

Essential AI for Leaders

(special edition abridged workshop for InnoLead members)

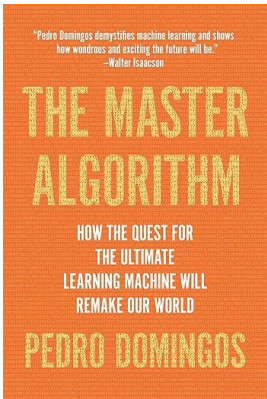
August 8th, 2024

Michael Hayes

PRACTICAL**AI**



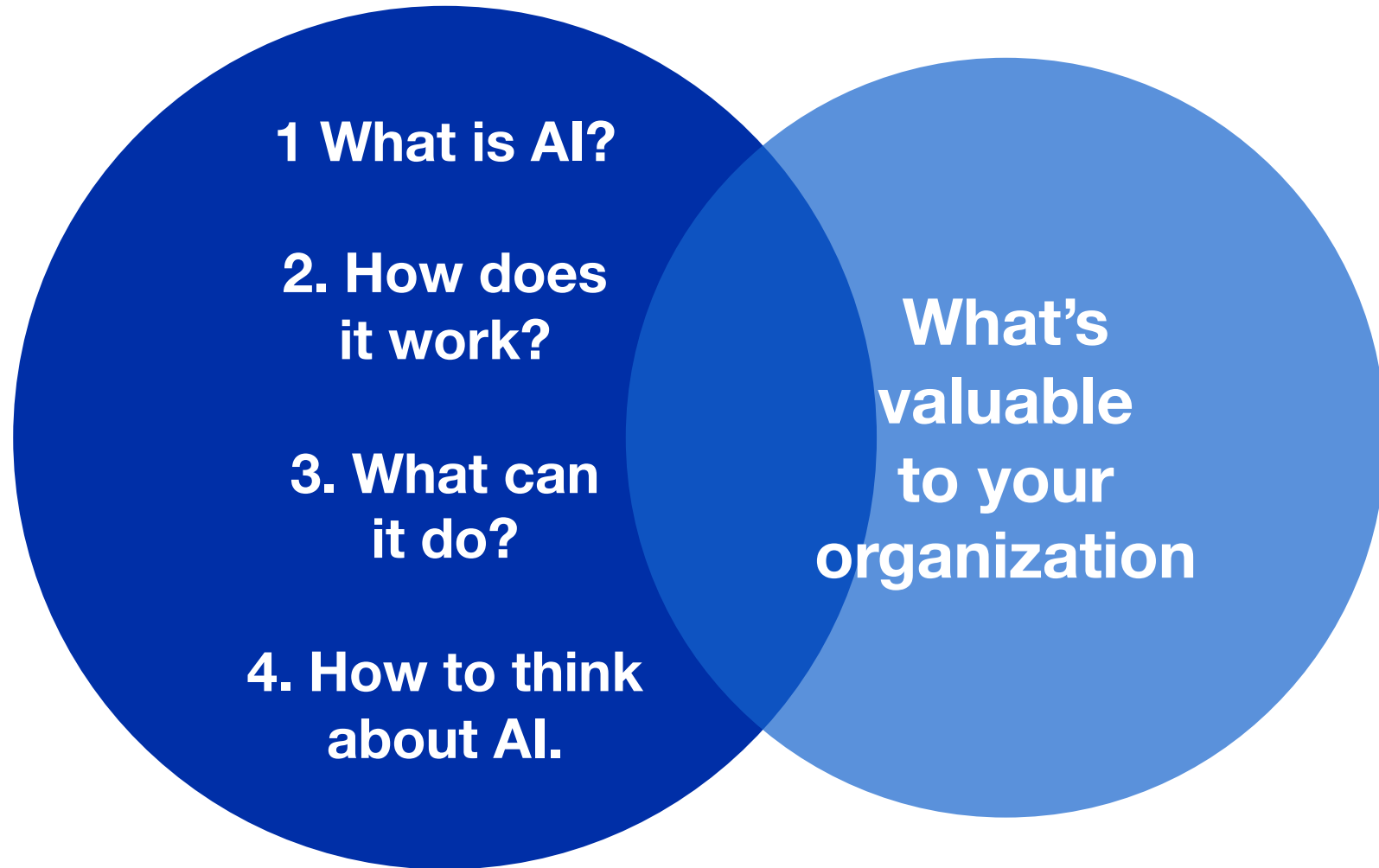
WIRED

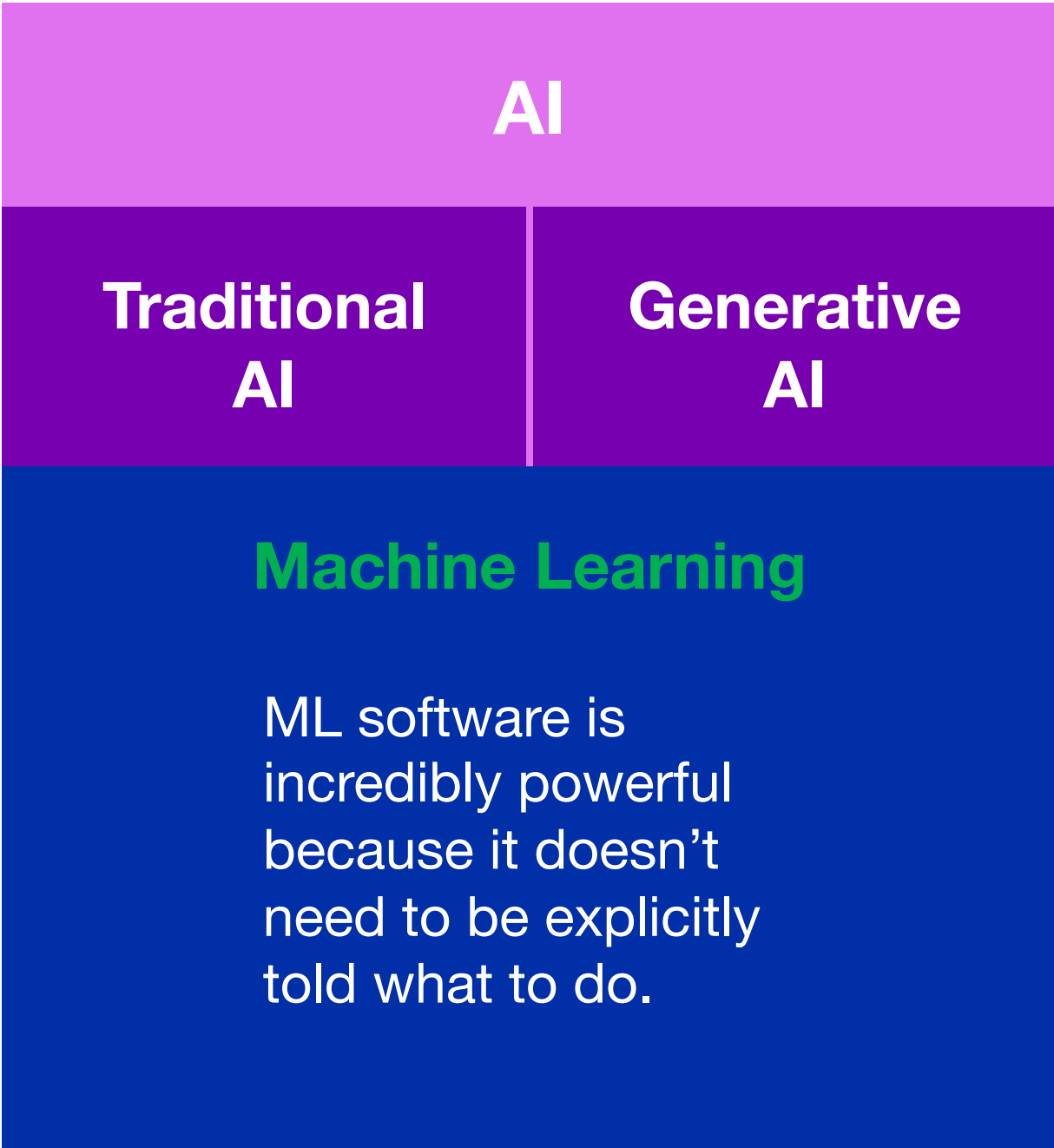


10 years in AI
20+ in B2B software
CS undergrad
MS MIT

1. Monitor AI developments in real time.
2. Cut through hype and curate for non-technical leaders.
3. Distill complex AI concepts down to easily understood essentials.
4. Create and conduct customized, in-person, 2-3 hour PracticalAI workshops.

What can AI do for my organization?





Practically speaking,
AI today = ML

How does ML software learn?

ML software learns sort of like we do:

- By example
- By trial and error

How does ML software learn?

ML software learns sort of like we do:

- **By example**
 - You don't teach a child what a dog is by saying "if it has four legs and a tail and two eyes and fur or hair and a wet nose, then it's a dog."
 - You point at dogs and pictures of dogs and say "dog."

Software learning from example data

	A	B	C	D	E	F
1	DRINK DATA					
2	Glass Stem	Bubbles	Color	Temp	Alcohol %	Drink (label)
3	Y	Y	Yellow	57 F	5.4	Beer
4	Y	Y	Yellow	47 F	12.2	Champagne
5	N	N	Red	59 F	10	Wine
6	N	Y	Amber	53 F	4.2	Beer
7	Y	Y	Brown	54 F	6.2	Beer
8	N	Y	Yellow	57 F	5.6	Beer
9	Y	N	Yellow	57 F	12.4	Wine
10	Y	N	Red	65 F	9.4	Wine
11	Y	N	Red	67 F	10.5	Wine
12	Y	Y	Yellow	45 F	11.4	Champagne
13	N	Y	Yellow	57 F	5.4	Beer
14	Y	Y	Yellow	47 F	12.5	Champagne

Software learning from example data

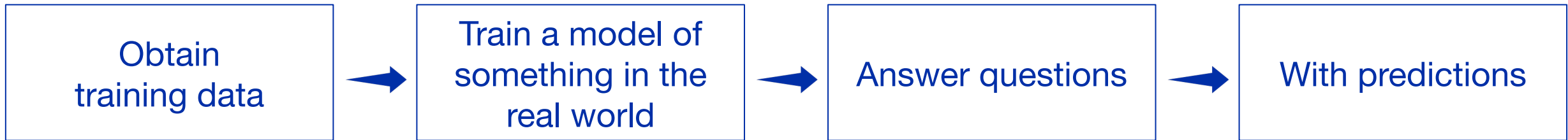
“supervised learning”

	A	B	C	D	E	F
1	DRINK DATA					
2	Glass Stem	Bubbles	Color	Temp	Alcohol %	Drink (label)
3	Y	Y	Yellow	57 F	5.4	Beer
4	Y	Y	Yellow	47 F	12.2	Champagne
5	N	N	Red	59 F	10	Wine
6	N	Y	Amber	53 F	4.2	Beer
7	Y	Y	Brown	54 F	6.2	Beer
8	N	Y	Yellow	57 F	5.6	Beer
9	Y	N	Yellow	57 F	12.4	Wine
10	Y	N	Red	65 F	9.4	Wine
11	Y	N	Red	67 F	10.5	Wine
12	Y	Y	Yellow	45 F	11.4	Champagne
13	N	Y	Yellow	57 F	5.4	Beer
14	Y	Y	Yellow	47 F	12.5	Champagne

“labeled
training
data”



Supervised learning paradigm



Supervised learning is in many ways the core of AI today. It can be incredibly useful in many business settings.

AI

Traditional
AI

Machine Learning

Supervised
Learning

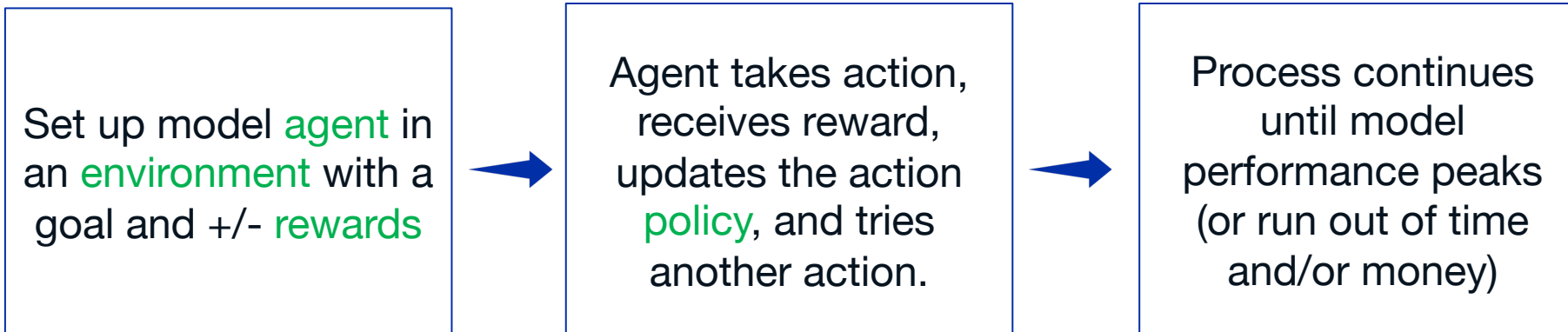
How does ML software learn?

ML software learns sort of like we do:

- By example (supervised learning)
- **By trial and error**
 - If you want to teach a dog a new trick, you don't explicitly tell the dog what to do. You reward it with treats when it performs the desired behavior.



Learning by trial and error – reinforcement learning





Learning by trial and error – reinforcement learning

Examples

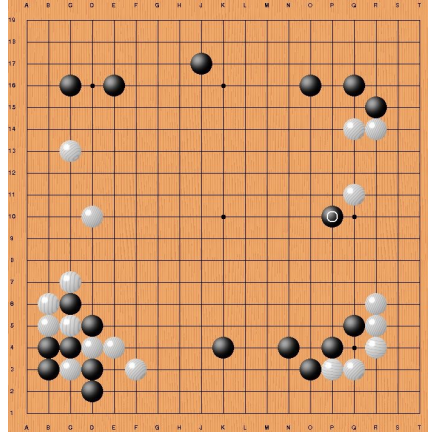
Car learning to park; https://www.youtube.com/watch?v=VMp6pq6_QjI

Robot dog learning to walk: <https://www.youtube.com/watch?v=xAXvfVTgqr0>



Three very cool things about reinforcement learning

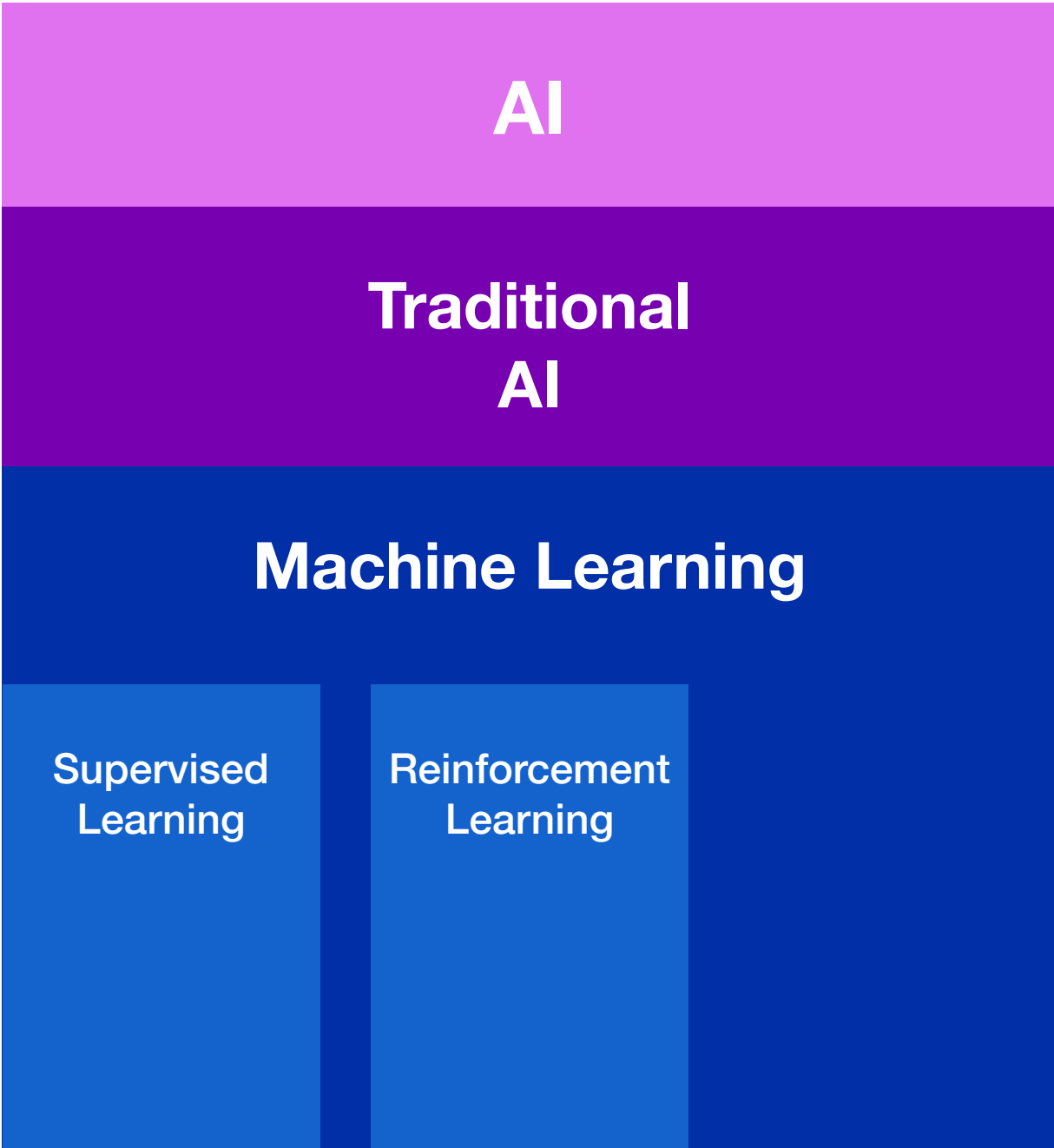
1. Doesn't need training data.
2. Doesn't learn from human examples, so it is not limited by what the best humans know. It can achieve superhuman performance.
3. Starts from scratch, so it doesn't hesitate to try things that humans "know not to try." It ends up "discovering" things.

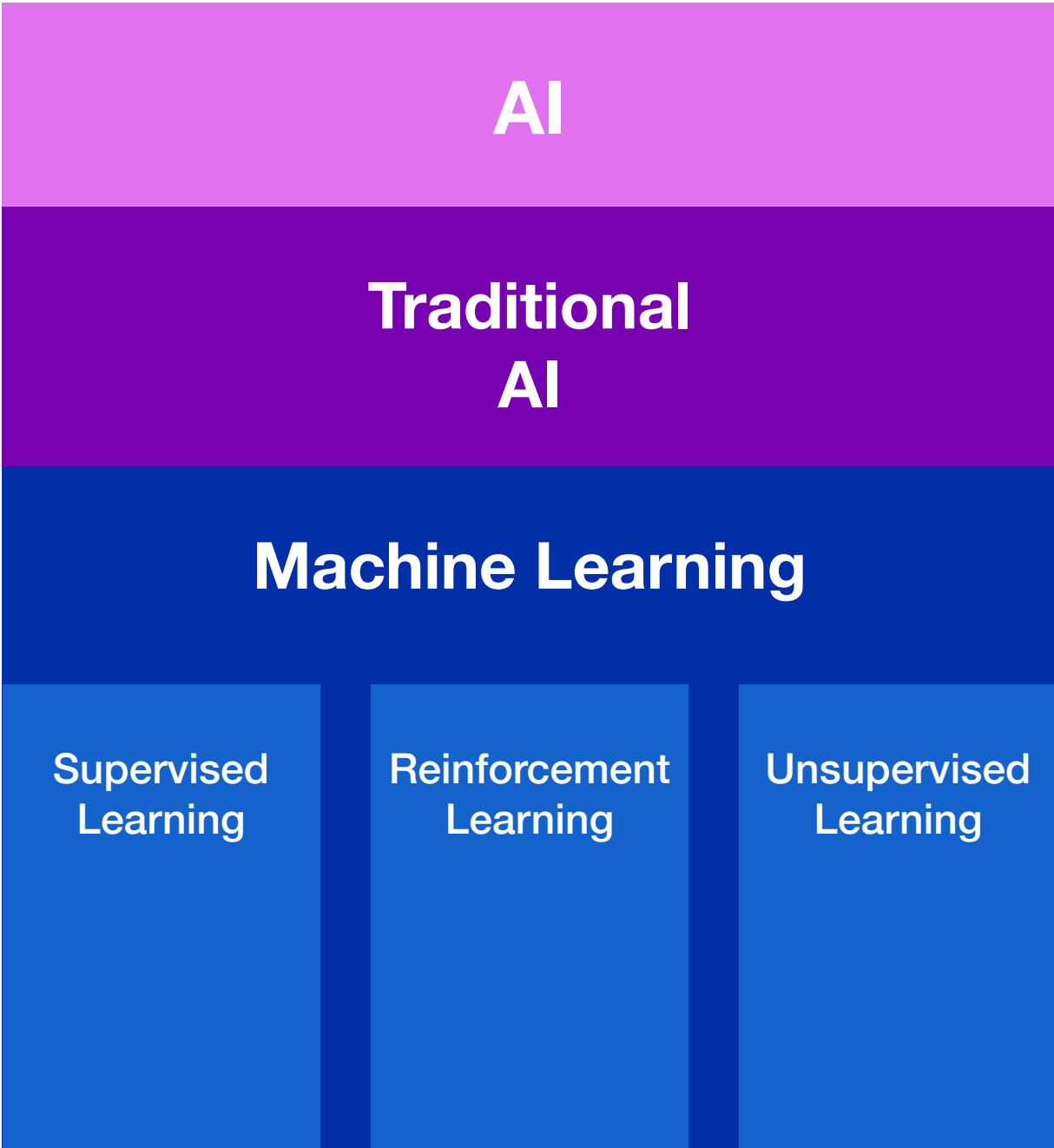


Move 37, or how AI can change the world.

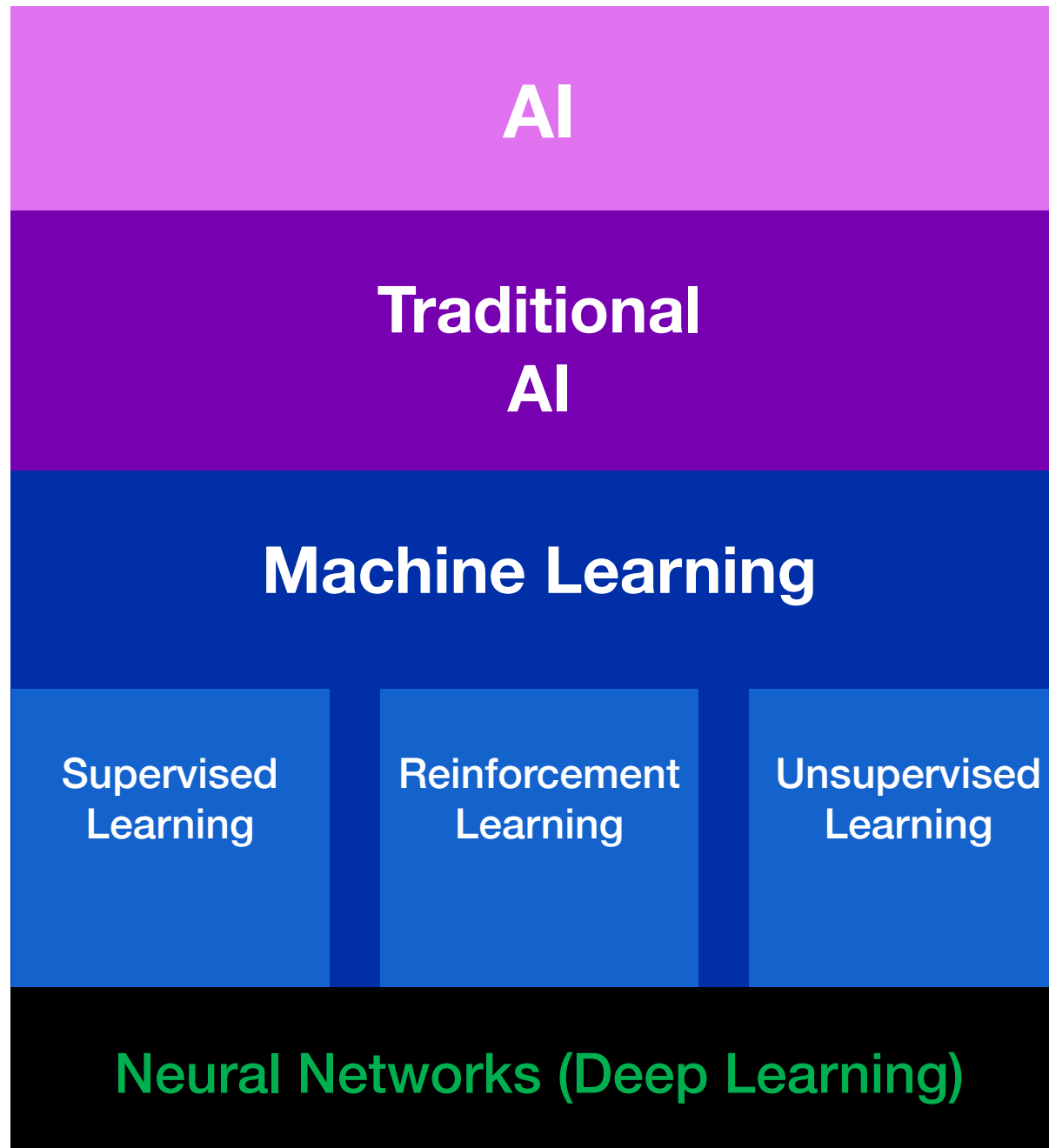


DeepMind AI Reduces Google Data Centre Cooling Bill by 40%



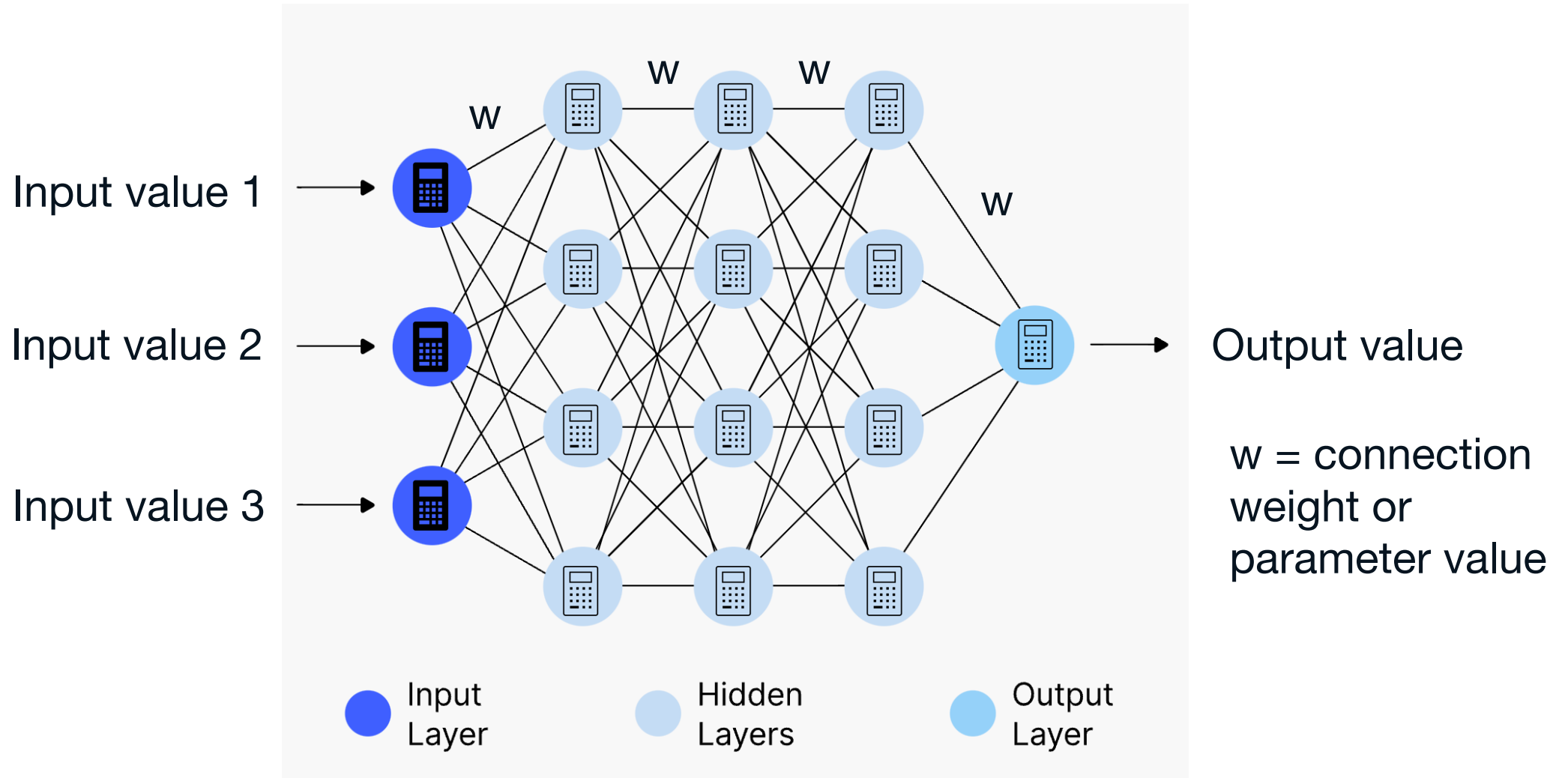


Practically speaking,
AI today = ML



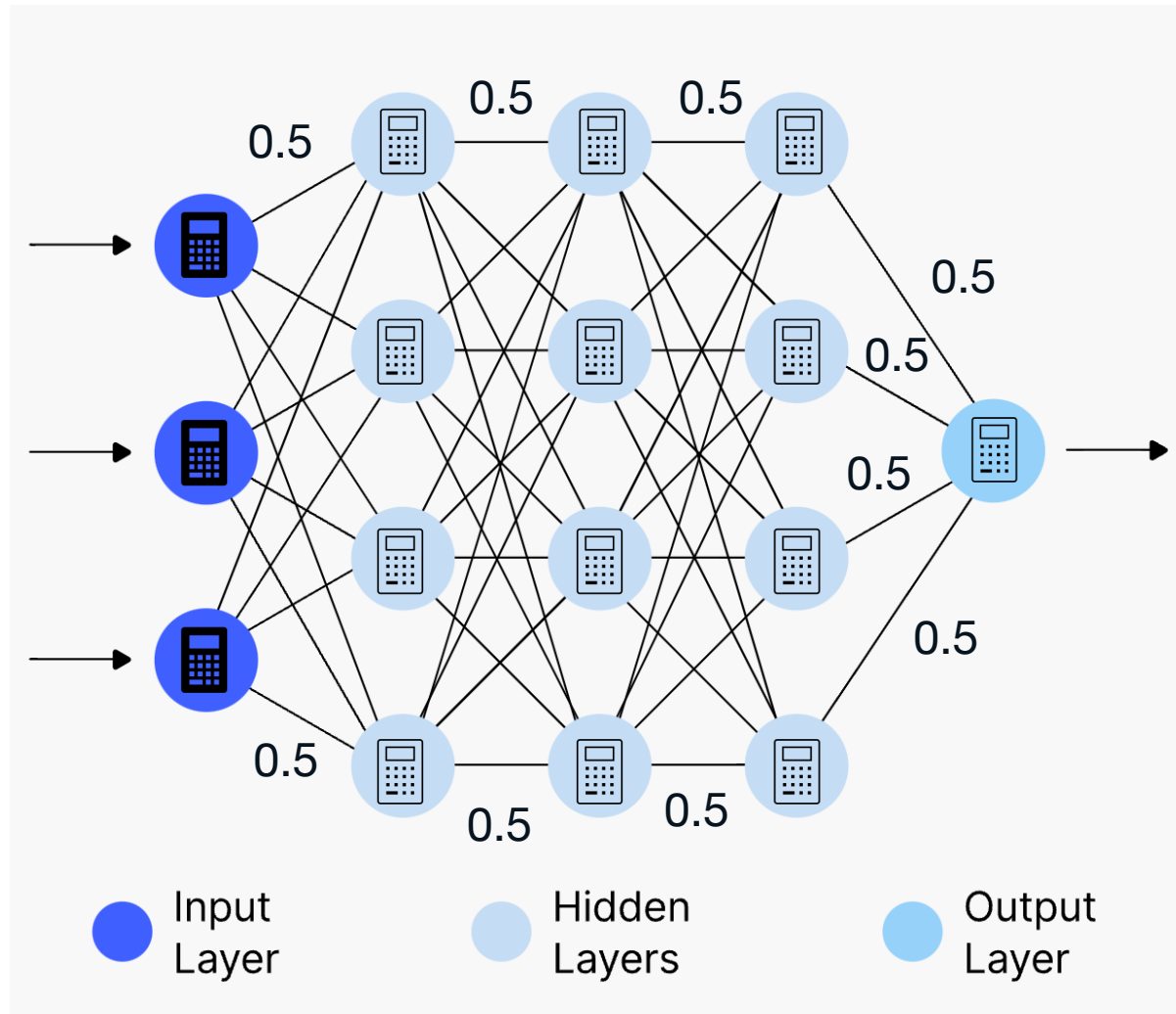
Practically speaking,
AI today = ML = Deep Learning with Neural Networks

Deep Learning with Neural Networks



1. Design network architecture and initialize parameters

INPUT



OUTPUT

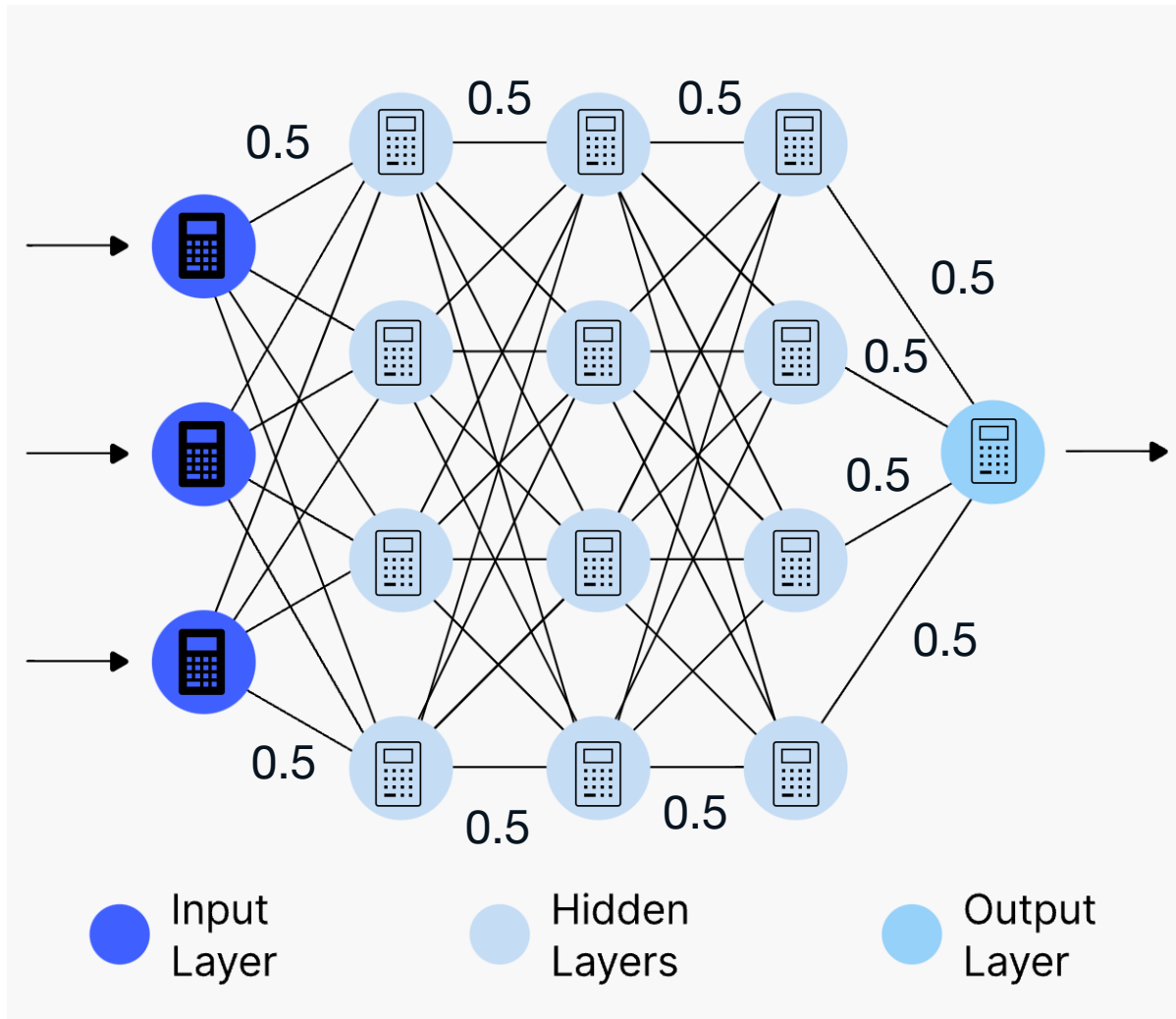
2. Train with input and output pairs (labeled training data)

INPUT

Square feet

Bedrooms

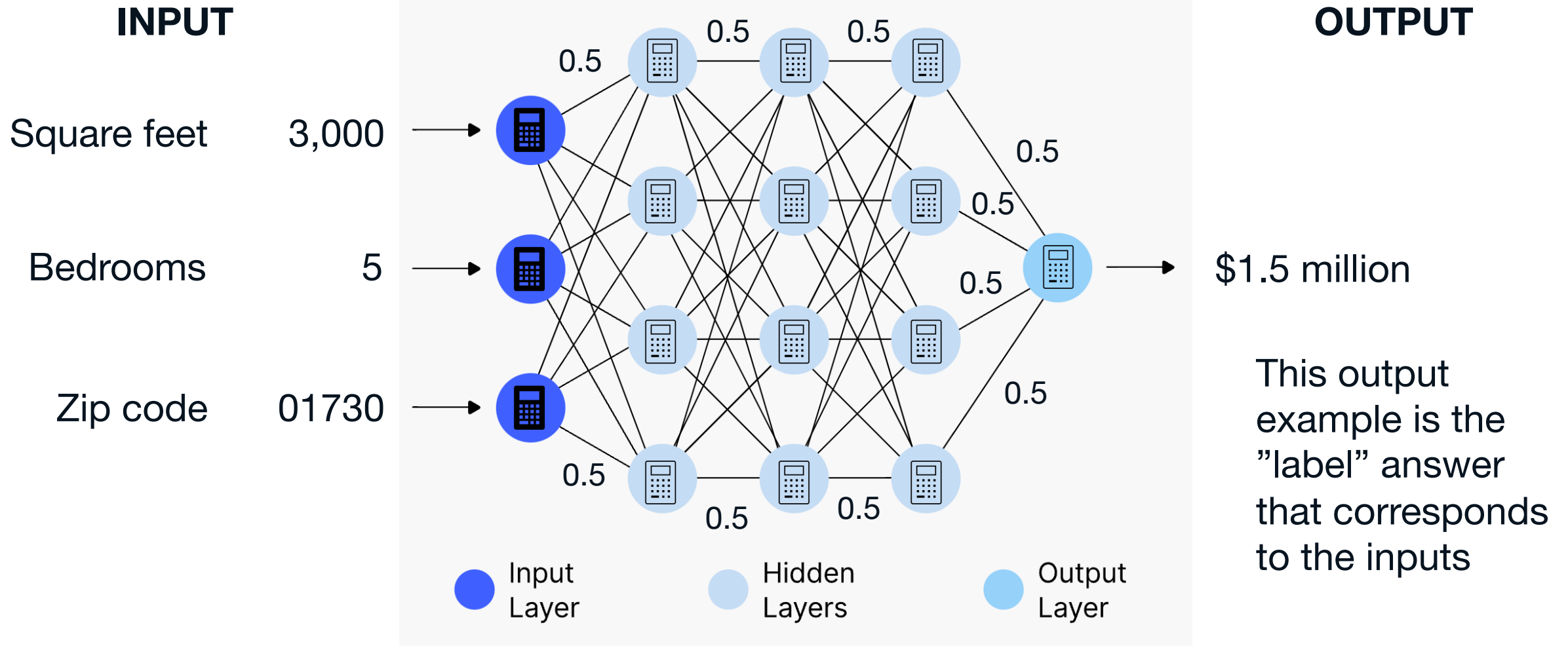
Zip code



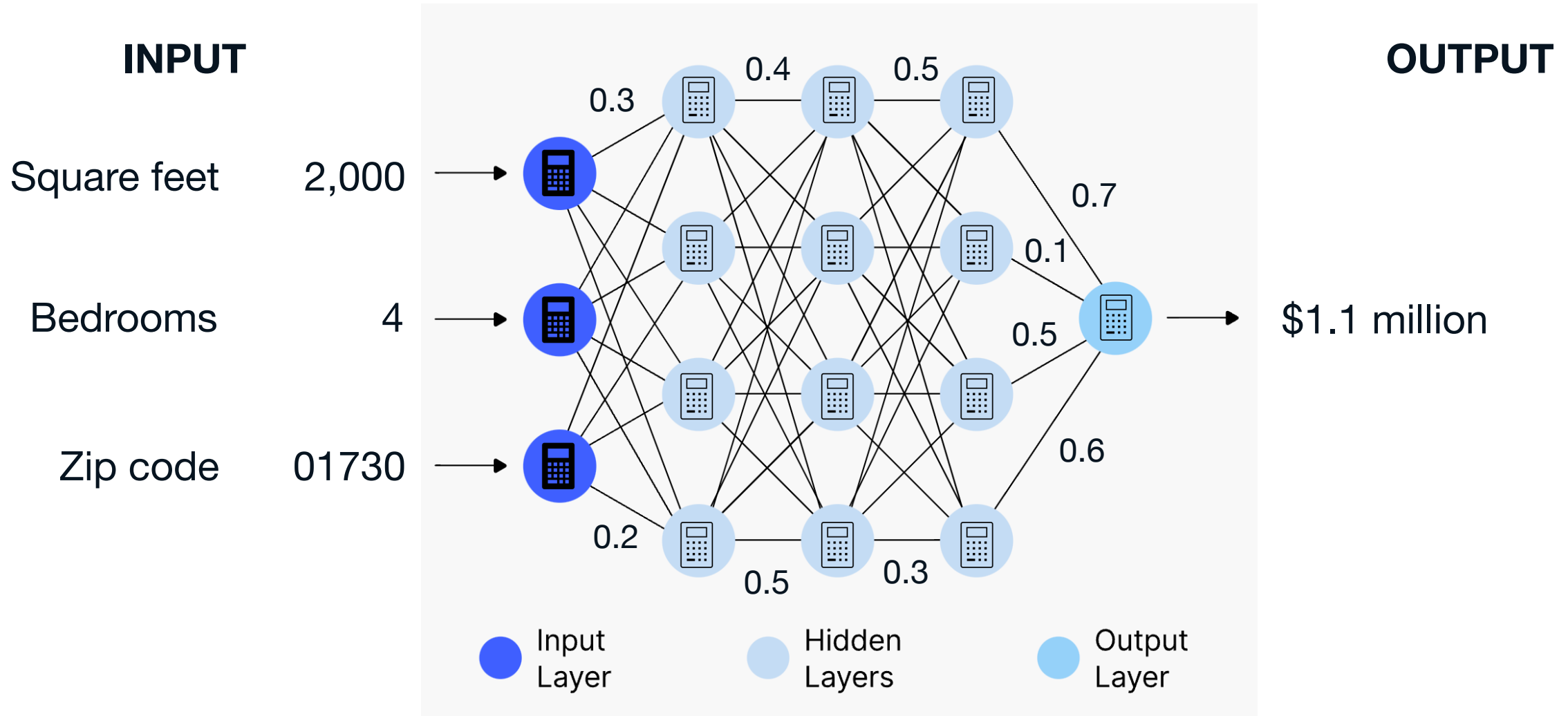
OUTPUT

House Price

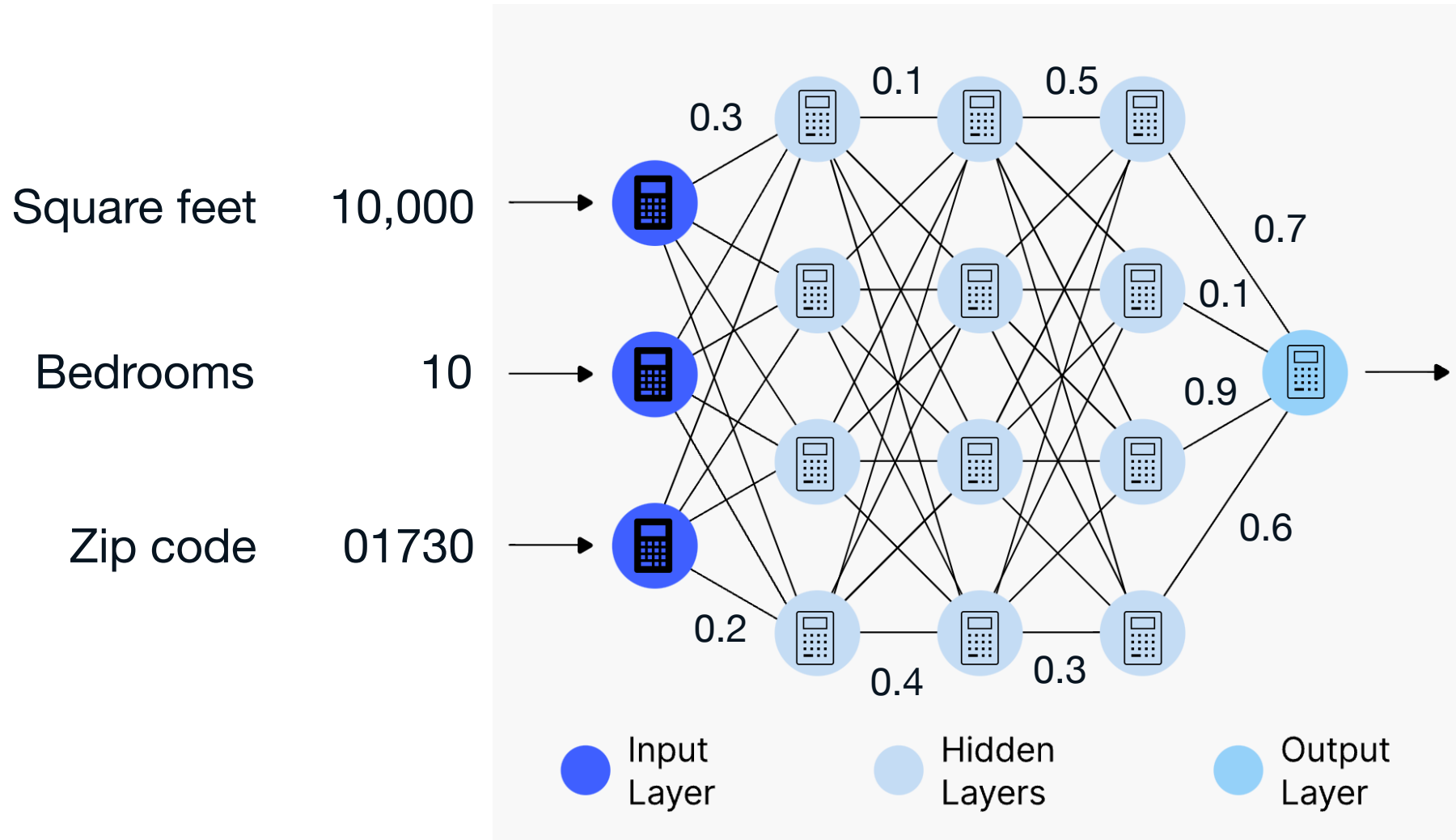
2. Train with input and output pairs



2. Every example pair adjusts parameters



3. When training complete, parameters are set

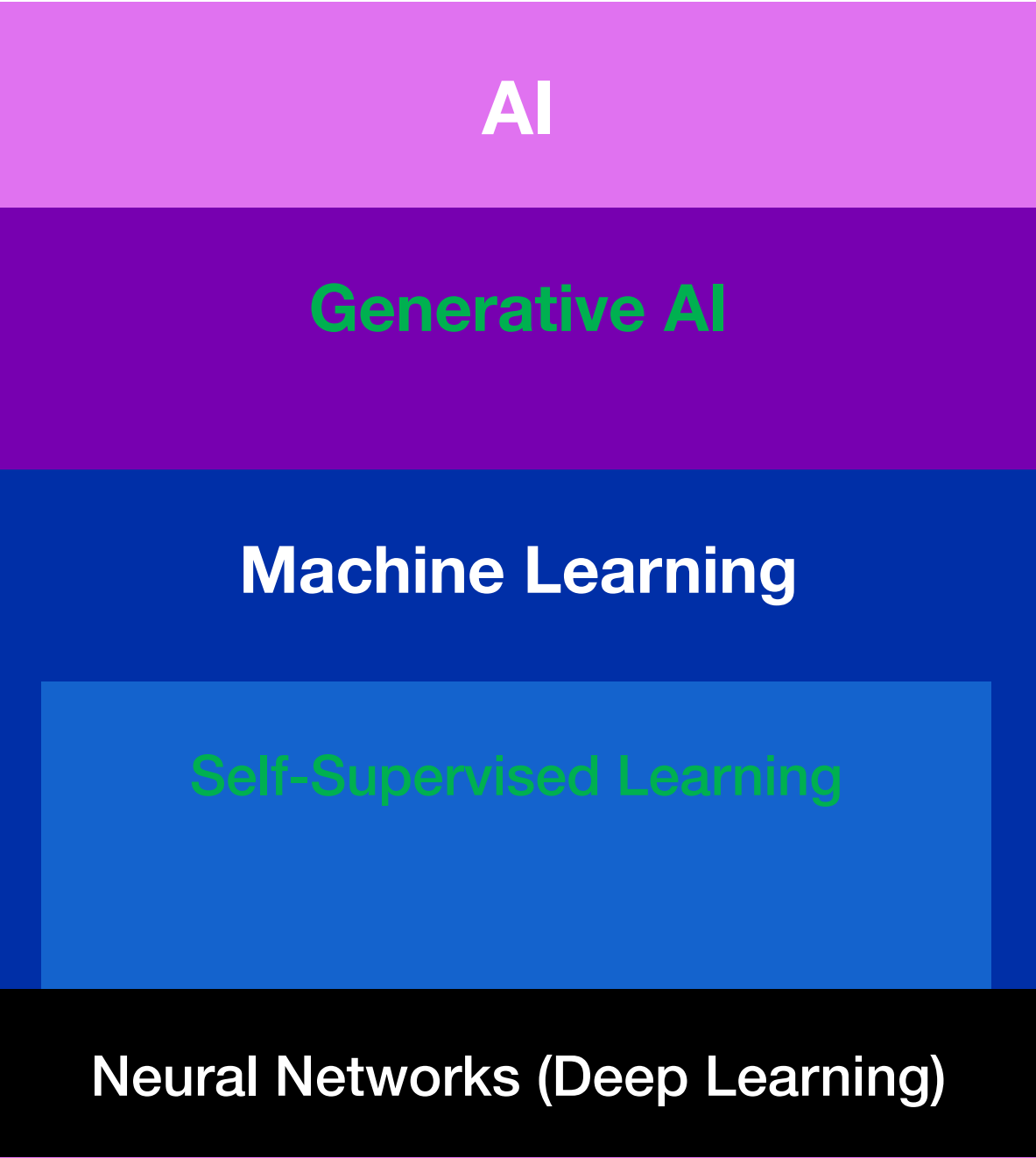


Given a new input the model has never seen, it will predict a price.

\$5.1 million

95% certain

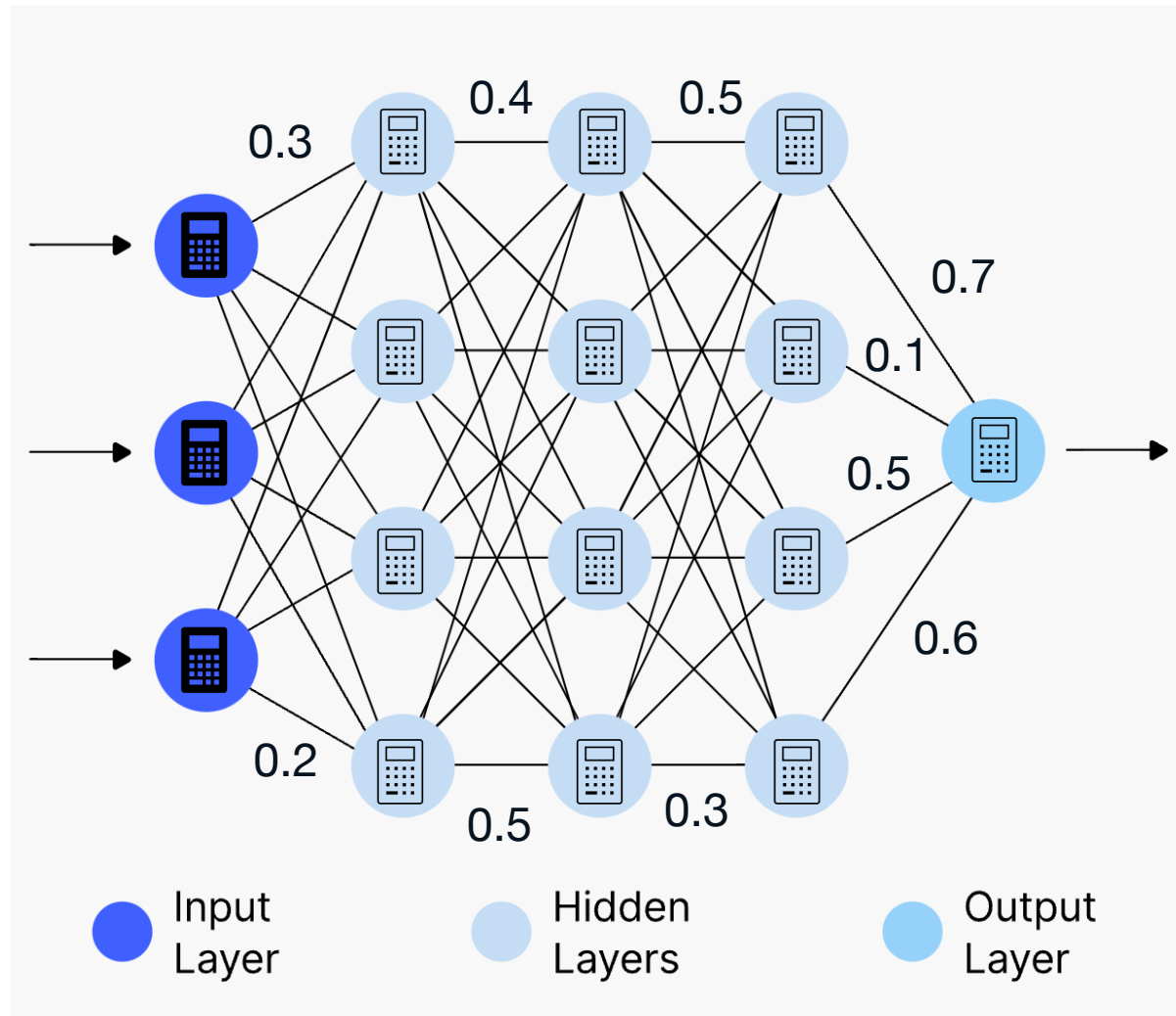
Generative AI



Still true:
practically
speaking,
**AI today = ML =
Deep Learning
with Neural
Networks**

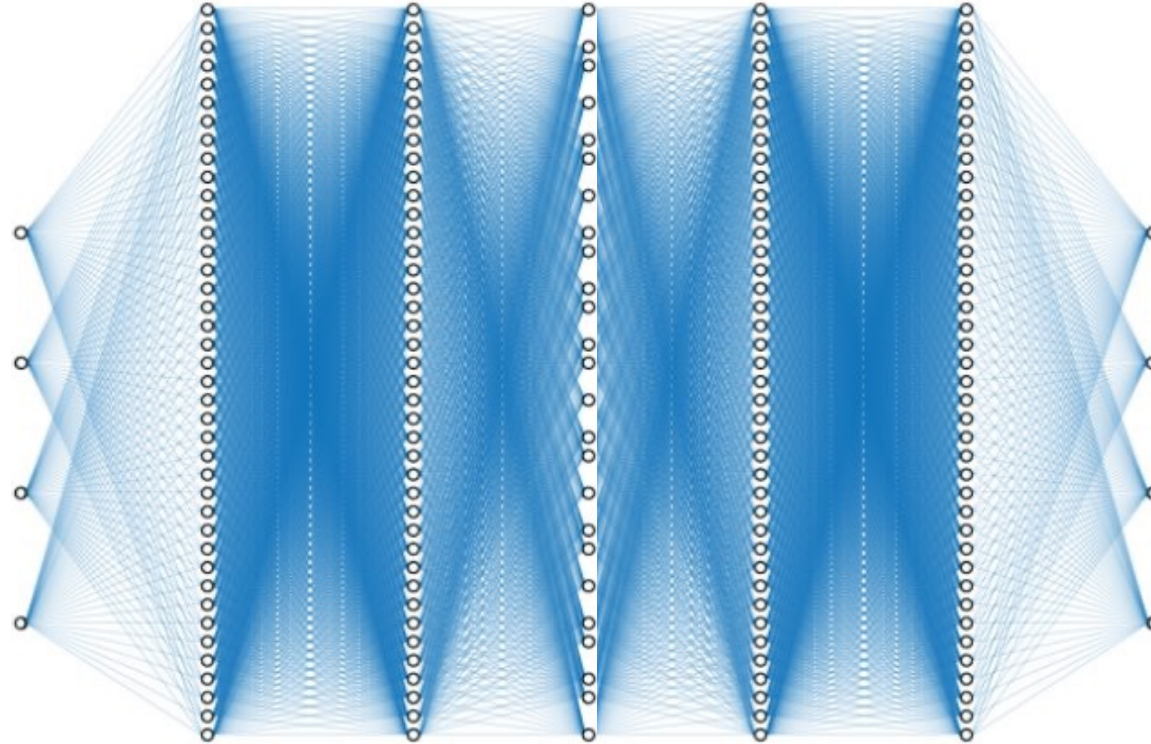
GPT models are **very** large

This model has four layers and 92 parameters



GPT models are **very** large

The largest **open source** GPT model has 120 layers and 405 **billion** parameters



It can capture very complex relationships between inputs and outputs

The network architecture is called a transformer, which is the T in GPT

Let's build a Large Language Model (LLM)

Stage I: Pre-training (the P in GPT)

1. Download all the text on the internet (~100TB e.g. 100K GB). This is the training data.
2. Buy or rent thousands of GPUs.
3. Use **Self-Supervised Learning (SSL)** to train a **model** of the text internet.
 - costs ~\$10M-\$100M, takes weeks-months
4. Obtain **Base Model**

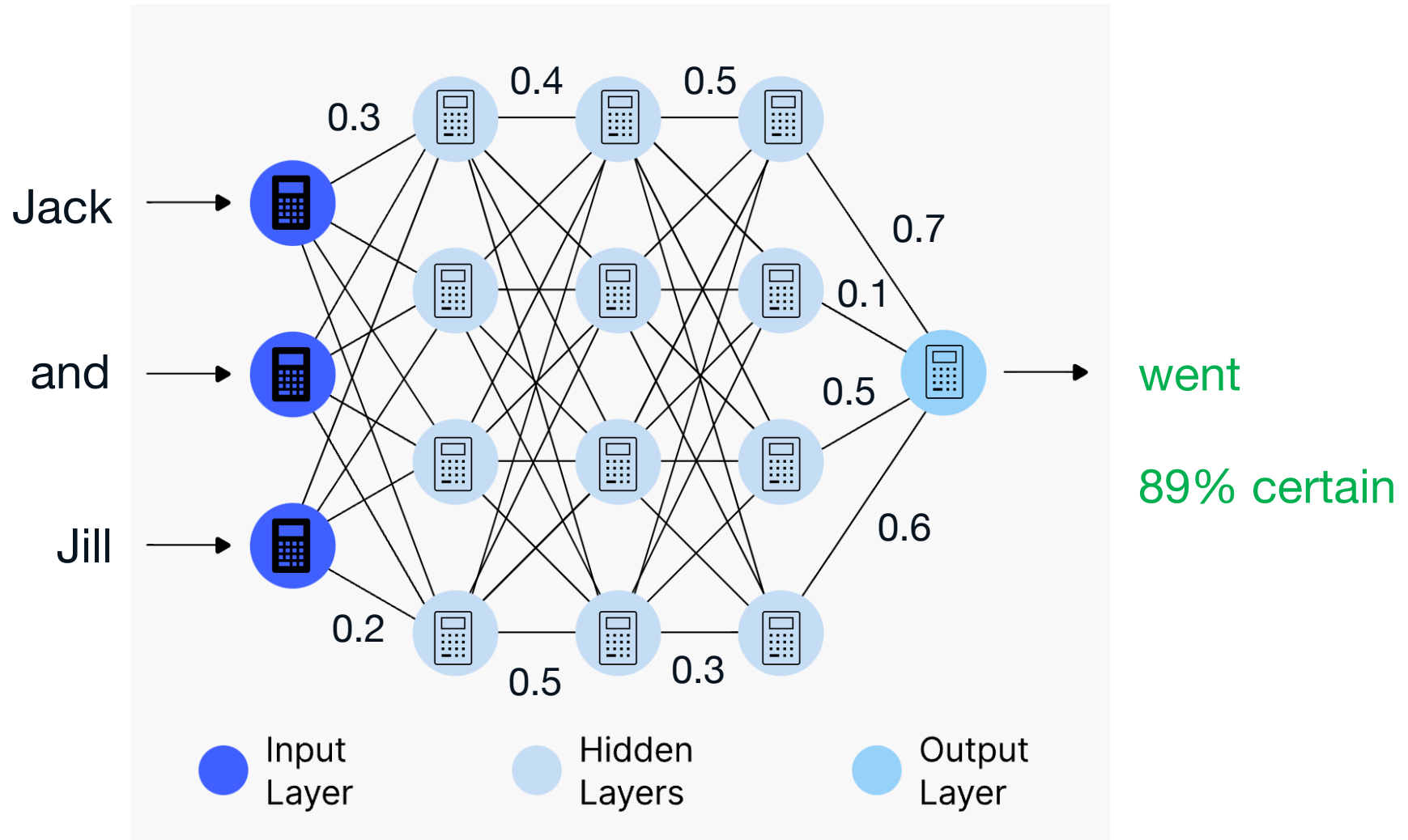
Model = a mathematical representation of underlying patterns in data.

LLM uses SSL to learn relationships between words in strings of text, including importance of each word to nearby words.

Input	Answer*
My favorite	food
My favorite food	is
My favorite food is	pizza

*somewhat like supervised learning, but the “labels” are the next word in a string of words.

Predicts next word using previous words



Base models are amazing, but not super useful

Because there is so much information about the world captured in language, when the base model learns a representation of languages, it **indirectly learns an “understanding” of how much of the world works.**

Unfortunately, it doesn't know much about communicating with humans and answering questions.

It will just generate or “dream up” documents that look a lot like a combination of information from web sites that it has seen. Lots of “hallucinations” at this point.

It also learns a lot of incorrect information that's out there.

Let's build a Large Language Model (LLM)

Stage II: Fine Tuning

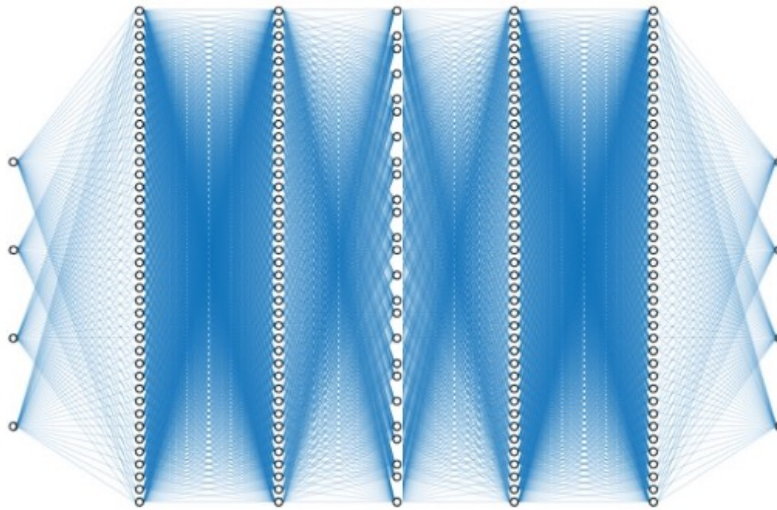
1. Decide how you want your LLM to communicate with humans. Helpful, polite, truthful etc. Elon Musk added a Snark Mode to his.
2. Hire people or outsource, collect ~100K (or as many as you can) **high quality examples of Q&A responses**.
3. Fine-tune base model on this data. Similar training as base model – predicting next word in a well-formed answer based on examples you gave it.
4. Insert **guardrails** so LLM will not explain how to kill someone, to make nerve gas, swear, make racist comments, etc.
5. Obtain **Assistant Model**.
6. Run a lot of tests.
7. Deploy your model.
8. Monitor, collect mistakes, go to Step 3

Let's build a Large Language Model (LLM)

Stage II: Fine Tuning

high quality Q&A response pairs

Explain how high-speed trading works on Wall Street



High-frequency trading (HFT) firms use powerful computers and algorithms to execute rapid trades, profiting from tiny price differences. They gain advantages through:

Colocation: Placing servers near the exchange for faster order execution.
Market Data Feeds: Accessing real-time market information.
Algorithms: Analyzing data and making split-second trading decisions.

HFT strategies include:

Market Making: Providing liquidity by simultaneously placing buy and sell orders.

Arbitrage: Exploiting price differences between markets or securities.

Statistical Arbitrage: Identifying patterns in stock prices to predict movements.

HFT impacts the market through:

Increased Liquidity: More continuous trading activity.

Reduced Trading Costs: Narrower bid-ask spreads due to competition.

Potential Volatility: Rapid trades can amplify price swings.

Fairness Concerns: Advantages for HFT firms raise questions of equity.

Open versus closed LLMs

Open (e.g. Llama)

1. Parameters and code necessary to run the model are available to download.
2. They come in different sizes, from around 1B parameters to 400B parameters.
3. LLMs up to around 5B can run on a very high-end laptop. There will be some latency in responses.
4. Some provide both the Base Model and the Assistant Model, so you can fine-tune a Base Model yourself.

Closed (e.g. ChatGPT 4o)

1. Access via web browser, computer / phone app, or API.
2. Incredibly easy to use - up and running in minutes.
3. Free versions are very good for many applications.
4. Paid versions usually score better than Open models on performance.
5. Enhanced functionality, such as ability to create a custom GPT and share it with others.



LMSYS Chatbot Arena Leaderboard

Rank* (UB)	Model	Arena Score	95% CI	Votes	Organization	License	Knowledge Cutoff
1	Gemini-1.5-Pro-Exp-0801	1299	+4/-5	15244	Google	Proprietary	2023/11
2	GPT-4o-2024-05-13	1286	+3/-4	72589	OpenAI	Proprietary	2023/10
3	GPT-4o-mini-2024-07-18	1277	+4/-5	16064	OpenAI	Proprietary	2023/10
3	Claude 3.5 Sonnet	1271	+3/-4	42939	Anthropic	Proprietary	2024/4
4	Gemini Advanced App (2024-05-14)	1266	+3/-3	52126	Google	Proprietary	Online
4	Meta-Llama-3.1-405b-Instruct	1264	+5/-4	13831	Meta	Llama 3.1 Community	2023/12

Multi-Modal (non-text) GPTs

Awesome Generative AI awesome

A curated list of modern Generative Artificial Intelligence projects and services.

<https://github.com/steven2358/awesome-generative-ai>

GPTs/LLMs have two Achilles heels

- **Reliability.** They are never, ever, 100% sure of the answer. They always answer with authority, even when they are making things up. Easy to fall into habit of trusting them.
 - You should put an invisible “I’m not quite sure that this is correct, but:” in front of everything they say.
 - Image generation and other non-text modalities are still full of errors and can require significant back and forth to get a good result.
- **Privacy.** For closed models, the data in your prompts are leaving the building. You can turn off some data logging with some LLMs, but the best solution is to run locally with an open model if you can.

Better GenAI - Prompts

There are four main areas to consider when writing an effective prompt:

Persona, **Task**, **Context**, **Format**

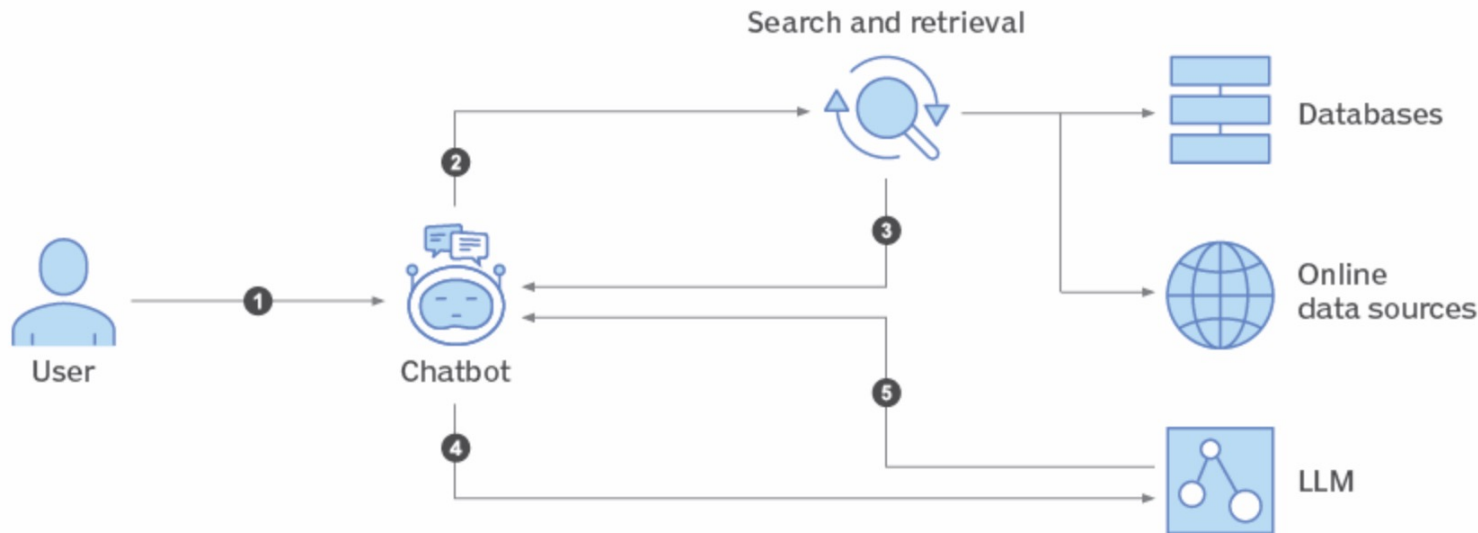
Example: “**You are a BCG Partner**. Read the attached. Draft an executive summary email to the CEO of a packaged goods company based on what you read. **Limit to bullet points.**”

Lots of info here <https://github.com/dair-ai/Prompt-Engineering-Guide>. Tips:

1. Use natural language. Write as if you're speaking to another person. Express complete thoughts in full sentences.
2. Be specific and iterate. Tell GPT what you need it to do (summarize, write, change the tone, create). Provide as much context as possible.
3. Be concise and avoid complexity. State your request in brief — but specific — language. Avoid jargon.
4. Make it a conversation. Fine-tune your prompts if the results don't meet your expectations or if you believe there's room for improvement. Use follow-up prompts and an iterative process of review and refinement to yield better results.

Better GenAI – Retrieval Augmented Generation

How an LLM using RAG works



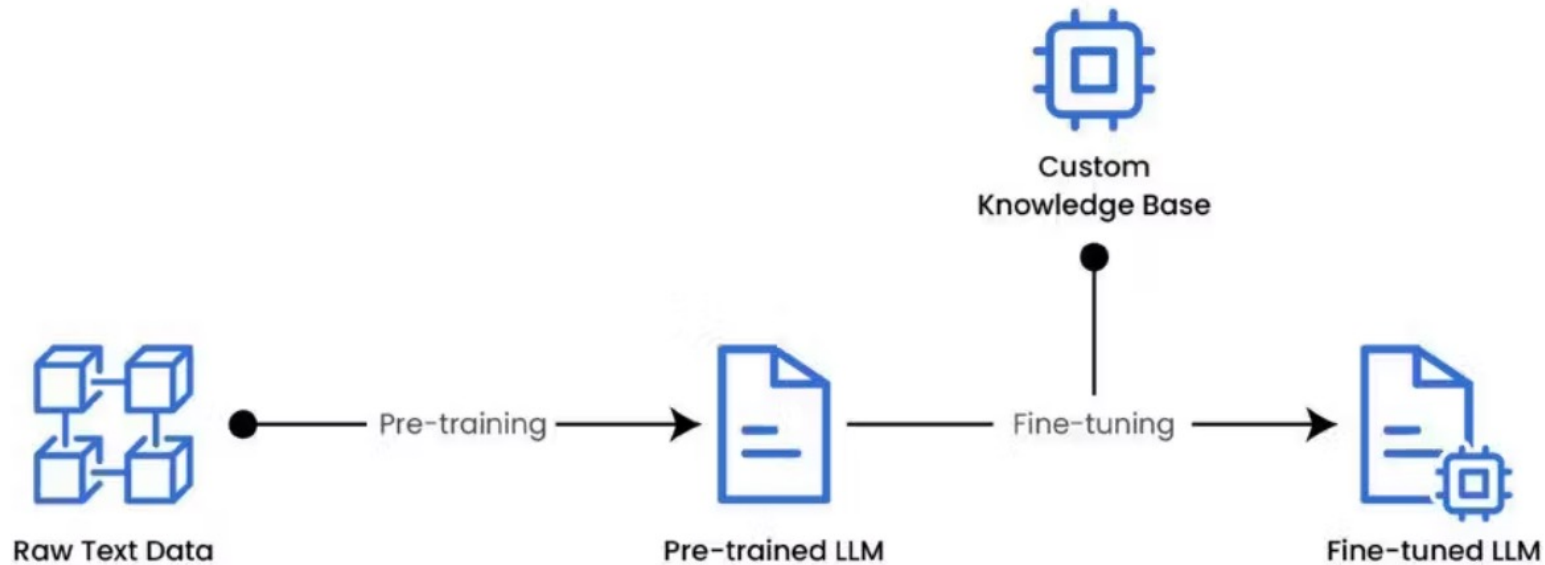
Examples:

Custom GPT for employees to “talk” to the employee handbook

Customer Support Chatbots: RAG can giving them access to a vast knowledge base of company information, product manuals, FAQs, and even previous customer interactions.

Personalized content: RAG can be used to create highly personalized content by analyzing a user's past interactions, preferences, and browsing history, the system can retrieve relevant articles, videos, or products from a large database

Better GenAI – Fine Tuning



Customization of both content and style of communication.

If run in house, maximizes privacy – e.g. financial or medical information.

Better GenAI – Agents

- [Auto-GPT](#) - An experimental open-source attempt to make GPT-4 fully autonomous.
- [babyagi](#) - An AI-powered task management system.
- [AgentGPT](#) - Assemble, configure, and deploy autonomous AI Agents in your browser.
- [GPT Engineer](#) - Specify what you want it to build, the AI asks for clarification, and then builds it.
- [MetaGPT](#) - The Multi-Agent Framework: Given one line requirement, return PRD, design, tasks, repo.
- [AutoGen](#) - AutoGen is a framework that enables the development of LLM applications using multiple agents that can converse with each other to solve tasks.
- [GPT Pilot](#) - Dev tool that writes scalable apps from scratch while the developer oversees the implementation.
- [Devin](#) - An autonomous AI software engineer by Cognition Labs.
- [OpenDevin](#) - An autonomous agent designed to navigate the complexities of software engineering. #opensource
- [Davika](#) - An agentic AI software engineer. #opensource

Moving forward with AI

1. Think in 3 dimensions.
2. Ask 4 questions.
3. Decide where on AI leadership spectrum you should be.

3 dimensions

1. Personal productivity with GenAI

- Everyone in your organization has free access to a
 - brainstorming partner
 - writing coach
 - research analyst
 - software engineer
 - executive assistant
- Requires developing new habit. Change your browser so there is always a GPT tool available to you. Set “can I use GPT?” reminders to pop up 2x per day.

2. Existing process improvement

- Break processes down into tasks – can any tasks be automated with AI?
- **Software development – your CTO should be all over this.**

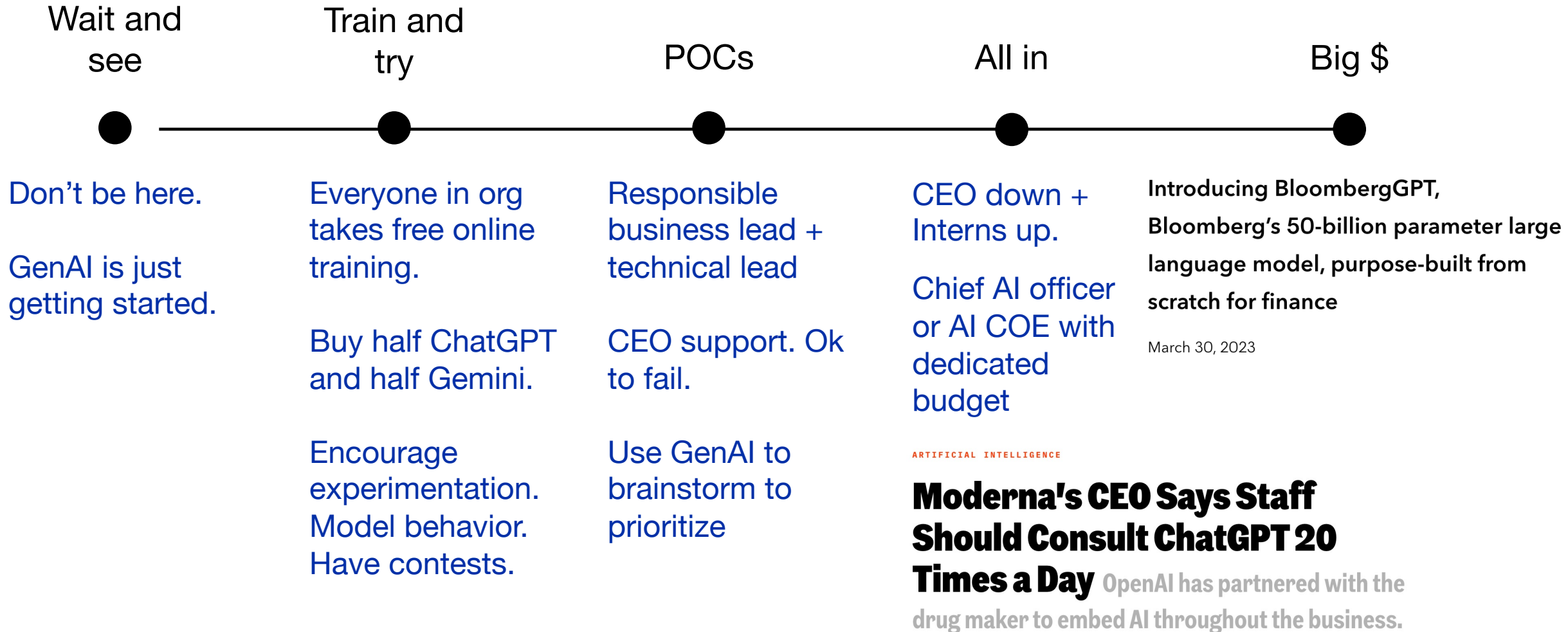
3. Think strategically

- Ask 4 questions

4 questions

- 1. What did you do that was valuable that is no longer valuable because it is now easy for everyone?**
 - Personalized customer service...
- 2. What was impossible before but you can now do?**
 - Provide freemium model for consulting...
- 3. Where can you move upmarket where you couldn't before?**
 - Customized marketing and sales for each prospect
- 4. Where can you move downmarket and democratize?**
 - Offering that would have taken 5 people and lost money now takes one person and makes money

The AI Leadership Spectrum



Hire smart

1. Pause hiring
2. Take every job description and break down to tasks
3. Can you automate any tasks with AI?
4. If yes, rewrite job descriptions
5. Re-evaluate hiring needs

Go down this rabbit hole

How AI Can Change the Way Your Company Gets Work Done

by Marc Zao-Sanders

July 03, 2024

<https://hbr.org/2024/07/how-ai-can-change-the-way-your-company-gets-work-done>

Good article, but more importantly, good links.

How to think about Traditional AI

- Think of Traditional AI software as your assistant. Your assistant learns from data and/or trial and error so it can answer your questions with predictions.
- If you train with proprietary data, your assistant can answer questions that no other AI can answer.
- Your assistant can be a genius on any given topic, from identifying manufacturing defects to diagnosing cancer.
- At the same time, your assistant doesn't understand the real world and can't put things in context. Your AI assistant depends on you to have common sense and make sure its work is used appropriately.

How to think about Generative AI

Think of GenAI software as your personal intern. Your intern:

- Has memorized the internet and knows a lot about the world.
- Works for free, never complains, and is very very fast.
- Will follow your prompt instructions quite literally and not do anything you don't specify.
- Works better when playing a role.
- Needs help planning a series of steps.
- Makes up answers instead of saying "I don't know."
- Produces work that you have to check before putting your name on it.

Thank you

Would your leadership team, your board, or your employees benefit from a better understanding of AI?

Would you like to receive a weekly newsletter with the top three AI stories from the previous week?

Please reach out at mhayes@getpractical.ai.