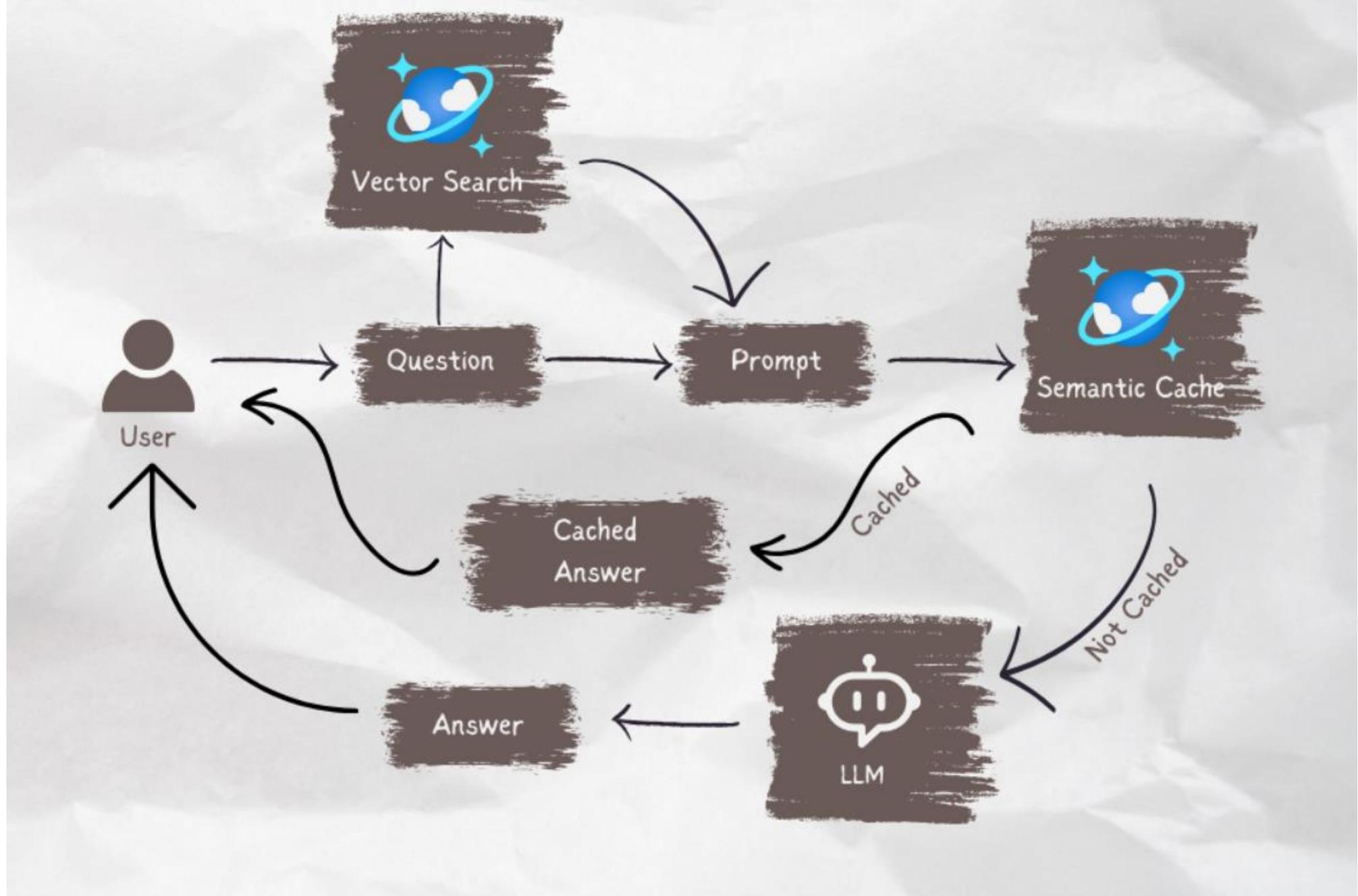
# RAG benefits

---

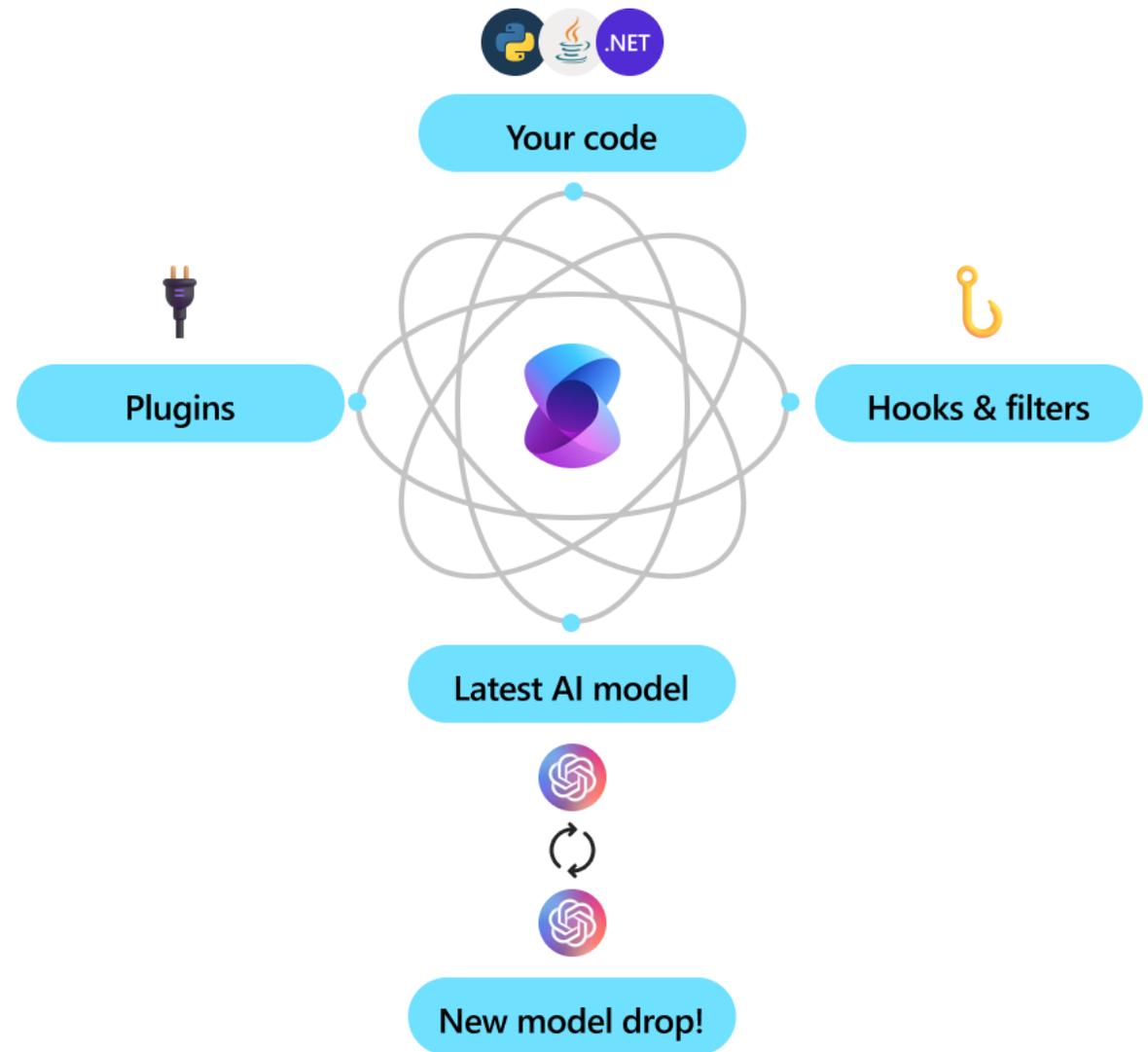
- Improved accuracy
- Reduced hallucinations
- Retrieve up-to-date info
- Get domain-specific knowledge

# Agents

- “Agentic”
- AI that doesn’t just generate, but takes action on behalf of the user
- Fulfill wide variety of tasks with some/minimal human intervention
- 3 parts
  1. LLM – model used
  2. State – context, helps guide decisions based on past results
  3. Tools – like the functions before with the convos, a bridge between the model and external systems



# Semantic Kernel

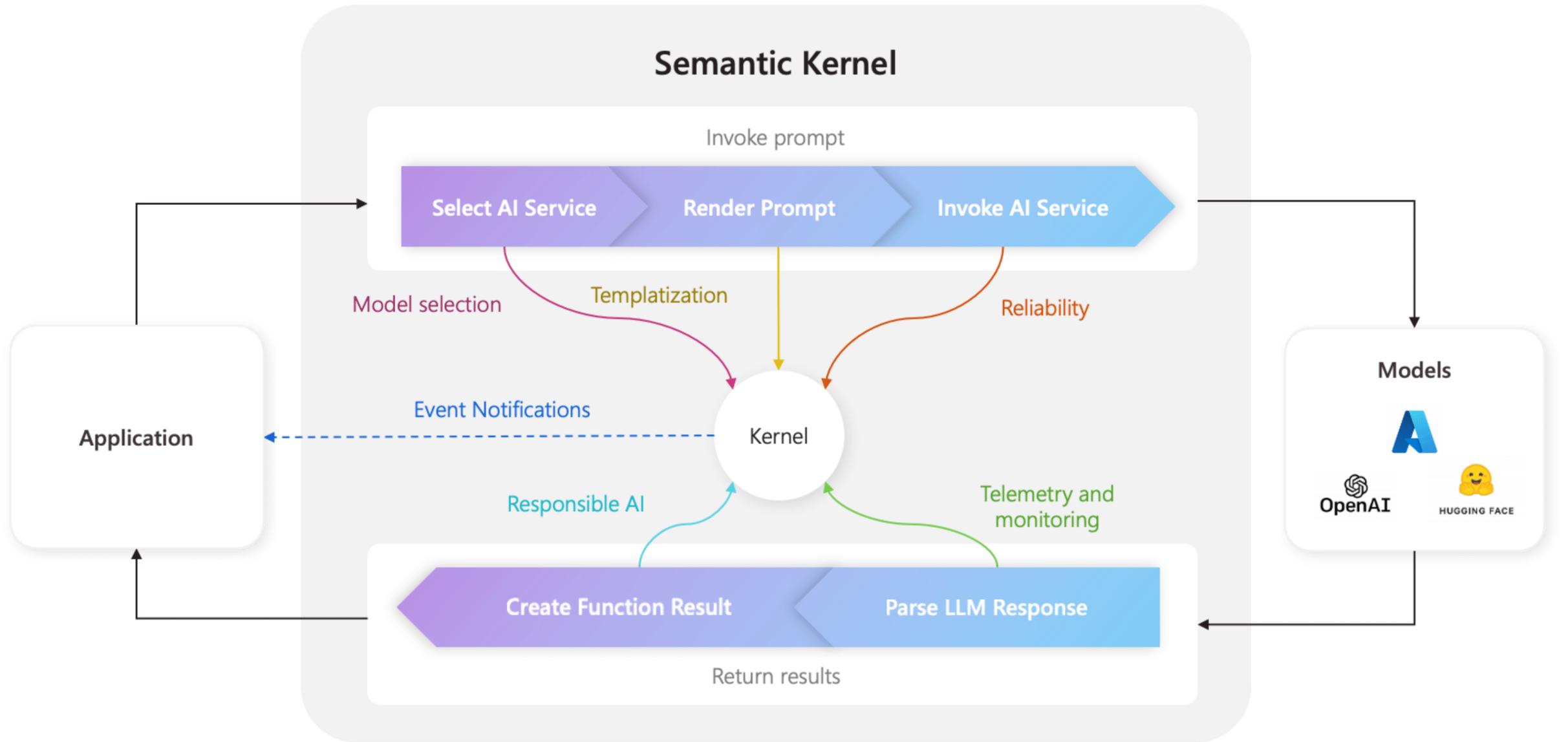


# Semantic Kernel

- SDK for building AI agents and integrating AI models into your apps
- Microsoft, open-source
- C#, Java, Python
- Bridge between normal AI usage and your code, which helps simplify process of developing AI-powered apps

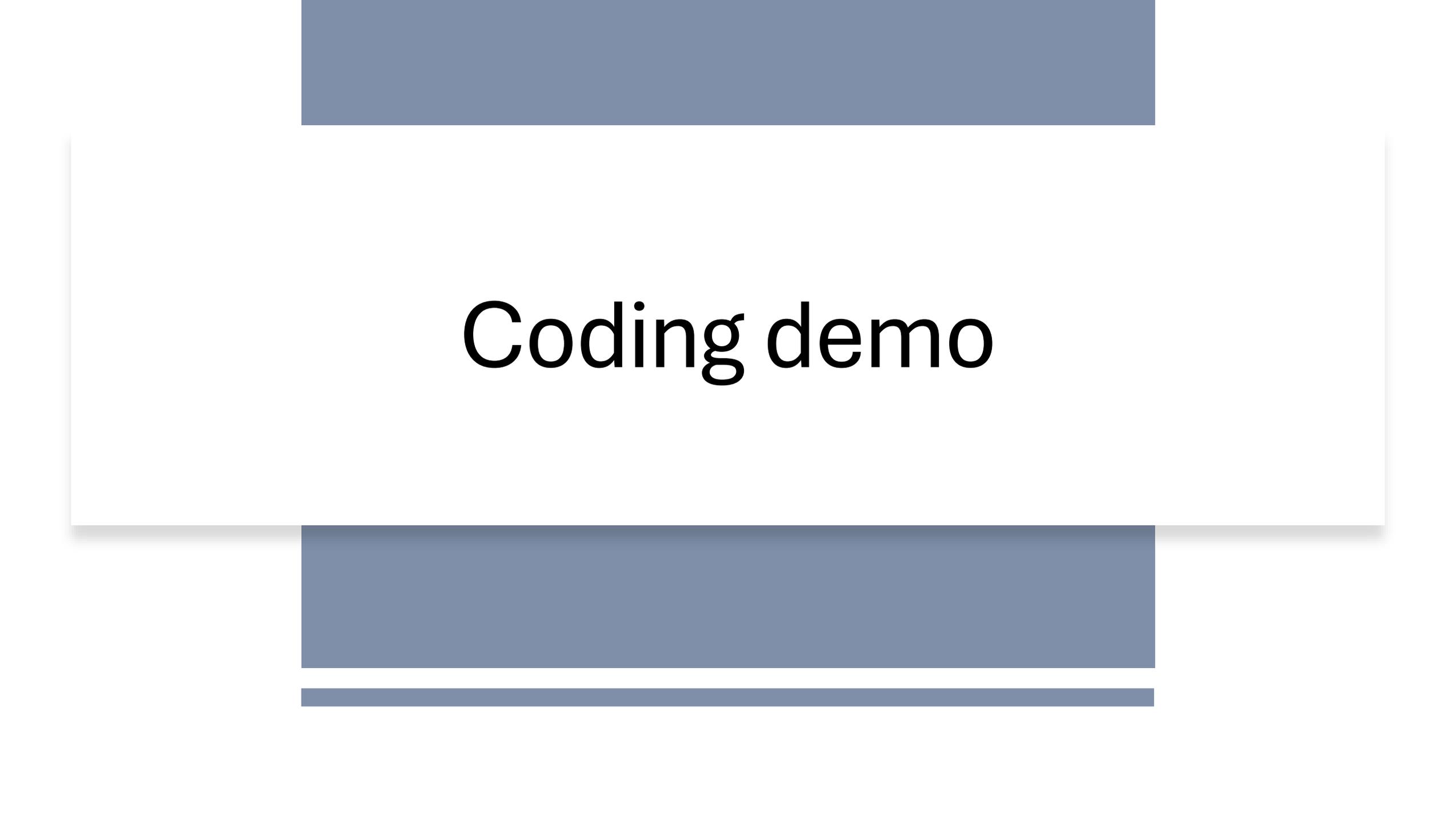
# What makes up Semantic Kernel?

- Kernel itself is basically just a DI container. Central location to configure/monitor agents
- Two types of components
  1. Services. AI services like chat completion, but also logging, HTTP clients, etc. Modeled after existing service provider in .NET
  2. Plugins. Components used by your AI services to perform work, like retrieving data from a database or calling an external API



# Why Semantic Kernel?

- Devs normally have to learn the different APIs/quirks of each service
- SK abstracts and unifies APIs
- Must support function calling
- Bridge between your code and LLM



# Coding demo

Function calling example with get\_weather function

javascript ↕

```
1 import { OpenAI } from "openai";
2
3 const openai = new OpenAI();
4
5 const tools = [{
6   "type": "function",
7   "name": "get_weather",
8   "description": "Get current temperature for a given location.",
9   "parameters": {
10    "type": "object",
11    "properties": {
12      "location": {
13        "type": "string",
14        "description": "City and country e.g. Bogotá, Colombia"
15      }
16    },
17    "required": [
18      "location"
19    ],
20    "additionalProperties": false
21  }
22 }];
23
24 const response = await openai.responses.create({
25   model: "gpt-4.1",
26   input: [{ role: "user", content: "What is the weather like in Paris today?" }],
27   tools,
28 });
29
30 console.log(response.output);
```



# Resources

- [3Blue1Brown Neural networks playlist](#)
  - [Andrej Karpathy Intro to Large Language Models](#)
  - [Microsoft Github samples - Getting Started](#)
  - [Semantic Kernel Microsoft docs](#)
  - [Semantic Kernel course on Microsoft Learn](#)
-

Thank you!



Daniel Ward

