

Information and Intelligence

A whirlwind tour of AI

Disclaimer:

- I'm not an AI researcher.
- This is for entertainment purpose.
- Take everything I say with a grain of salt.
- Opinions are solely my own.

“The Hardest Logic Puzzle”



- Three gods: Truth, Lie, Random (in unknown order)
- They answer yes/no using “Da” and Ja (meaning unknown)
- Exactly three yes-no questions; one question to one god
- Goal: identify each god



- Three gods: Truth, Lie, Random (in unknown order)
- They answer yes/no using “Da” and Ja (meaning unknown)
- Exactly three yes-no questions; one question to one god
- Goal: identify each god

- 6 positions
- 3 questions $\rightarrow 2^3=8$ states at most
- Random’s answer does not provide any information $\rightarrow 2^2=4$ states

$4 < 6$, this is not solvable?

- We assumed Random ALWAYS provides one of the answers
- What if we can remove Random after the first question?

$$\rightarrow \frac{1}{3} \times 4 + \frac{2}{3} \times 8 = 6\frac{2}{3} > 6$$



- Three gods: Truth, Lie, Random (in unknown order)
- They answer yes/no using “Da” and Ja (meaning unknown)
- Exactly three yes-no questions; one question to one god
- Goal: identify each god

- We assumed Random ALWAYS provides one of the answers
- What if we can remove Random after the first question?

$$\rightarrow \frac{1}{3} \times 4 + \frac{2}{3} \times 8 = 6\frac{2}{3} > 6$$

We must be able to identify which god is NOT Random from the first question.

LLMs could not solve this puzzle correctly

ChatGPT 4o

✓ The Steps (in brief)

1. Pick any god (say, A) and ask:

"If I asked you 'Are you Random?', would you say 'da'?"

- If the answer is inconsistent in later steps, then A is Random.
- If it makes logical sense, A is either Truth or Lie.



Claude Sonnet 4

Step 1: Find a non-Random god

Ask god B: "If I asked you 'Is A Random?', would you say 'da'?"

This complex question works because:

- If B is True or False, they'll give a consistent answer about A
- If B is Random, the answer is meaningless anyway
- The double-layering ("if I asked you...would you say...") makes True and False gods equivalent in their response pattern

Analysis of B's answer:

- If B says "da": Then A is Random, so B and C are True/False
- If B says "ja": Then A is not Random, so A is True/False



LLM could not solve this puzzle correctly until the reasoning model

Gemini 2.5 Pro

1. **Ask God B:** "If I were to ask you 'Is A Random?', would you answer 'ja'?"
 - If B answers 'ja', C is the non-Random Pivot God.
 - If B answers 'da', A is the non-Random Pivot God.
2. **Ask the Pivot God (identified in step 1):** "If I were to ask you 'Are you True?', would you answer 'ja'?"
 - If the Pivot God answers 'ja', they are True.
 - If the Pivot God answers 'da', they are False.
3. **Ask the Pivot God:** "If I were to ask you 'Is B Random?', would you answer 'ja'?" (Or ask about the other non-Pivot god).
 - If the Pivot God answers 'ja', then B (or the god asked about) is Random. The final god has the remaining identity (True/False).
 - If the Pivot God answers 'da', then B (or the god asked about) is *not* Random (having the remaining True/False identity). The final god must be Random.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Identify Smaller Problems to Solve

5. Initial Brainstorming & Challenges:

- **The Language Problem:** Not knowing 'da'/'ja'. This means a simple "Are you True?" won't work directly, as you don't know if 'da' means yes or no.
- **The Random Problem:** Random's answers are useless for gaining direct information. We need a way to either identify Random early or exclude them from informative questioning.
- **The Liar Problem:** False inverts the truth. How to handle this combined with the language problem?
- **Limited Questions:** Only three questions to identify three gods. Efficiency is crucial.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

One problem at a time

6. Develop Key Strategies:

- **Handling the Language Problem:** Find a question structure where the *answer* ('da' or 'ja') reliably correlates with a specific underlying truth, regardless of which word means 'yes'. This often involves embedding a condition within the question. The classic "If I asked you Q, would you say 'ja'?" structure comes to mind. Let's analyze this:
 - Ask god X: "If I asked you 'Is P true?', would you say 'ja'?"
 - Case 1: 'ja' means yes.
 - True: Would say 'ja' if P is true, 'da' if P is false. Answers 'ja' if P is true.
 - False: Would *actually* say 'da' if P is true (lie about P), 'ja' if P is false (lie about P). But the question asks

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Prove a theorem and use it repeatedly

- *Conclusion:* This embedded question structure ("If I asked you P, would you say 'ja'?") works! The answer is 'ja' if and only if P is true, and 'da' if and only if P is false, *regardless* of what 'ja'/'da' mean, and *regardless* of whether you are asking True or False. This is a crucial breakthrough! Let's call this the "reliable question structure".

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Exploration with many false starts

- **Handling the Random Problem:** How to avoid Random or identify them? If we ask a question whose answer *must* be predictable for True and False, any god who answers differently *could* be Random. Or, can we formulate a question that forces Random to reveal themselves, or forces a non-Random god to be identified?
 - Idea: Ask a question where True and False *must* give the same answer ('ja' or 'da'). If we get that answer, the god *could* be True or False. If we get the *other* answer, the god *must* be Random.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Exploration with many false starts

- Alternative approach: Identify a non-Random god first. If we can ask a question that identifies *one* god as definitely *not* Random, we can use that god for subsequent questions.
- This doesn't quite isolate a non-Random god yet. Let's rethink. We need a question whose *truth value* is linked to the identity of the god being asked.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Double check and self review

- Hold on, is that right? Let's re-verify the reliable structure for True/False applied to $P = \text{"Is A Random?"}$.
 - If A is Random ($P=\text{True}$): Both True and False will answer 'ja'.
 - If A is not Random ($P=\text{False}$): Both True and False will answer 'da'.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Double check and self review

- This isn't working to identify a non-Random god *for certain*. The issue is that Random B could mimic the True/False answer.
- **Refined Strategy - The Pivotal Insight:** The key is to devise a question whose answer points you towards a god who is *guaranteed not to be Random*.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Double check and self review

- **This isn't isolating a non-Random god.** The problem description *must* be solvable. What am I missing? Ah, the *structure* of the question can be even more complex. The standard solution uses a logical equivalence.
- Let's try a different first question, based on common solutions. Ask God A: "If I asked you 'Does B answer randomly?', would you say 'ja'?"
 - This feels similar.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Double check and self review

- So: If B answers 'ja', B could be True (if ja=yes) or False (if ja=yes or ja=no) or Random. If B answers 'da', B must be True (and ja=no) or Random. This doesn't seem to isolate a non-Random god either.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Reference knowledge

- **Let's revisit the Rabern & Rabern simpler solution structure, as it's often cited.**
 - Core Idea: Use a question that points towards a non-Random god based on the answer.
 - Question 1 (to God B): "If I were to ask you 'Is A Random?', would you say 'ja'?"
 - Let's analyze the *implication* of the answer.

LLM could not solve this puzzle correctly
until the reasoning model

Gemini 2.5 Pro

Finally reached the correct solution for Q1

To God B: "If I were to ask you 'Is God A the Random god?', would you answer 'ja'?"

If God B answers "ja":

- Possibility 1: B is True or False. This means the proposition "God A is Random" is true. So, A is Random.
- Possibility 2: B is Random. The answer is meaningless.
- In either possibility, God C cannot be Random. If B is Random, C isn't. If A is Random, C isn't. So, you have found a god (C) who is definitely not Random.

- **Success!** Question 1 identifies a god (either A or C) who is guaranteed *not* to be Random. Let's call this god the "identified non-Random god" or "Pivot God".

How does it do that?

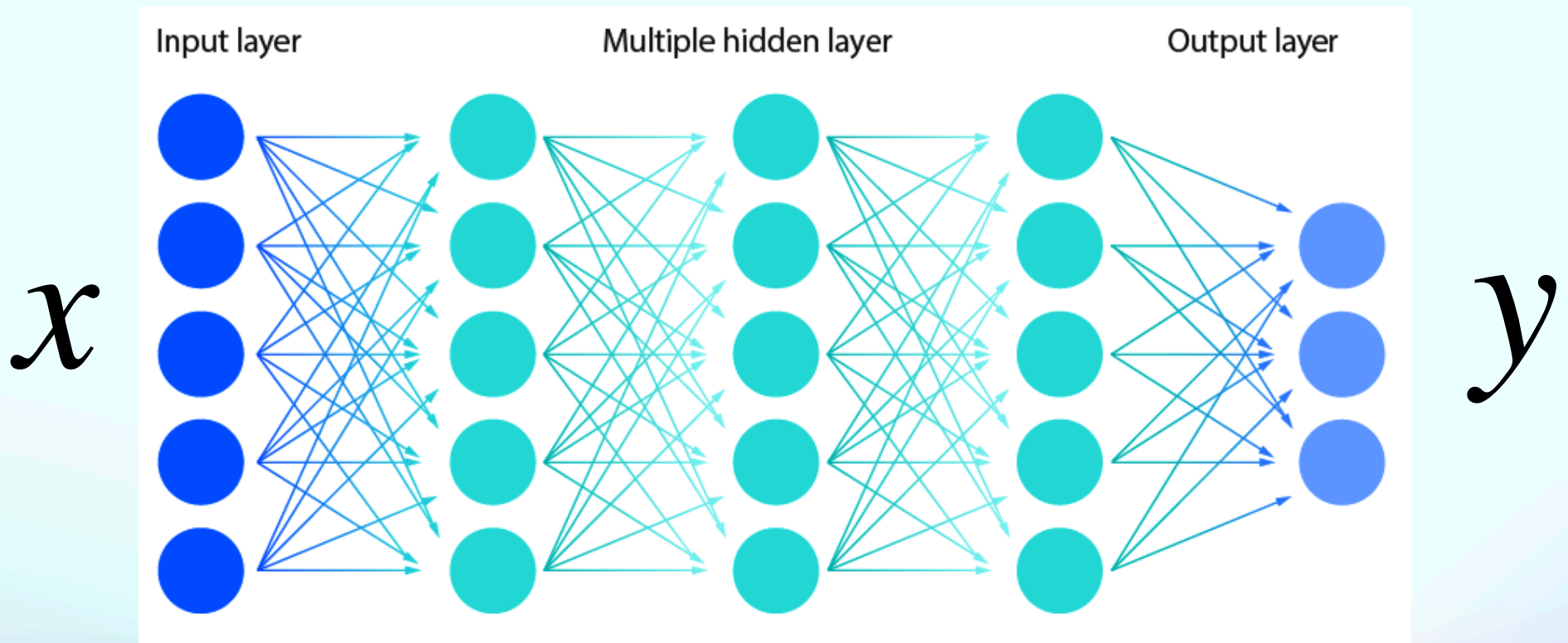
How does it do that?

- Under the Hood
- Supervised Learning
- Reinforcement Learning
- Minds vs Machines

Act 1: Under the Hood

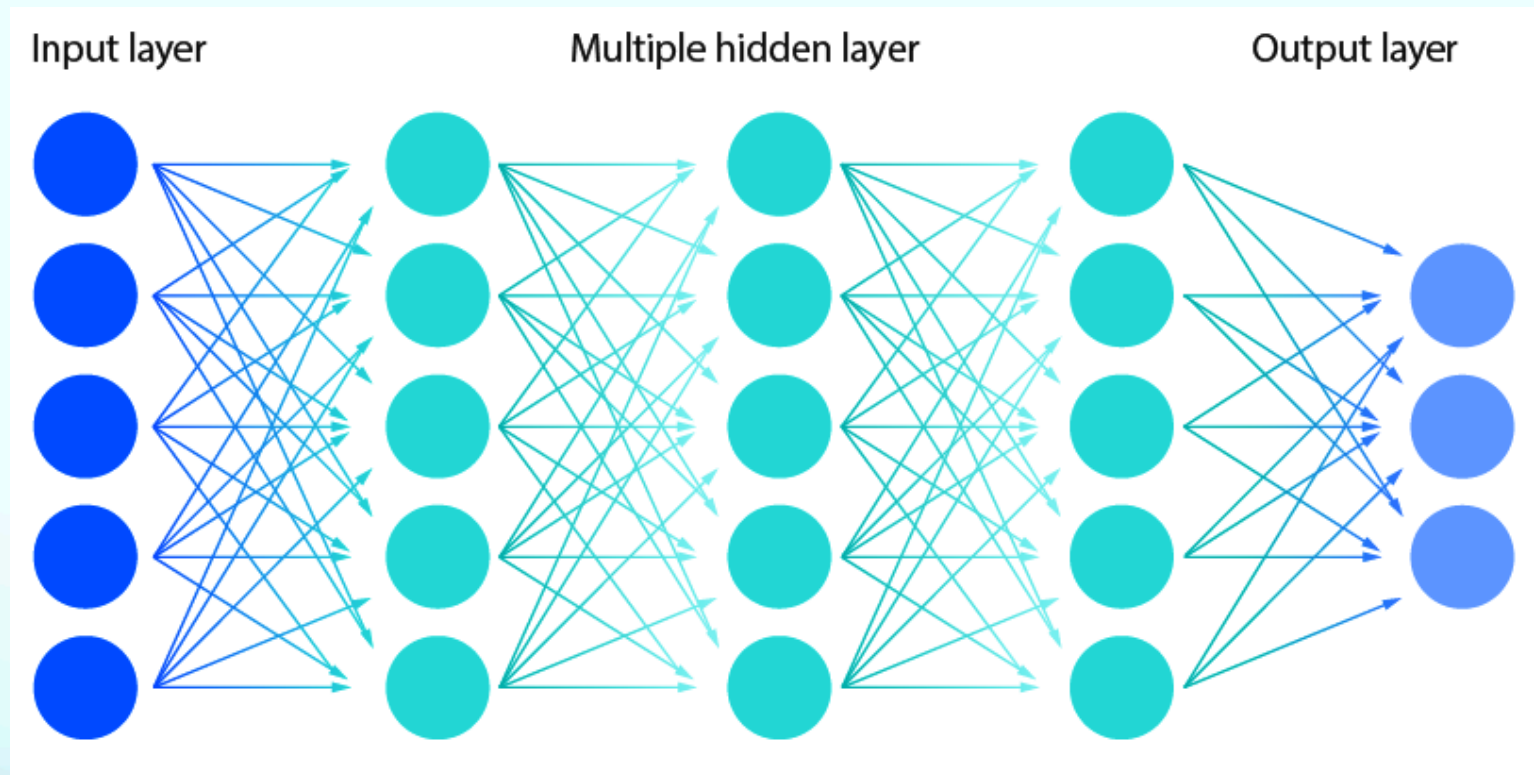
Neural Network & Universal Approximation

(Supervised) Learning is Fitting the Function



$$x \rightarrow f(x) \rightarrow y$$

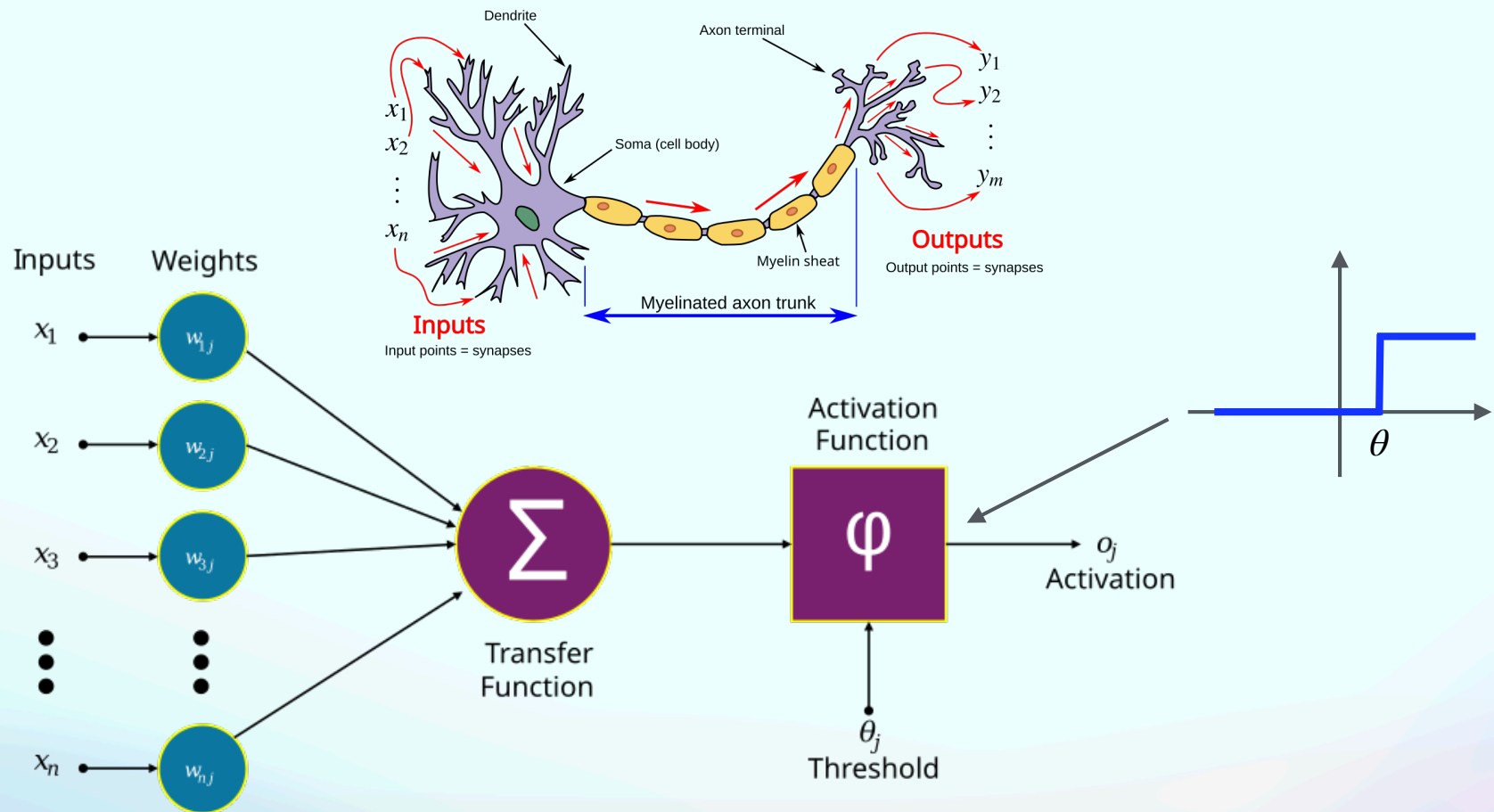
Universal Approximation



A deep neural network with ... , can approximate any **measurable function** to arbitrary accuracy.

Universal Approximation

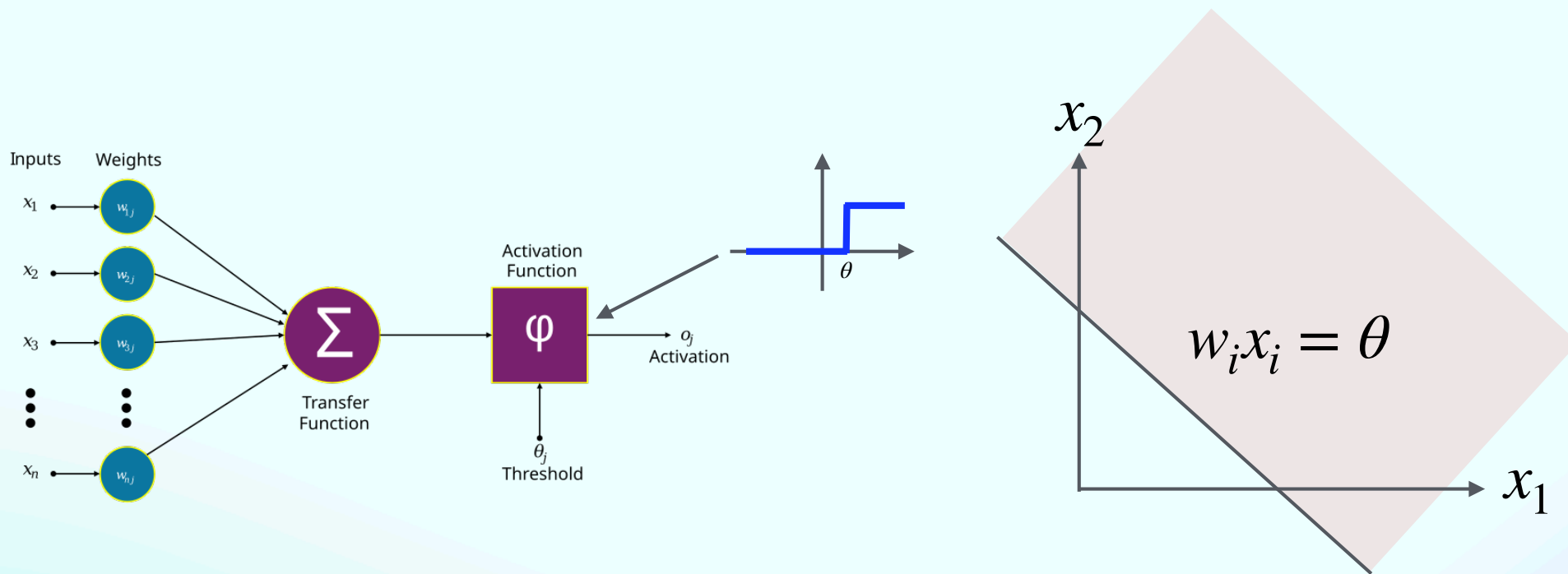
Single Neuron



$$y = \varphi \left(\sum_i w_i x_i \right)$$

Universal Approximation

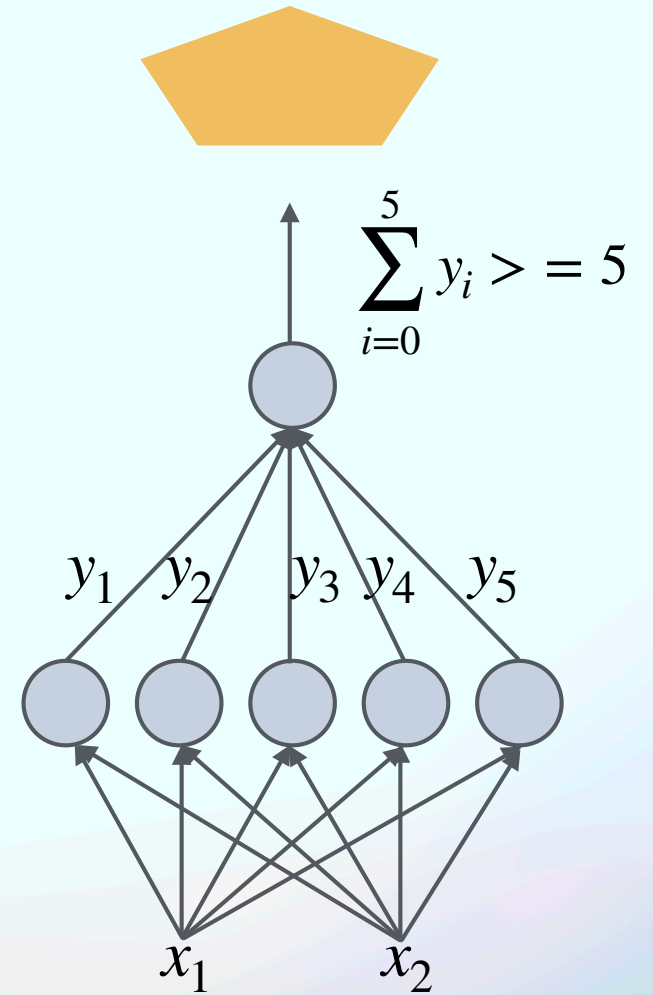
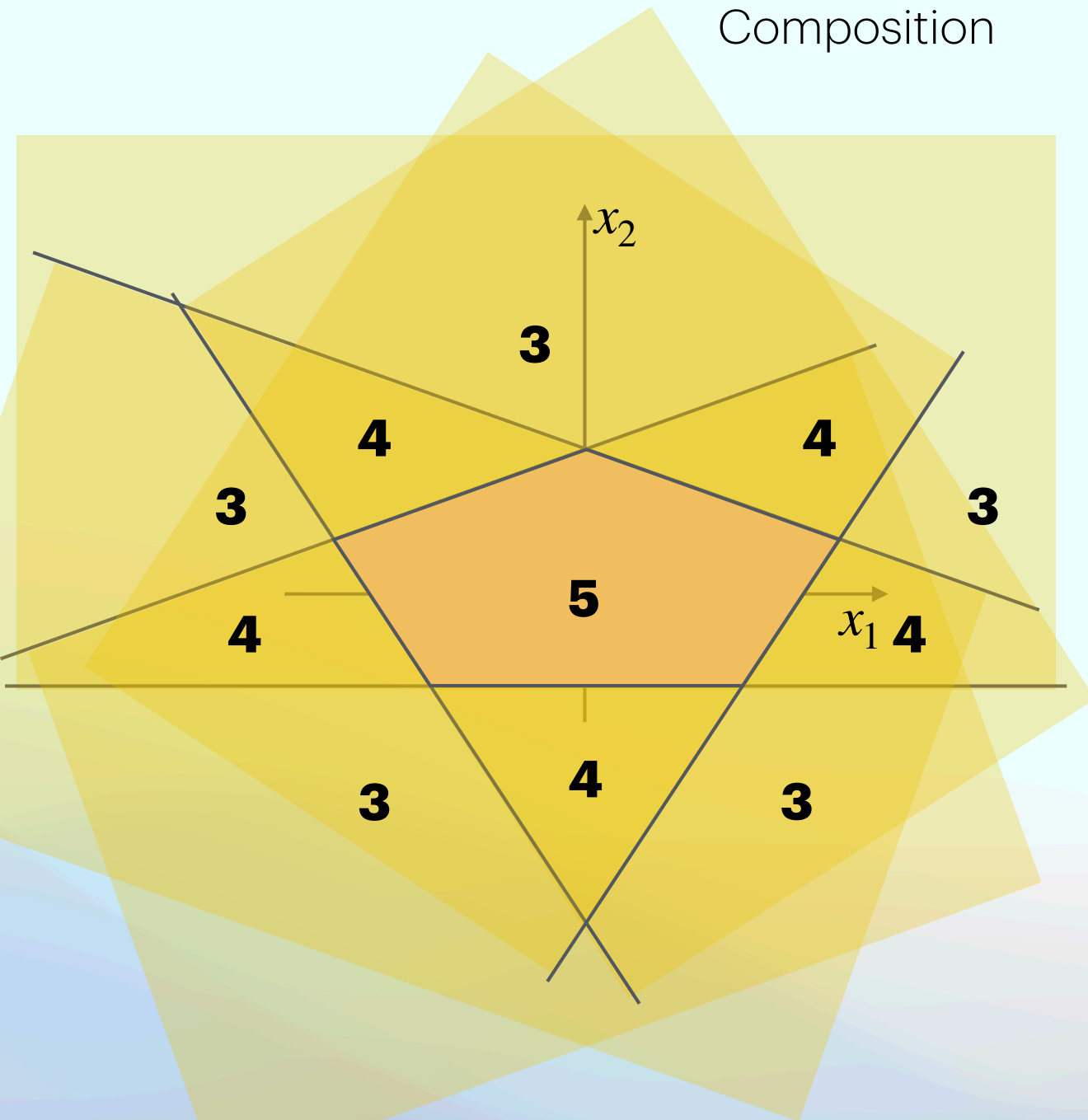
Single Neuron



$$y = \varphi \left(\sum_i w_i x_i \right) = \begin{cases} 1 & \text{if } \sum_i w_i x_i \geq \theta \\ 0 & \text{else} \end{cases}$$

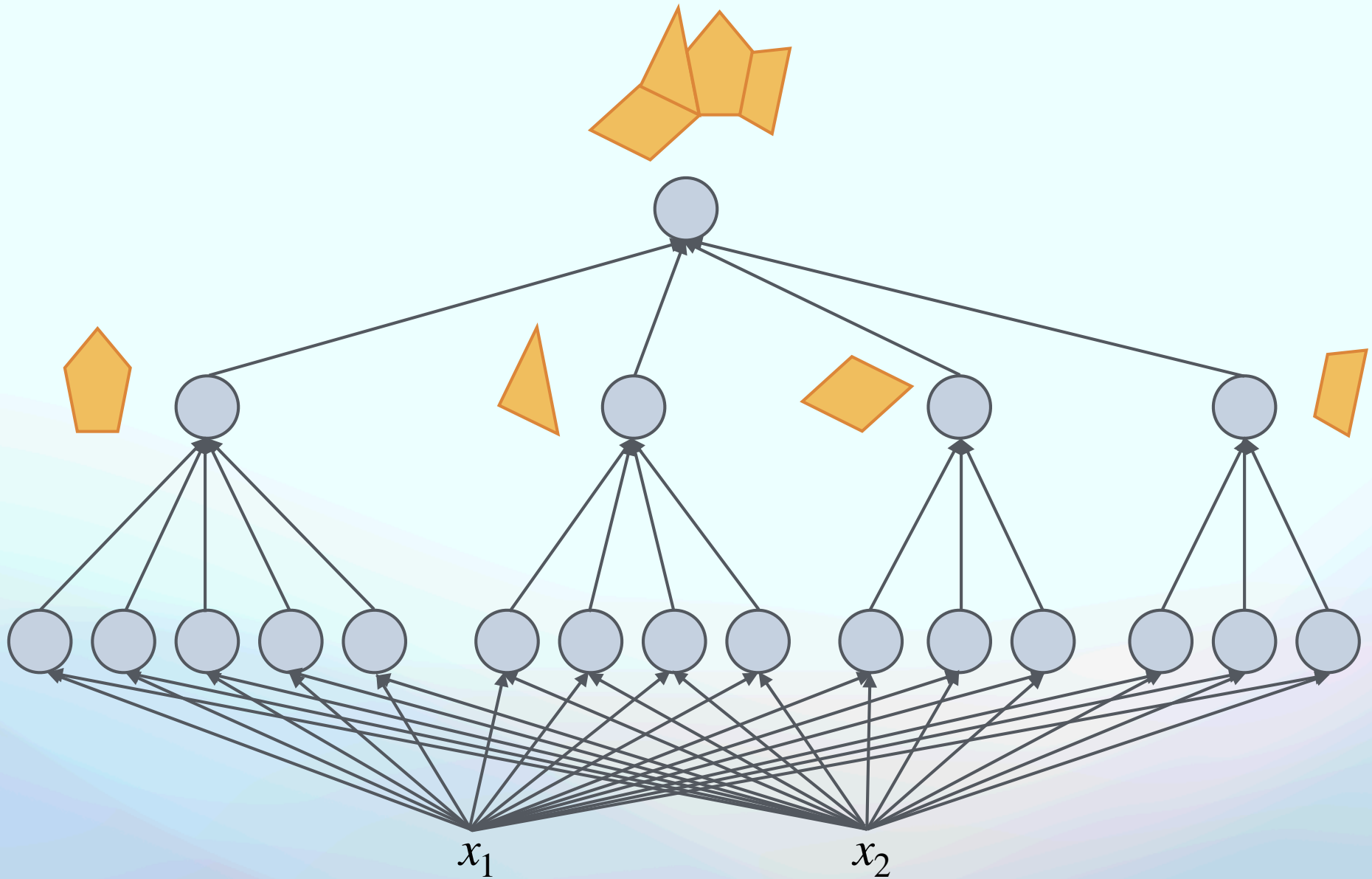
Universal Classification

Composition



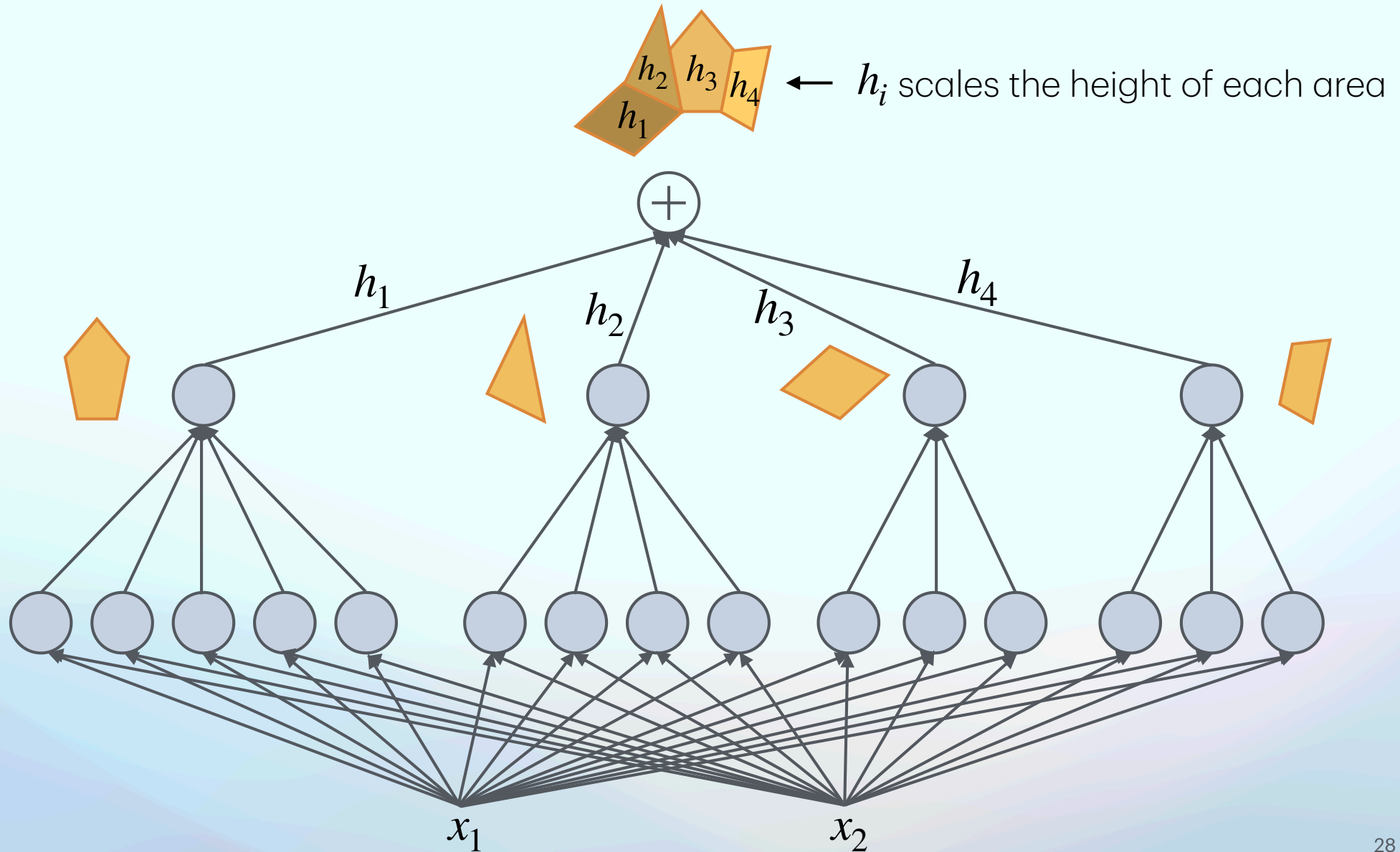
Universal Classification

Universal Classifier: can compose arbitrary decision boundaries



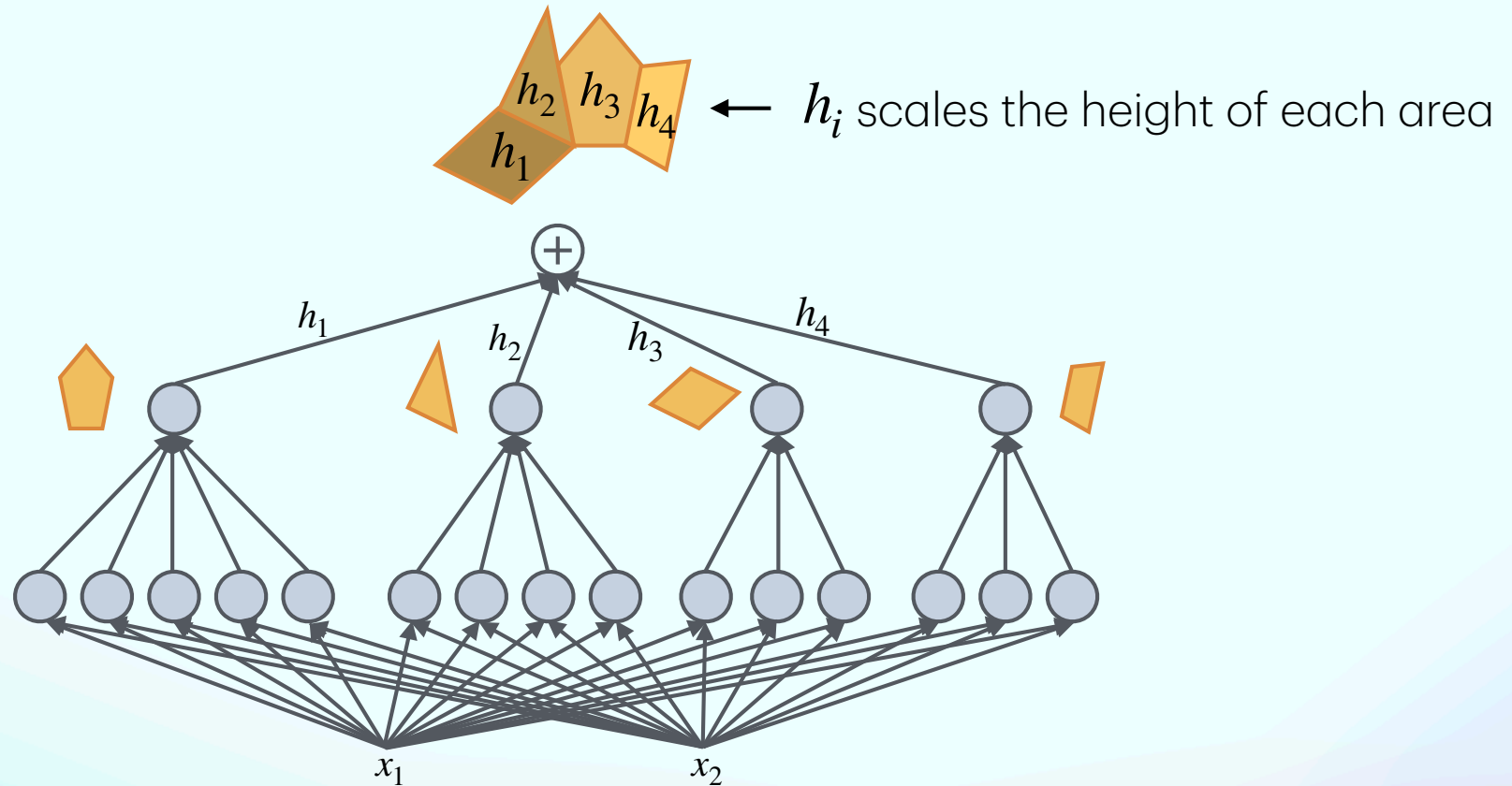
Universal Approximation

Slight change of the last layer can approximate any 2-D function



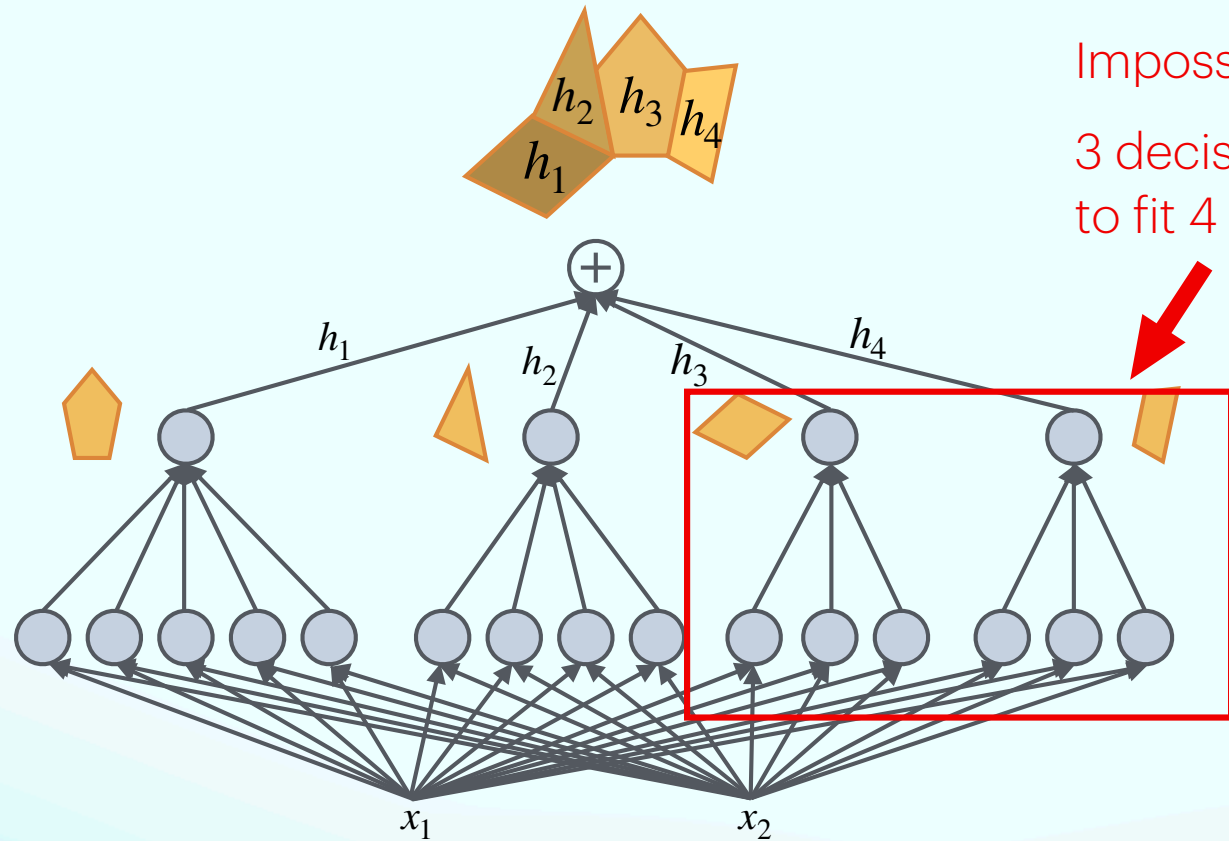
Universal Approximation

Slight change of the last layer can approximate any 2-D function



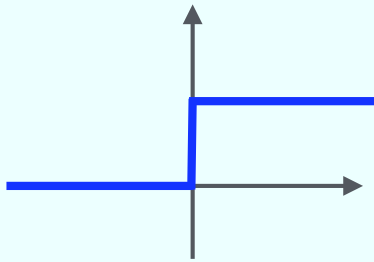
The universal approximation theorem **holds for** functions of any finite number of dimensions.

Deep Neural Network



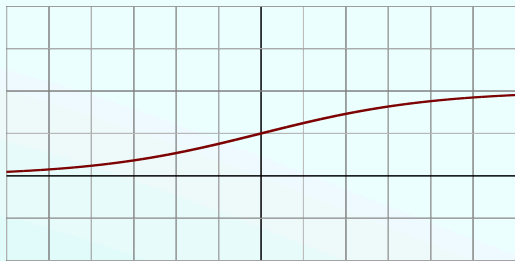
layers have to be sufficiently wide to capture the feature

Deep Neural Network

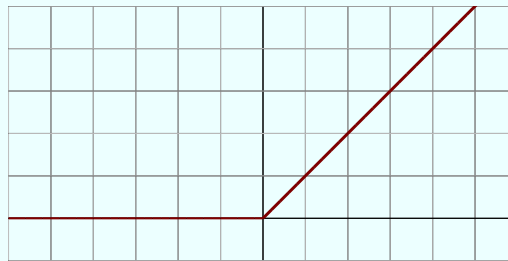


This is because we used the threshold activation function that discards information about the inputs

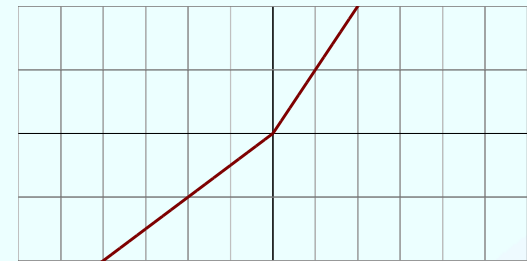
There are many other **monotonic non-linear** activation functions that pass more information to the next layer



Sigmoid



ReLU



Leaky ReLU

Can reduce layer width, but will require greater depth, to permit later layers to capture patterns

Deep Neural Network

HOW A DEEP NEURAL NETWORK SEES

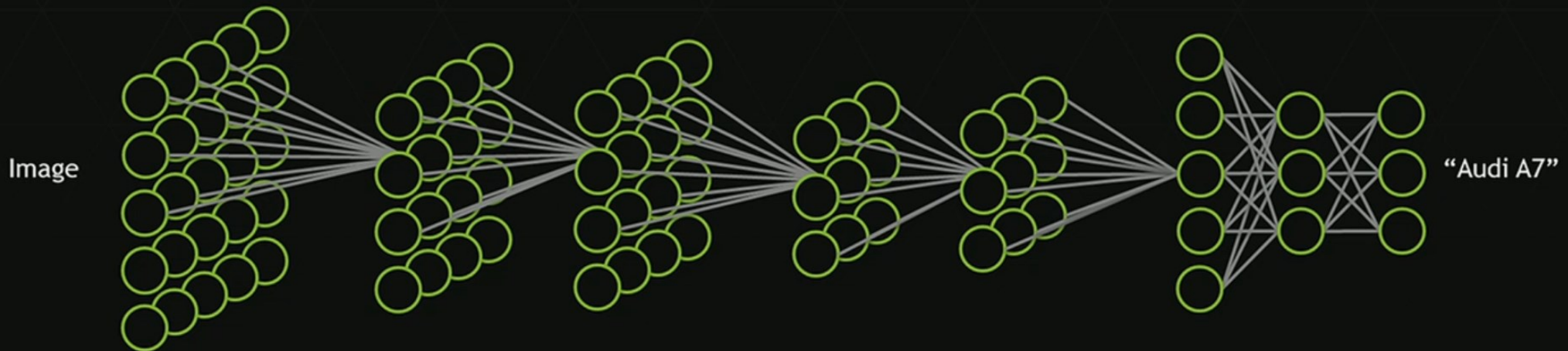
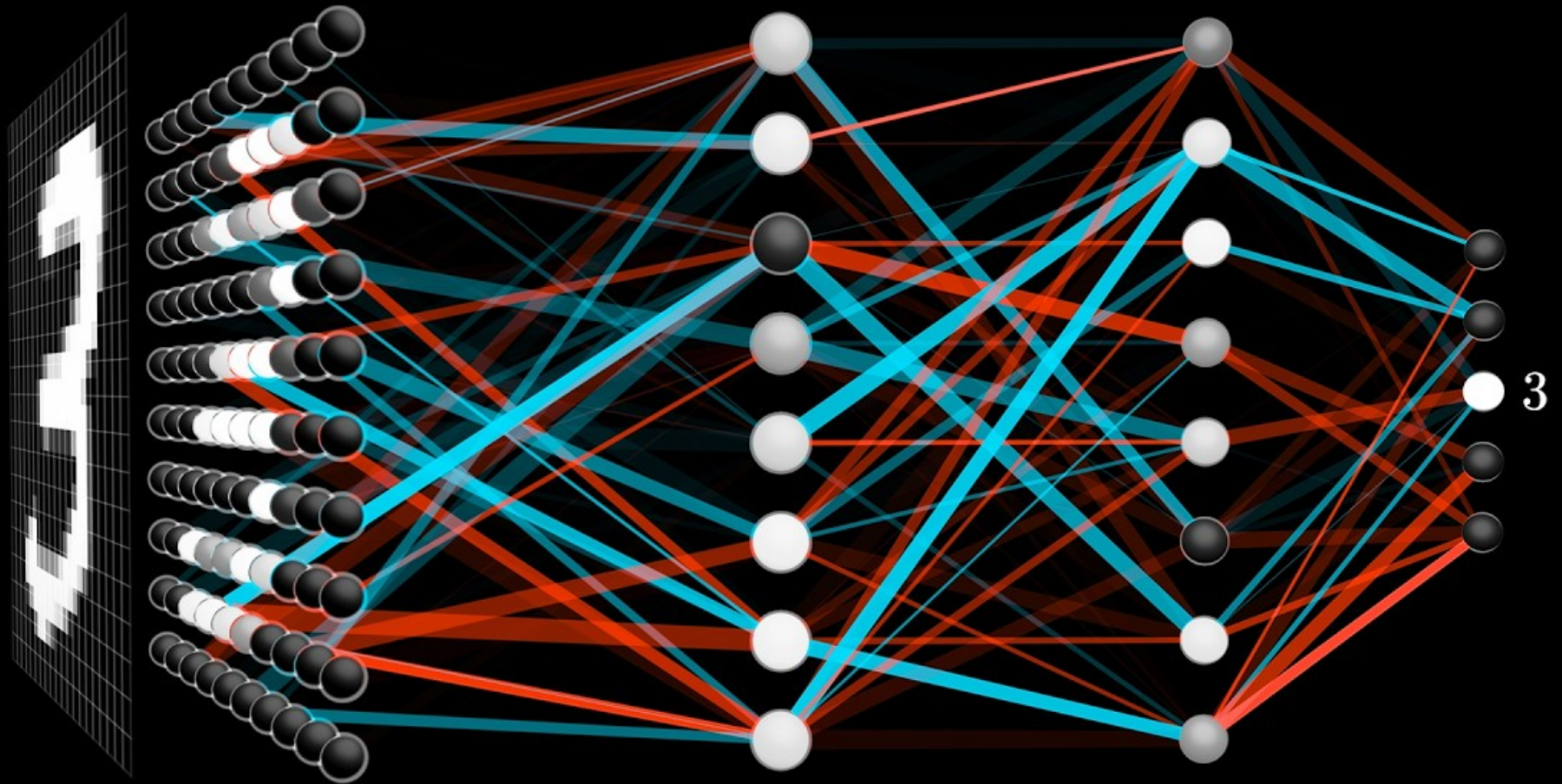


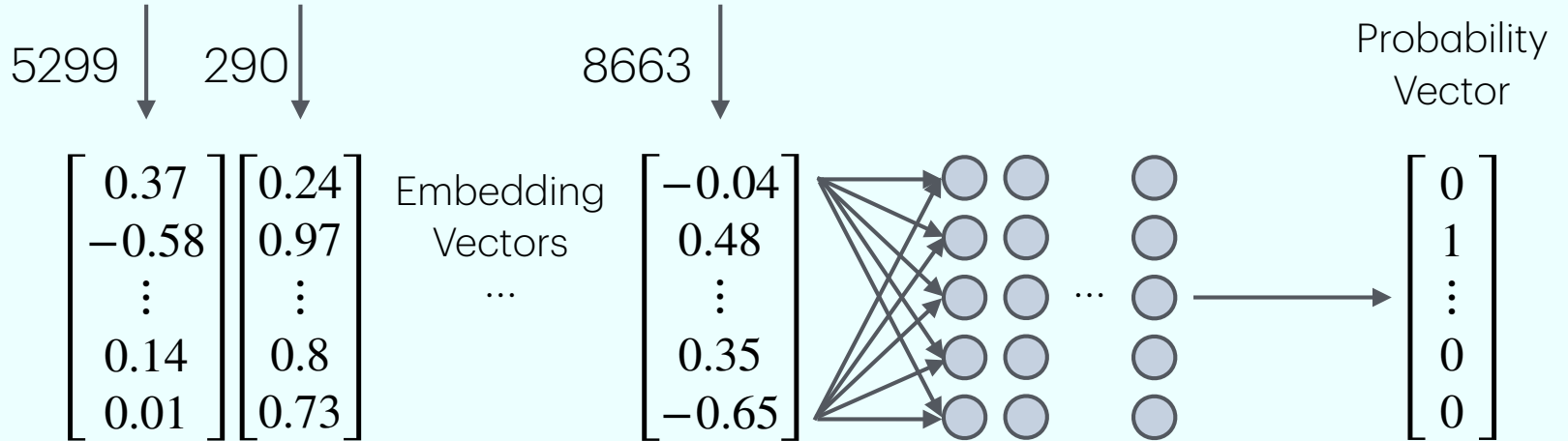
Image source: "Unsupervised Learning of Hierarchical Representations with Convolutional Deep Belief Networks" ICML 2009 & Comm. ACM 2011. Honglak Lee, Roger Grosse, Rajesh Ranganath, and Andrew Ng.

How the world becomes numbers

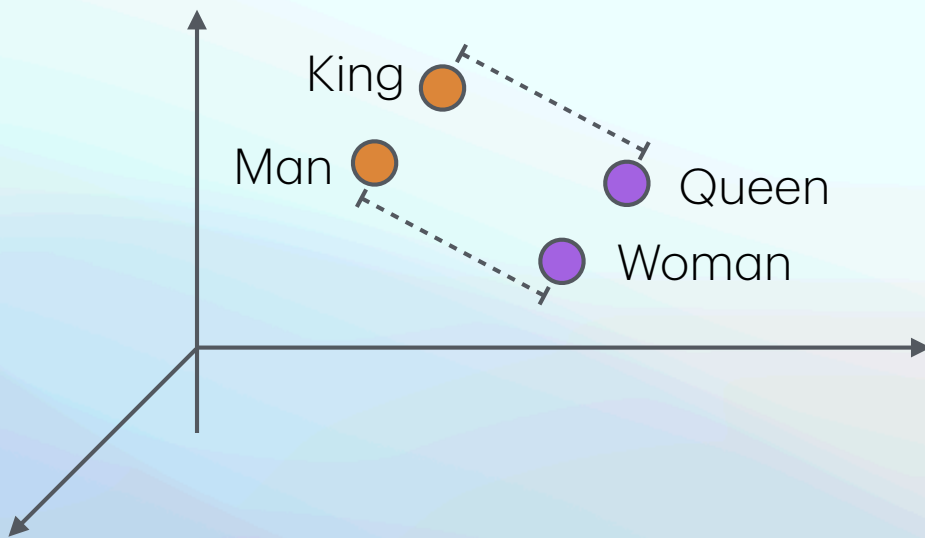


Credit: Youtube @3Blue1Brown
Video: But what is a neural network? | Deep learning chapter 1

How the world becomes numbers



LLM learns to predict the probabilities of the next token.



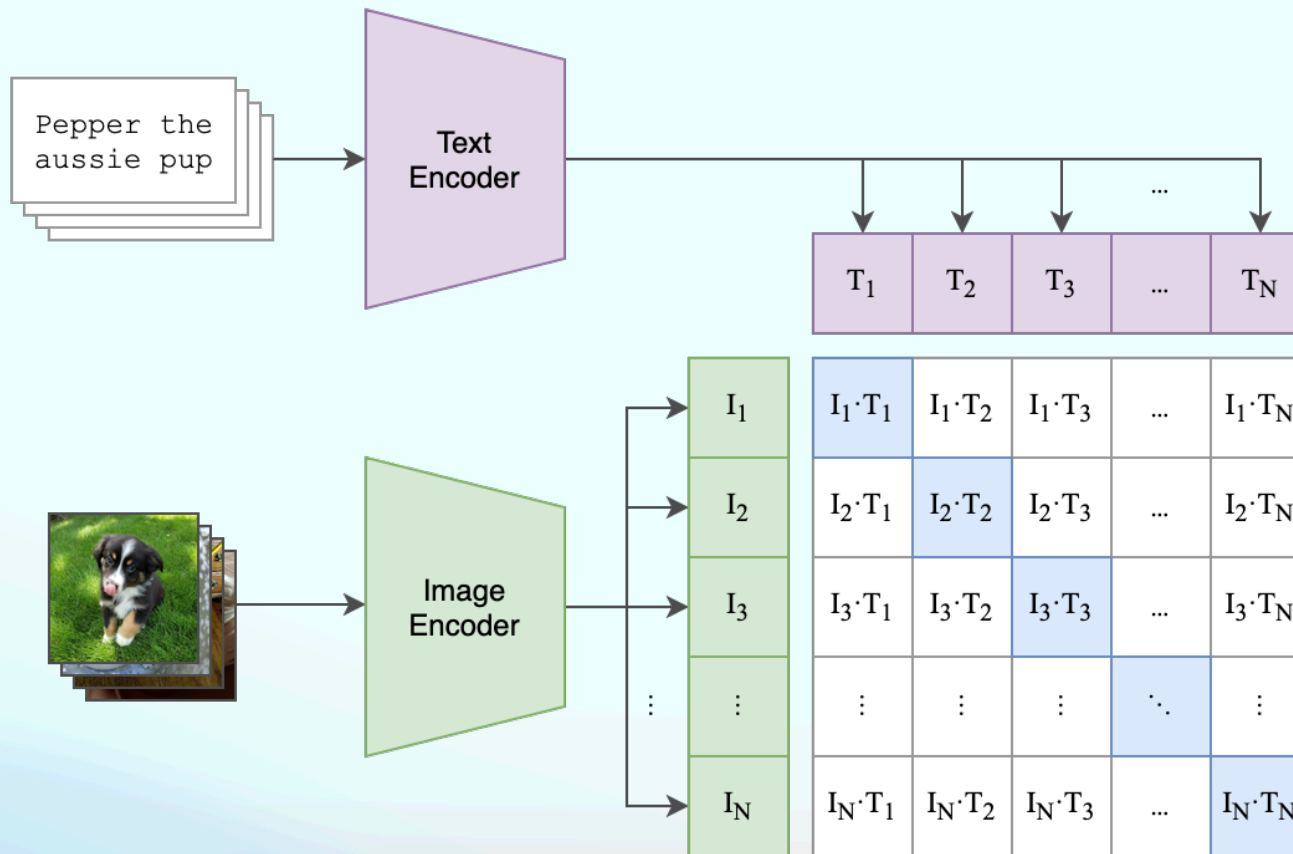
King - Man + Woman = Queen

Embeddings can capture semantics in vector space

How the world becomes numbers

Correlate Text with Images

(Credit: [CLIP model](#))

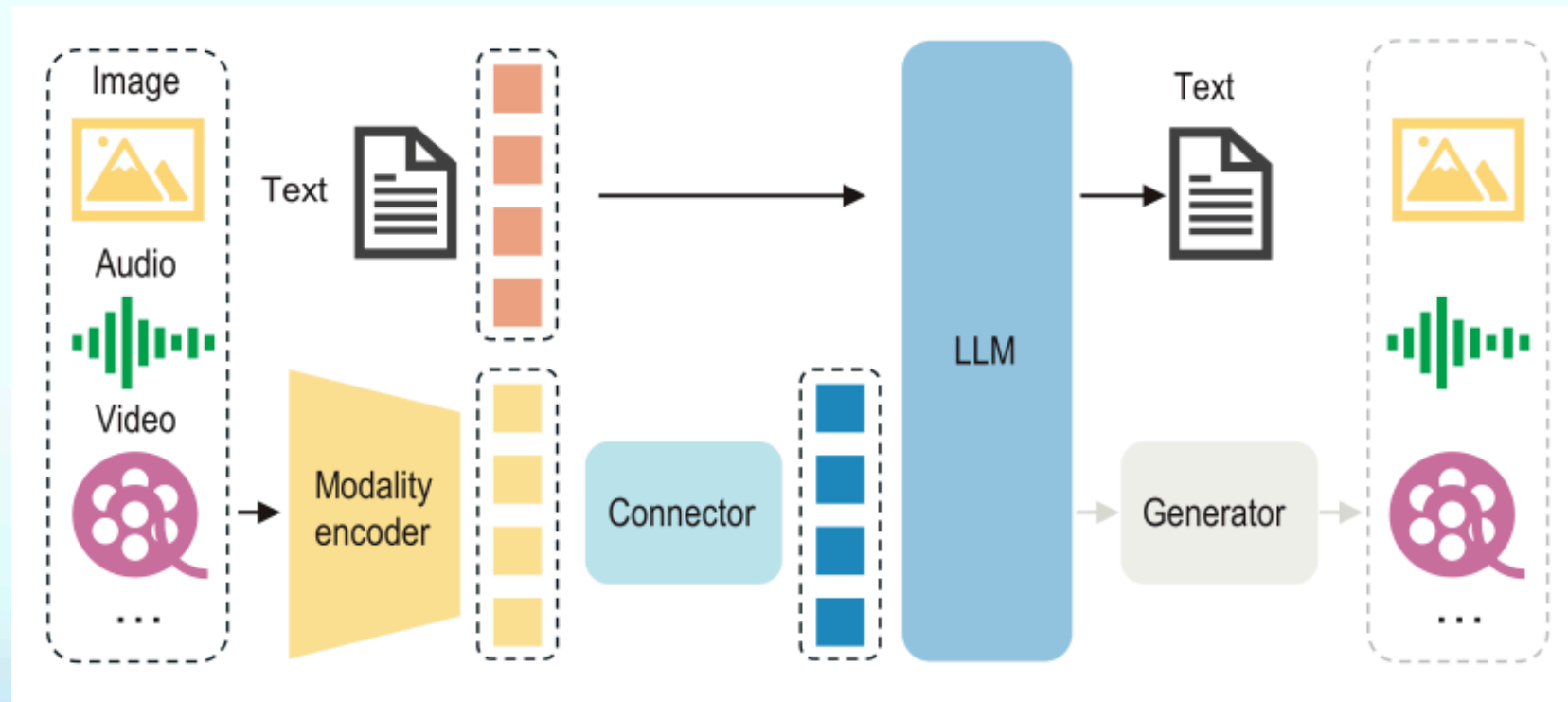


*Fundamentally, everything is information —
It can be represented by numbers and can be learned.*

How the world becomes numbers

Multimodal LLM

(Credit: [A survey on multimodal large language models](#))

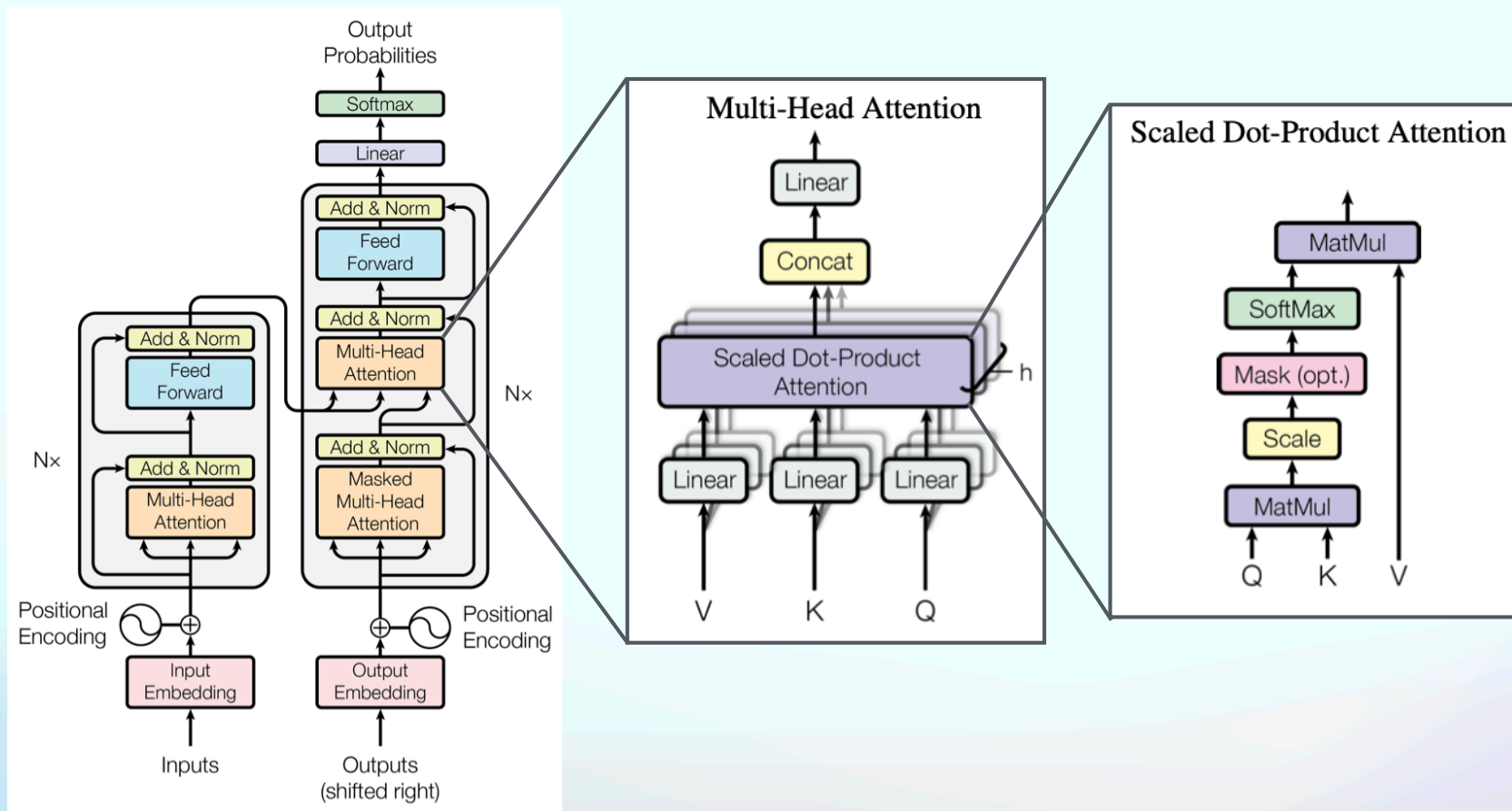


*Fundamentally, everything is information —
It can be represented by numbers and can be learned.*

How the world becomes numbers

Structure of the network is important for information flow

(Credit: *Attention is all you need*)



Act 1: Under the Hood

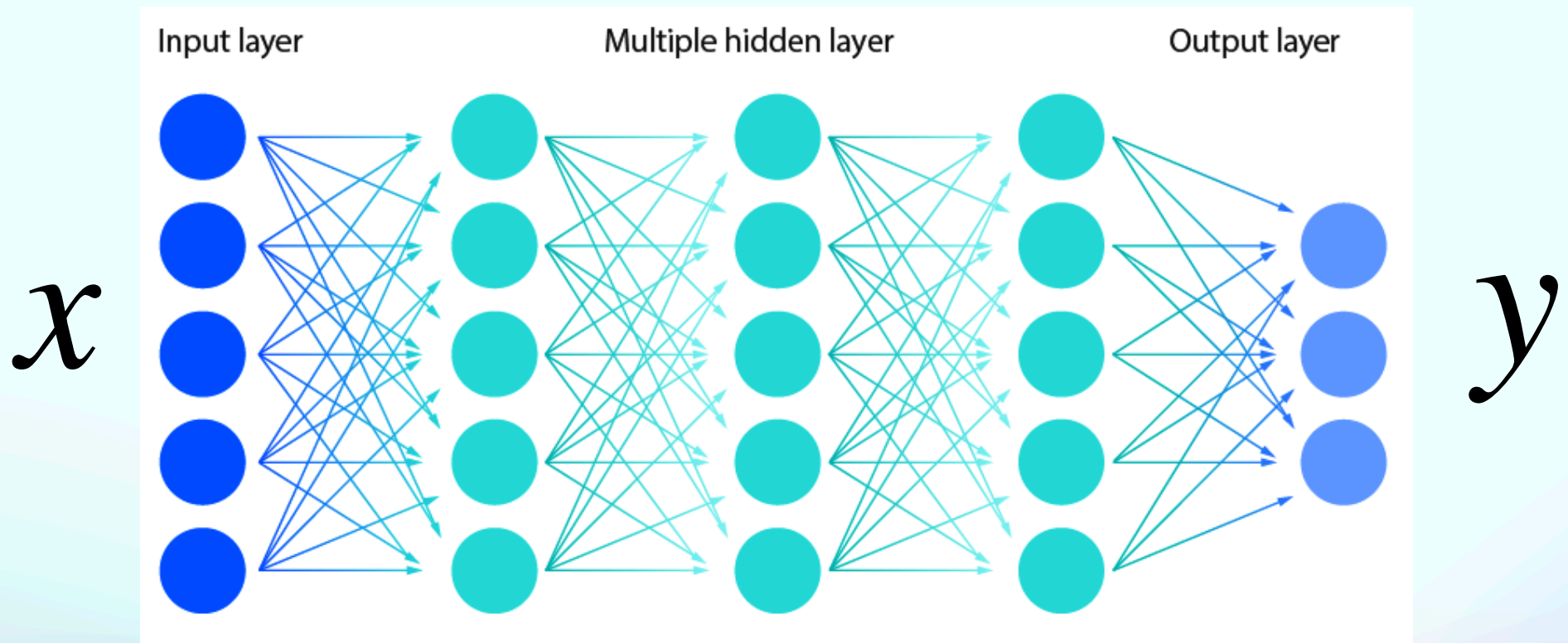
Neural Network & Universal Approximation

- Network with enough capacity can fit any measurable function
- Structure influences learning efficiency
- Embeddings let us represent the world as vectors
- Neural networks transform and reason with these representations

Act 2: Supervised Learning

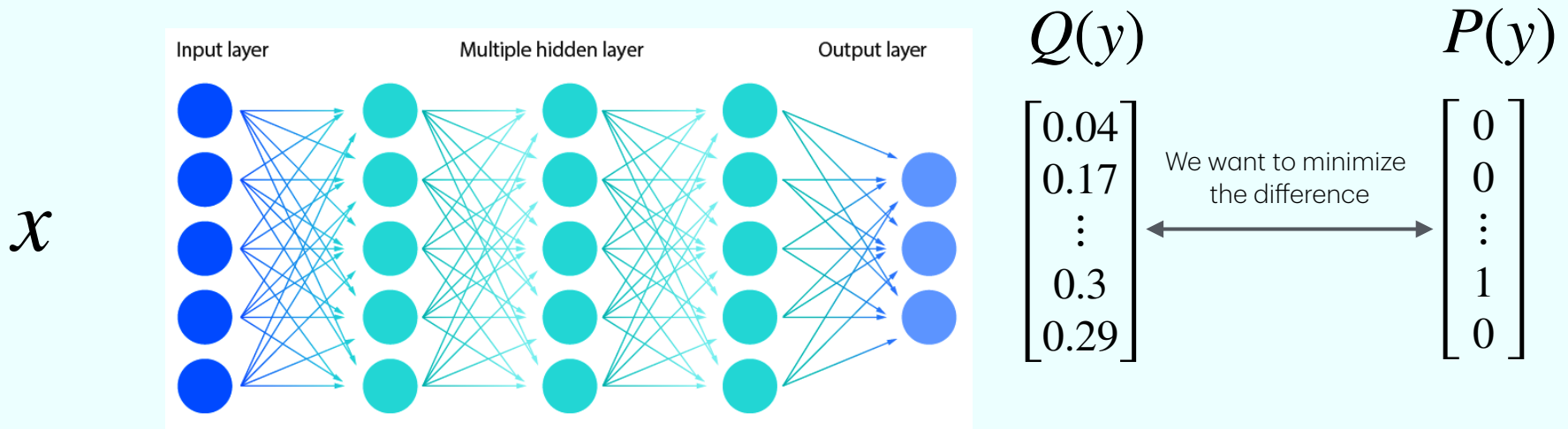
Bias is Necessary for Intelligence

Learning is finding the weights of the network



$$x \rightarrow f(x) \rightarrow y$$

Loss Function



Loss Function for LLM:

Gibbs Inequality: Minimum when $P = Q$

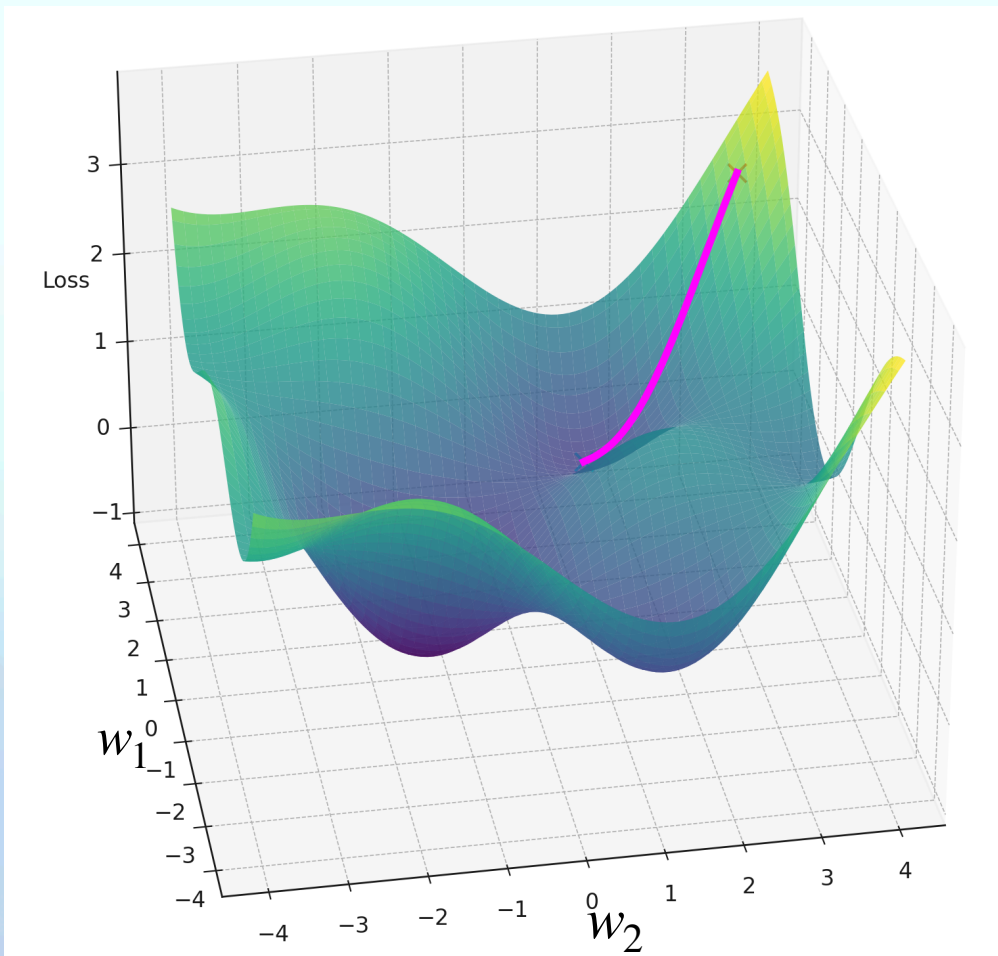
$$L(w) = - \underbrace{\sum P \log Q(y)}_{\text{Cross Entropy}} \leftarrow \underbrace{\sum P(y) \log P(y) - \sum P(y) \log Q(y)}_{\text{KL Divergence / Relative Entropy}} \geq 0$$

$$- \sum \log Q(y) \text{ for LLM}$$

When P is not a "one-hot" distribution
Used in "Model Distillation"

Gradient Descent

$$L(w) = \underbrace{- \sum P \log Q(y)}_{\text{Cross Entropy}}$$



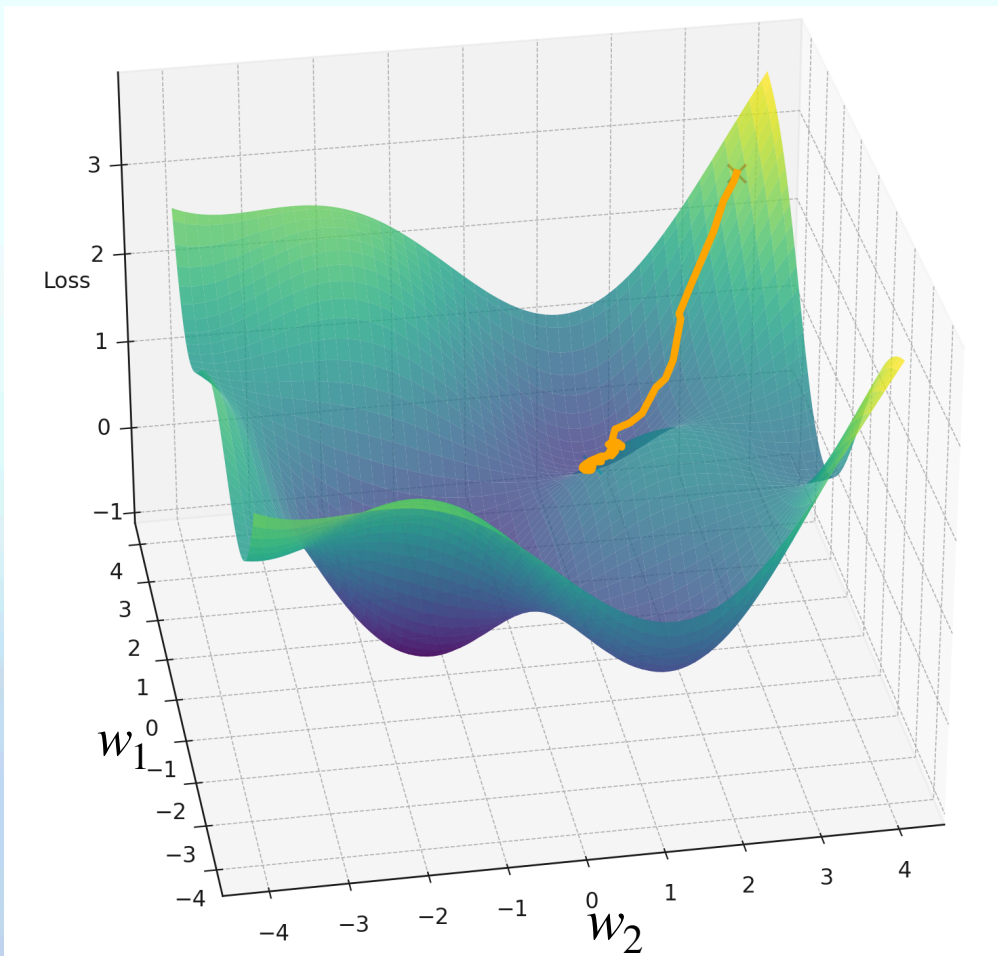
Gradient Descent

Update the weights:

$$w \leftarrow w - \eta \cdot \frac{\partial L}{\partial w}$$

Gradient Descent

$$L(w) = \underbrace{- \sum P \log Q(y)}_{\text{Cross Entropy}}$$



Gradient Descent

Update the weights:

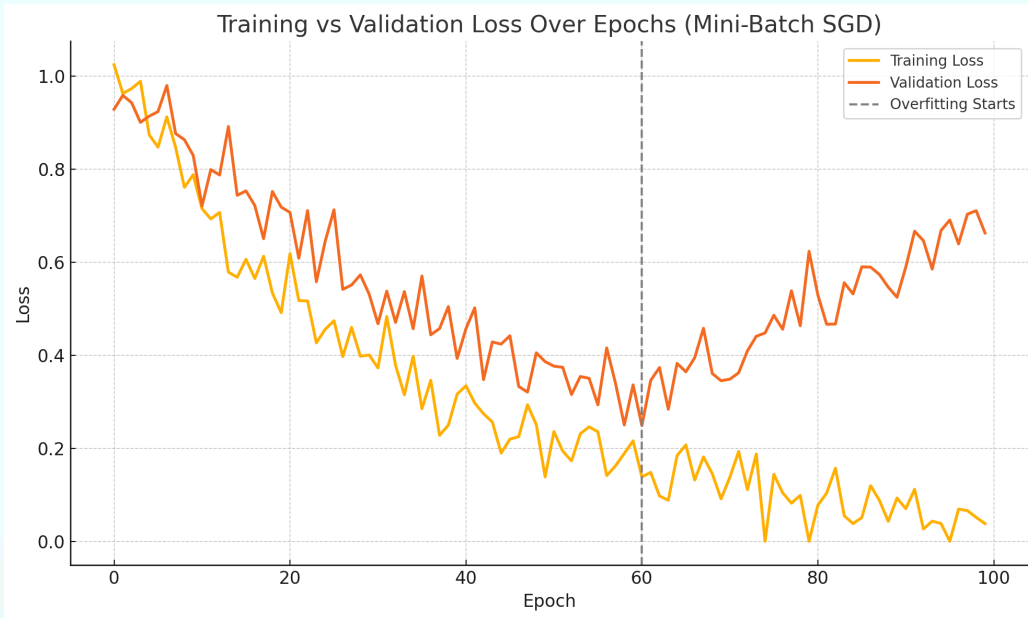
$$w \leftarrow w - \eta \cdot \frac{\partial L}{\partial w}$$

Stochastic Gradient Descent

In practice, we update the weights for mini-batches of the data

This adds “noise” to the gradient descent and makes the path more squiggly.

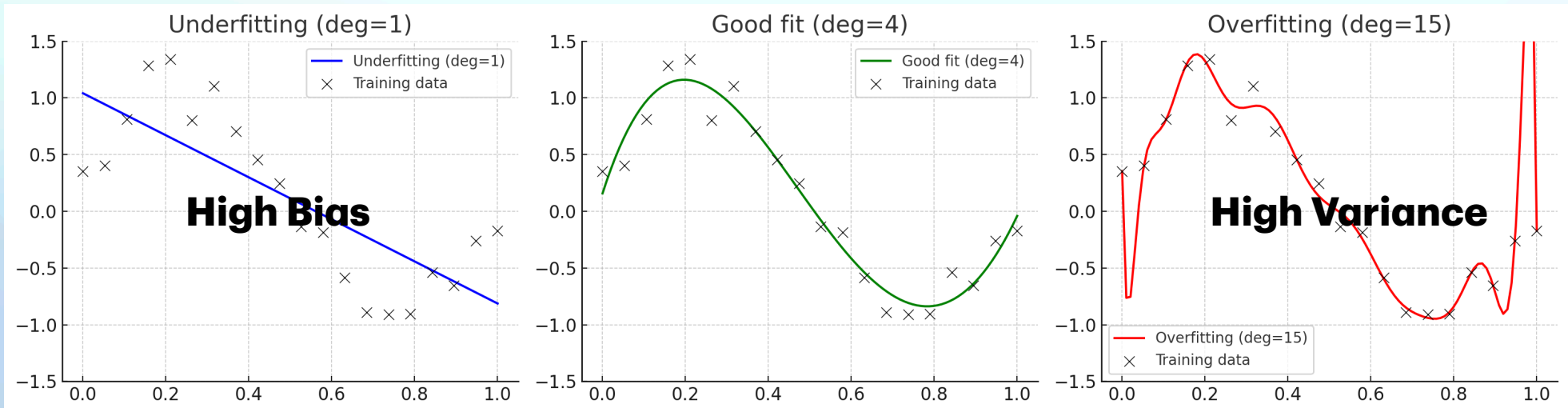
Bias-Variance Tradeoff



The best fit to the training data is not the best model

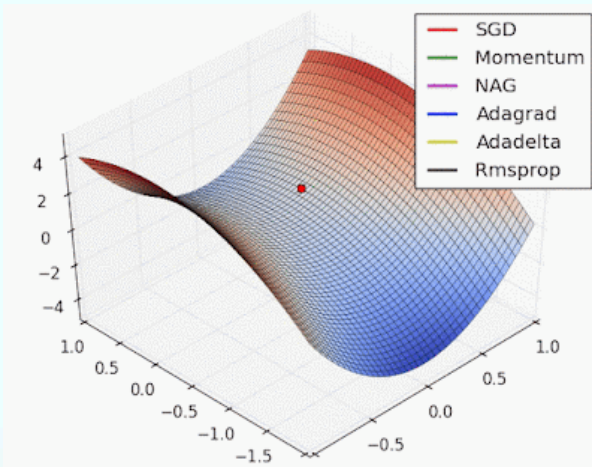
We don't want to be trapped by local minimums, but we are also not looking for the Global Minimum

We are looking for a minimum that's good for generalization



Mini-batch Stochastic Gradient Descent

Some suggested theories for why it works so well (still an active research area)



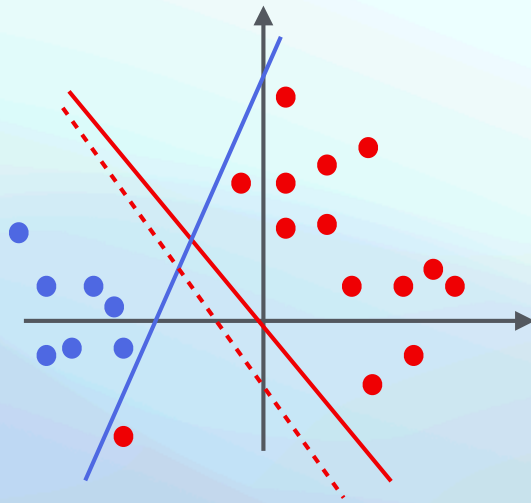
Credit: [Alec Randford](#)

A. Why does it not underfit?

1. High dimensional space \rightarrow more saddle points
2. Over-parameterization smooths the landscape
3. Noise helps escape local minima

B. Why does it not overfit?

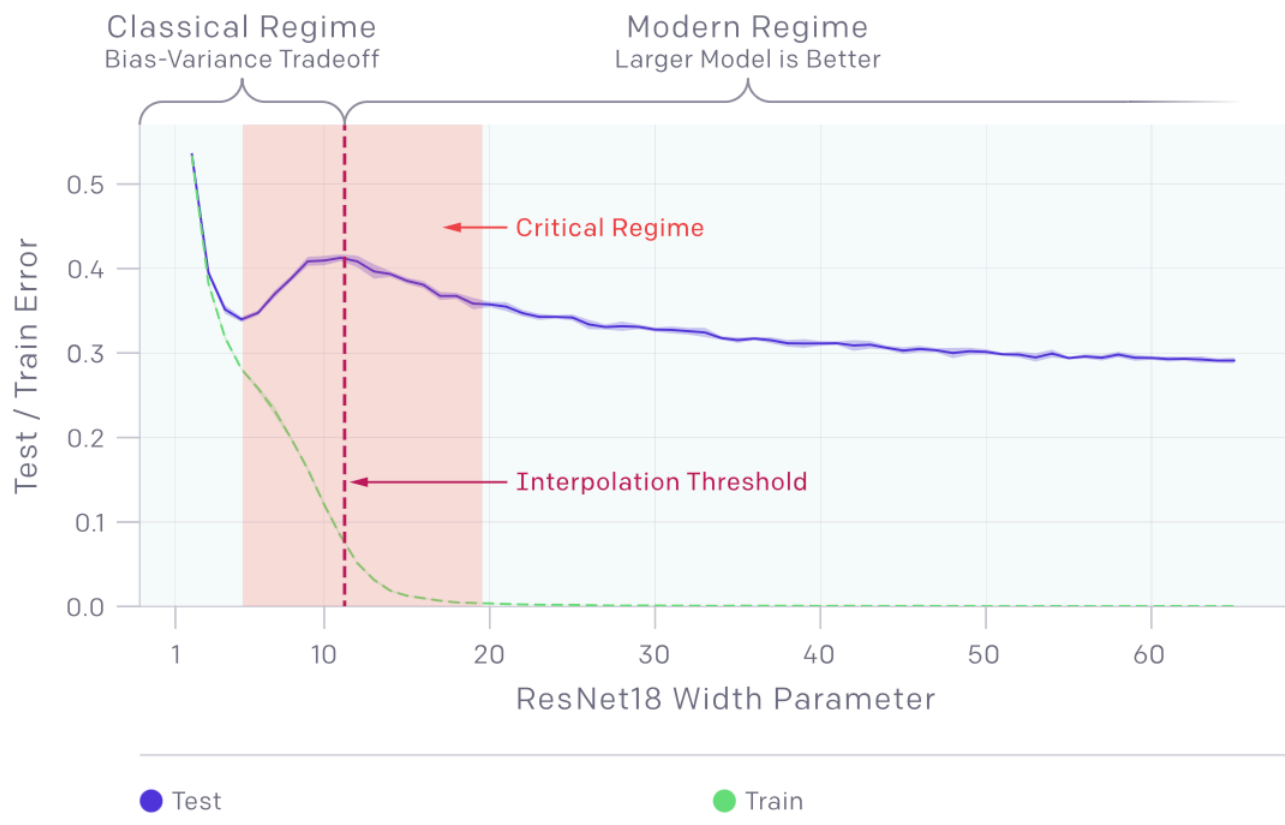
1. Its intrinsic “bias” prevents it from overfitting to outliers.
2. The noise also nudges the model toward flatter minima, which tend to generalize better.
3. The noise makes certain “saddle points” more attractive than the global minimum.



Avoiding Overfitting

Dropout **Massive and Diverse Datasets**
Weight Decay **Adding Noise of Corruption**
Mini-batch Stochastic Gradient Descent

Double Descent



Instead of overfitting, over-parameterized complex large models seem to generalize even better

Some suggest that these large models have more capacity to build internal structure and abstract representations that aid generalization

Act 2: Supervised Learning

Bias is Necessary for Intelligence

- Learning is about adjusting weights to minimize loss.
- We want to avoid both underfitting and overfitting.
- SGD is surprisingly effective—it finds generalizable solutions.
- Big models generalize well, contrary to intuition.

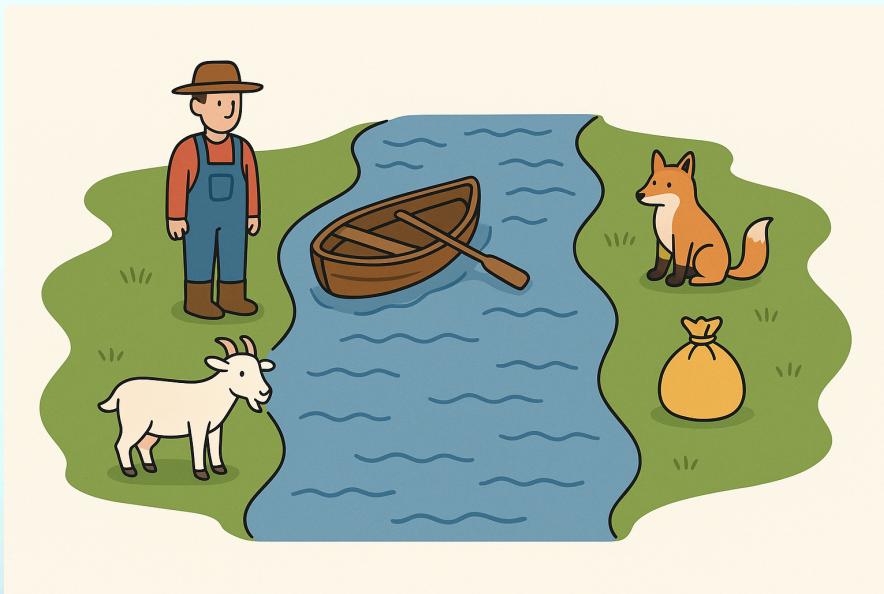
Act 3: Reinforcement Learning

The “Impossible” Reasoning Machine

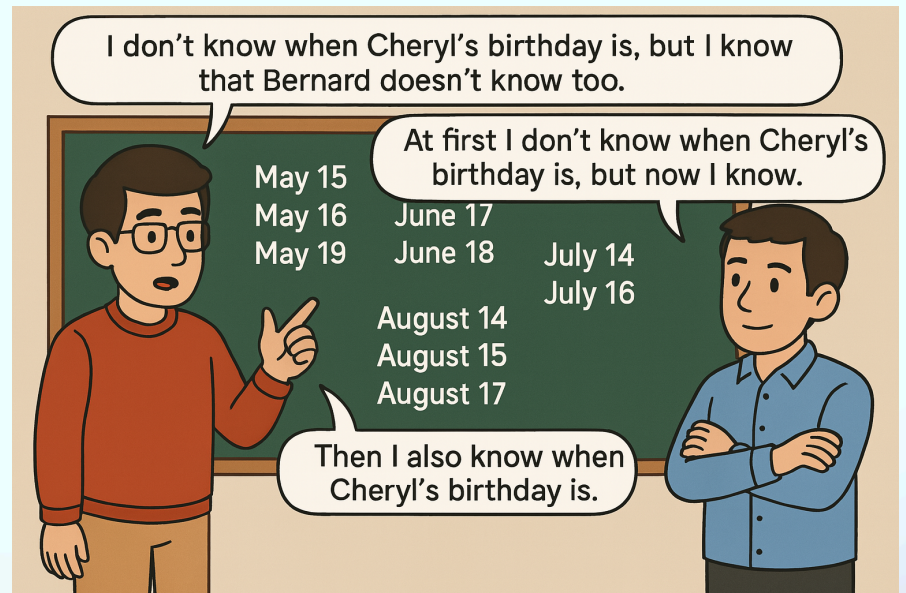
“Impossible”

LLM simply predicts the next token, how can it possibly know how to reason?

River Crossing Puzzle



Birthday Problem



Indeed, LLMs failed on these logic puzzles just a couple of years ago, and some people thought this was impossible.

Discovery: Chain of Thought

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had $23 - 20 = 3$. They bought 6 more apples, so they have $3 + 6 = 9$. The answer is 9. ✅

Credit: *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*

Simply prompting the model to “reason step by step” or providing chain-of-thought examples can dramatically improve its answers.

Discovery: Chain of Thought

I have 3 apples. My dad has 2 more apples than me. How many apples do we have in total?

5 apples. (Greedy Decoding)

I have 3 apples, my dad has 2 more apples than me, so he has 5 apples. $3+5=8$.

We have 8 apples in total.

You have 3 apples, your dad has 2 more apples than you, so he has 5 apples. $3+5=8$.

The answer is 5.

CoT

Stanford

Xuezhi Wang and Denny Zhou. [Chain-of-Thought Reasoning Without Prompting](#). NeurIPS 2024.

Credit: *Large Language Model Reasoning*

In fact, LLMs are “ready” to reason even without prompt engineering.

All we need is decoding. We just need to look at more candidates.

Train a Reasoning Model

Reshape the model output's distribution so that thoughtful responses naturally rank 1st.

1st Attempt: Supervised Fine Tuning with CoT examples created by **human**.

Did not generalize well. Did not scale well.

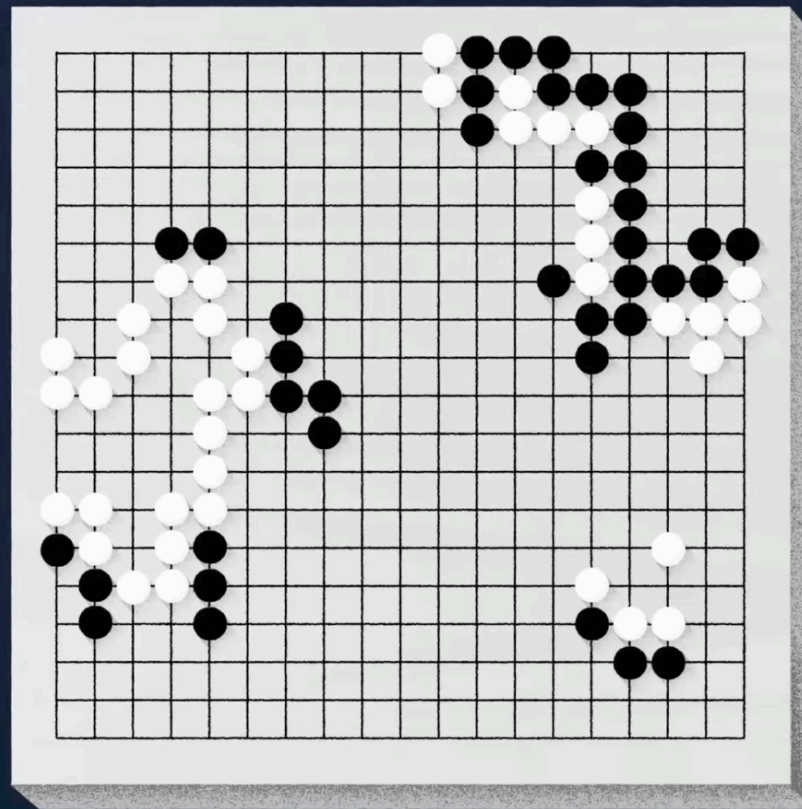
2ed Attempt: Reinforcement Learning with CoT generated by the model.

Train a Reasoning Model

Reshape the model output's distribution so that thoughtful responses naturally rank 1st.

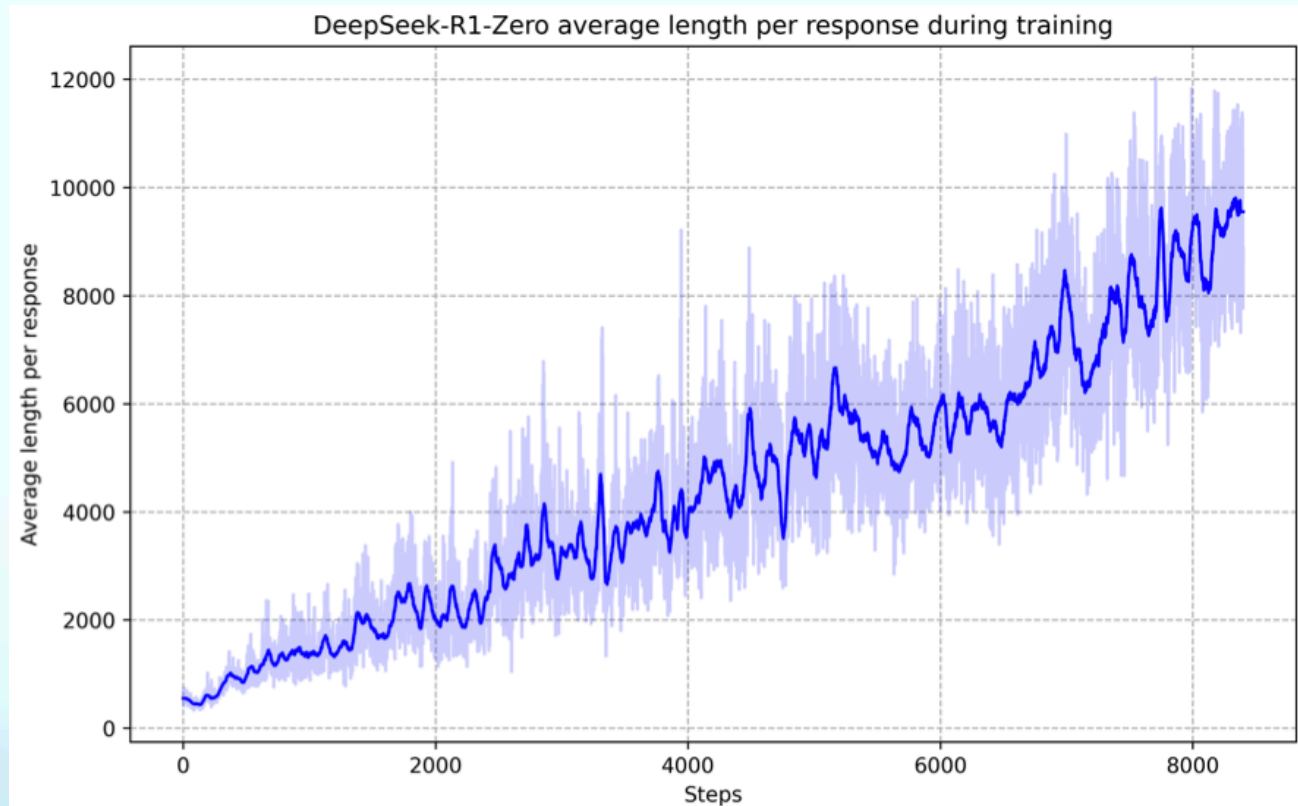
2ed Attempt: Reinforcement Learning with CoT generated by the model.

Similar to how Alpha Zero developed "intuition" about best move sequences.



Train a Reasoning Model

Reshape the model output's distribution so that thoughtful responses naturally rank 1st.



Credit: [DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning](#)

Model naturally learns to solve increasingly complex reasoning tasks by leveraging extended test-time computation.

Train a Reasoning Model

Reshape the model output's distribution so that thoughtful responses naturally rank 1st.

Question: If $a > 1$, then the sum of the real solutions of $\sqrt{a - \sqrt{a+x}} = x$ is equal to

Response: <think>

To solve the equation $\sqrt{a - \sqrt{a+x}} = x$, let's start by squaring both ...

$$\left(\sqrt{a - \sqrt{a+x}}\right)^2 = x^2 \implies a - \sqrt{a+x} = x^2.$$

Rearrange to isolate the inner square root term:

$$(a - x^2)^2 = a + x \implies a^2 - 2ax^2 + (x^2)^2 = a + x \implies x^4 - 2ax^2 - x + (a^2 - a) = 0$$

...

Wait, wait. Wait. That's an aha moment I can flag here.

Let's reevaluate this step-by-step to identify if the correct sum can be ...

We started with the equation:

$$\sqrt{a - \sqrt{a+x}} = x$$

First, let's square both sides:

$$a - \sqrt{a+x} = x^2 \implies \sqrt{a+x} = a - x^2$$

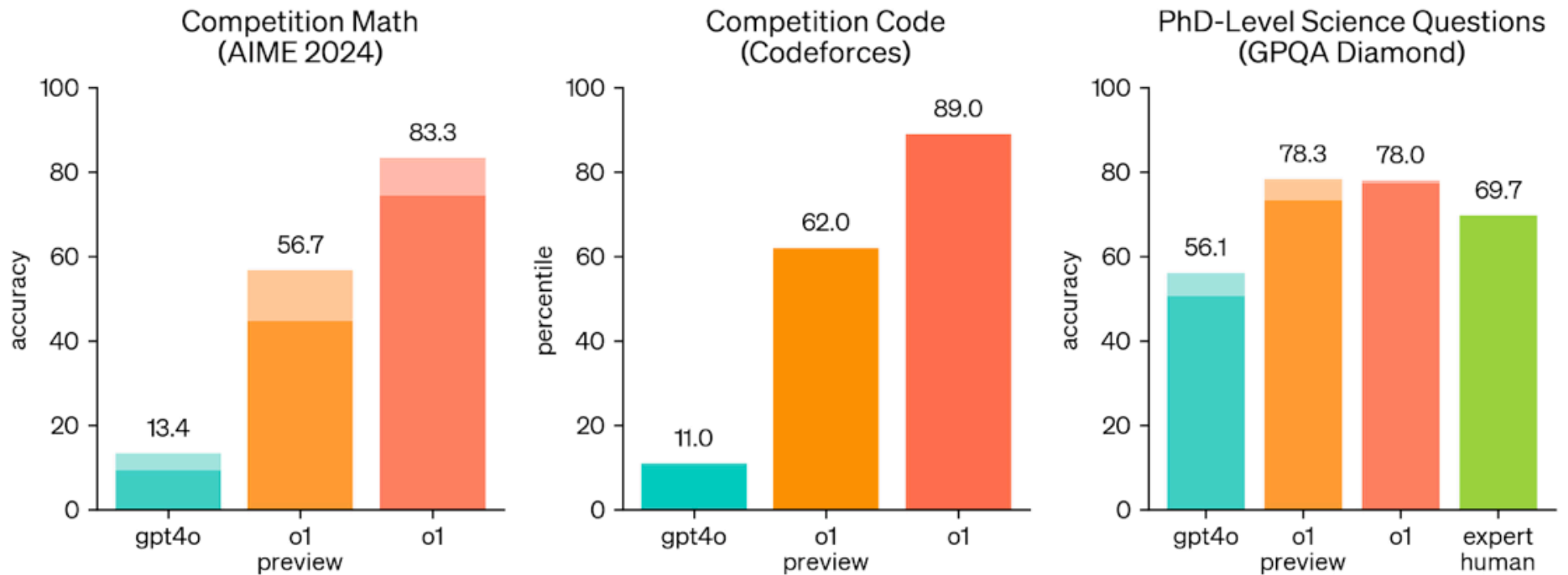
Next, I could square both sides again, treating the equation: ...

...

Credit: [DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning](#)

Behavior such as reflection, exploration of alternative approaches and "aha moment" arise spontaneously.

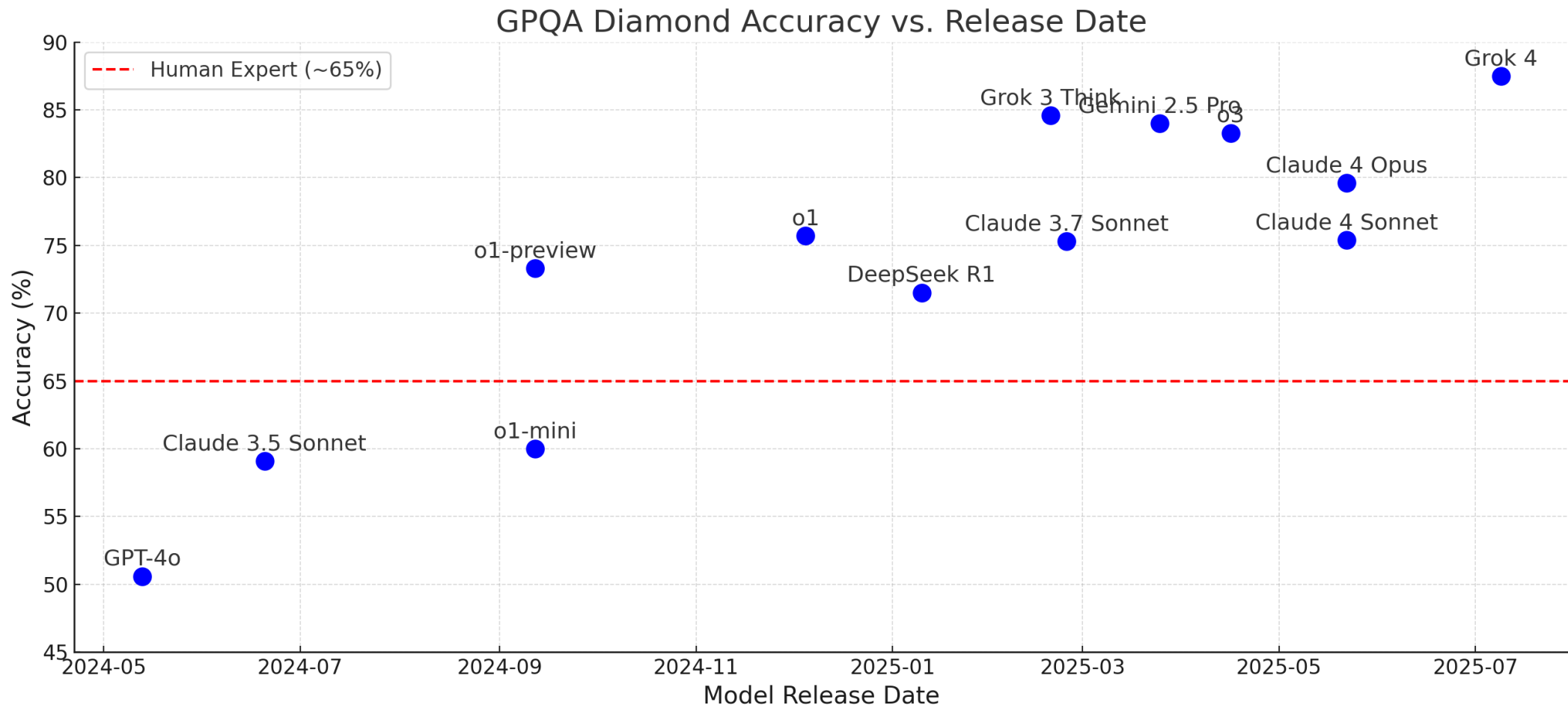
Train a Reasoning Model



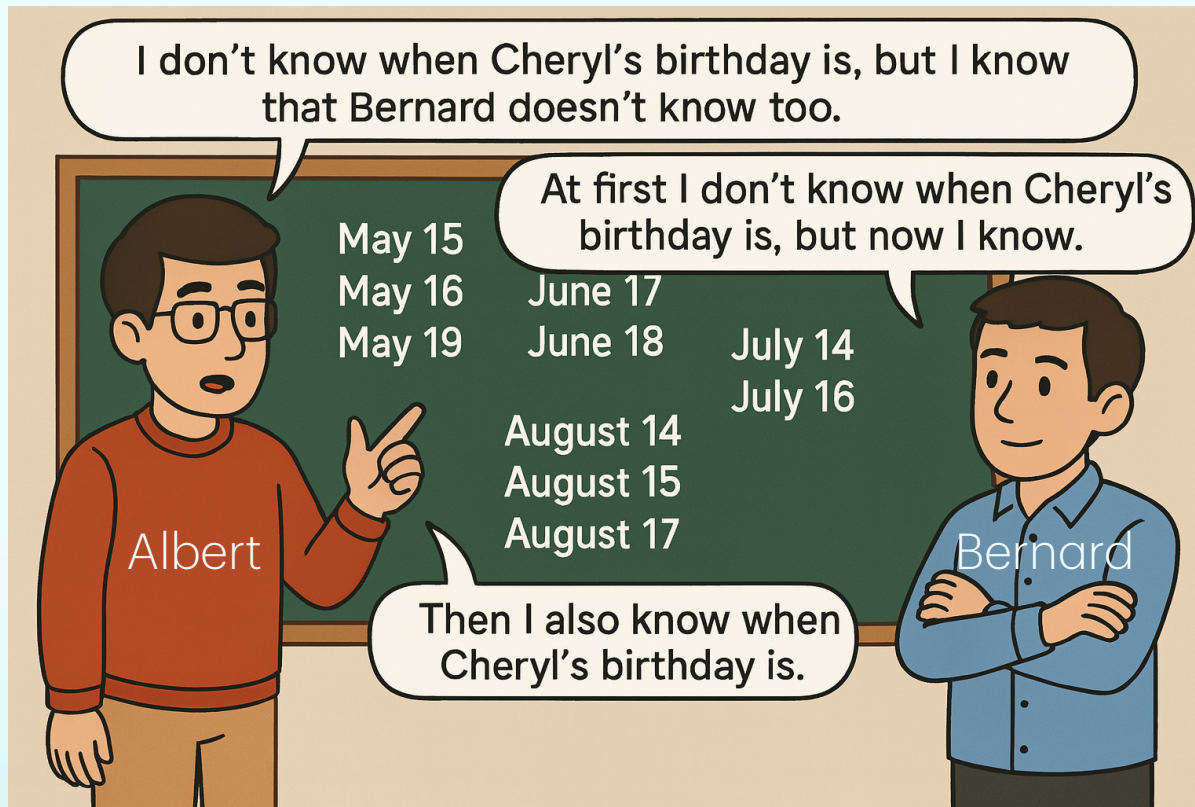
Credit: [Learning to Reason with LLMs \(OpenAI Blog Post\)](#)

Train a Reasoning Model

Graduate-Level Google-Proof Q&A Benchmark.



It's NOT just memorization



May		15	16			19
June				17	18	
July	14		16			
August	14	15		17		

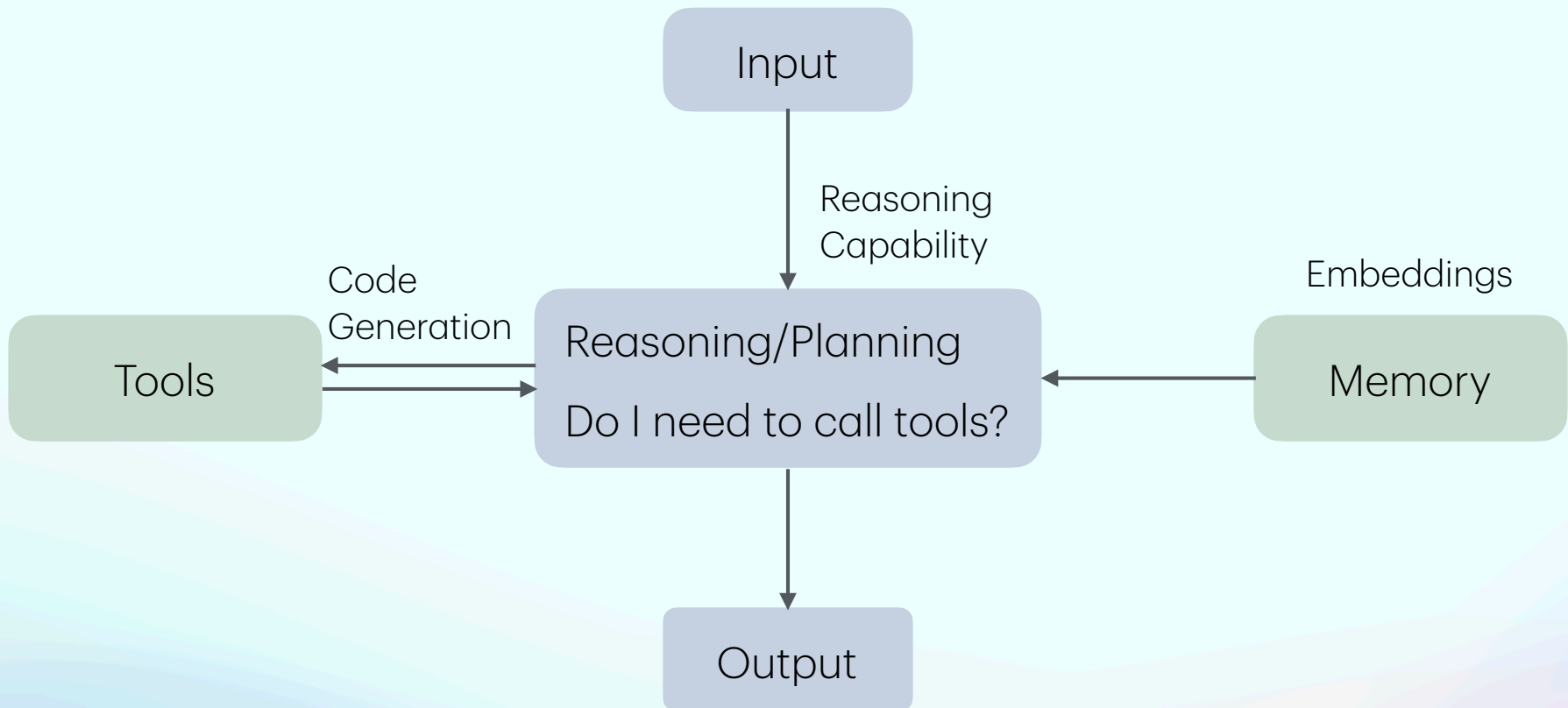
Albert only knows the month.
Bernard only knows the day.

The model can solve countless variations of this puzzle that are not in its training data.

It seems to have learned a deeper logic pattern for solving this puzzle.

AI Agent

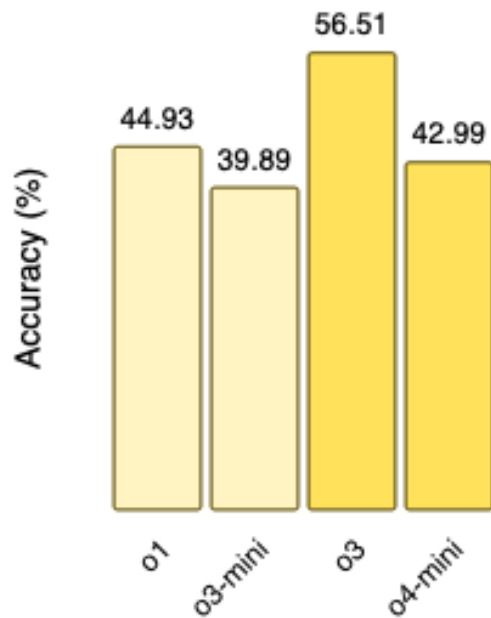
Reasoning capability makes more reliable AI Agent



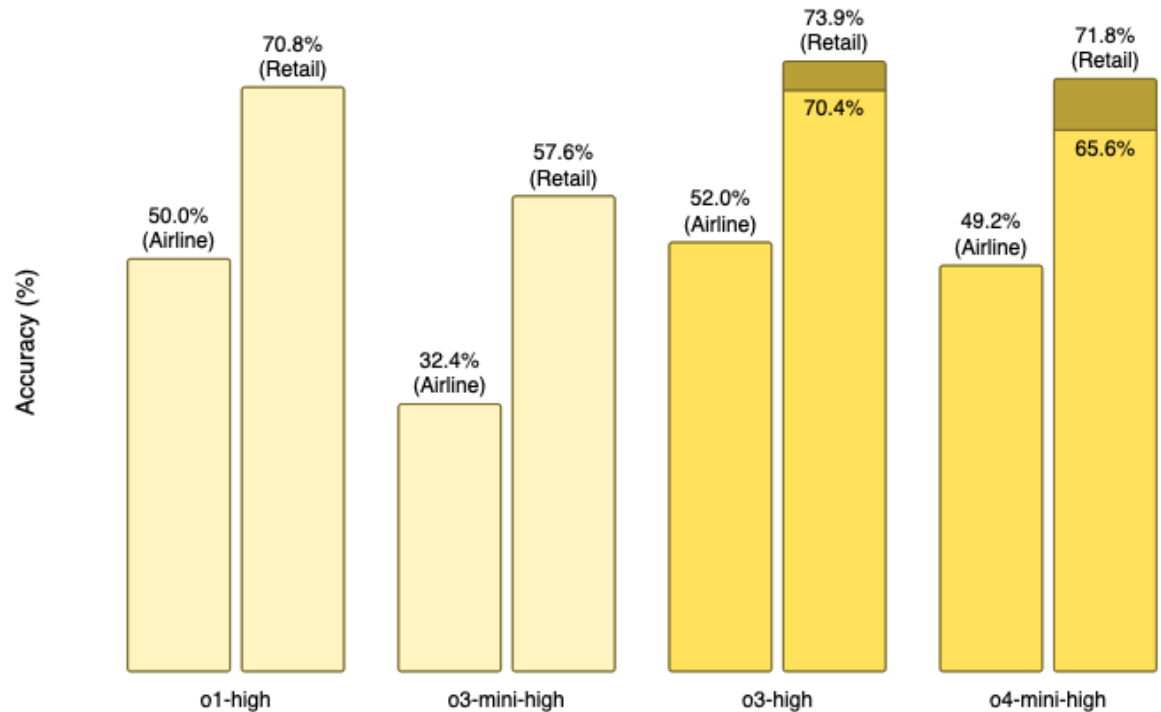
With AI agents becoming more capable and reliable, 2025 marks a year of accelerated AI adoption across diverse enterprise operations.

AI Agent

Scale MultiChallenge
Multi-turn instruction following



Tau-bench
Function Calling



Credit: [Introducing OpenAI o3 and o4-mini](#)

Instruction following and tool use still have **lots of room for improvement**.
In my opinion, this where the frontier models should spent more resources on.

Act 3: Reinforcement Learning

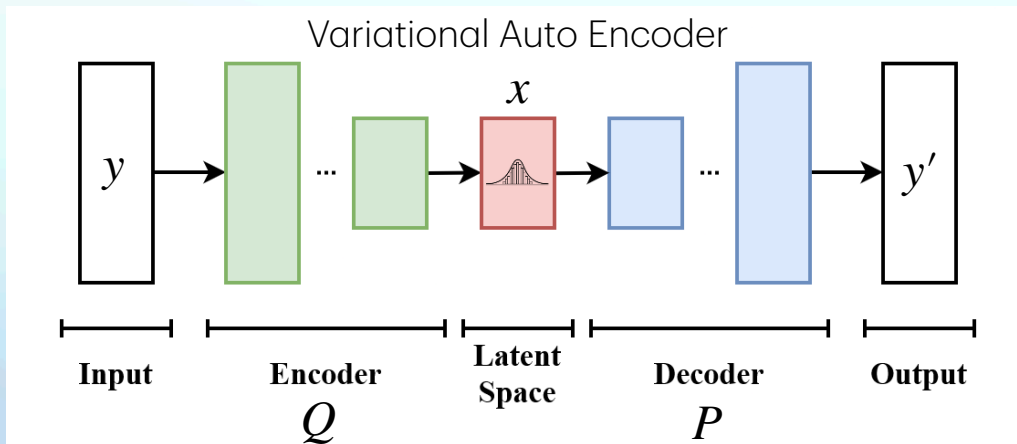
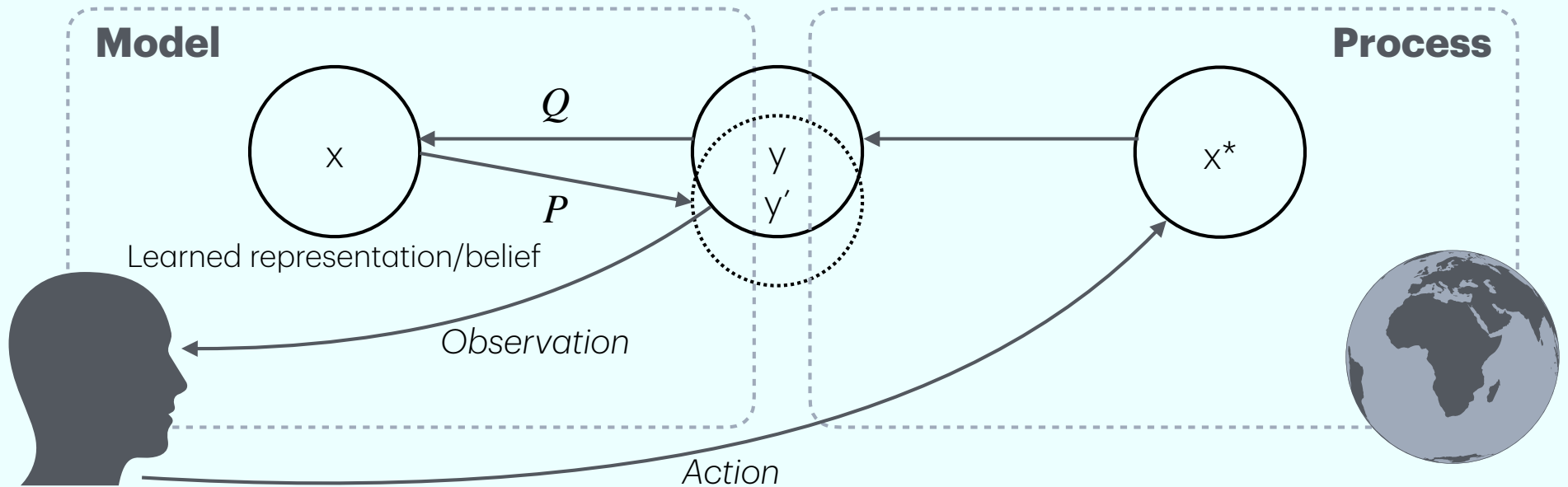
The “Impossible” Reasoning Machine

- CoT and reinforcement learning enabled reasoning in LLMs.
- Cogintion resembling behaviors emerged from optimization.
- LLMs can now tackle logic puzzles that once seemed impossible.
- AI Agents and robotics are the next frontier.

Act 4: Minds vs Machines

“Intelligence” without Consciousness

Variational Free Energy Principle



$$\begin{aligned}
 F[Q, y] &= \underbrace{-\mathbb{E}_Q [\ln P(y, x)]}_{\text{Energy}} - \underbrace{H[Q(x|y)]}_{\text{Entropy}} \\
 &= \underbrace{D_{\text{KL}} [Q(x|y) \parallel P(x|y)]}_{\text{Divergence}} - \underbrace{\ln P(y)}_{\text{Surprise}}
 \end{aligned}$$

According to this framework:
 Consciousness emerges through an internal process of active inference.

Consciousness May Emerge as AI Develops



On the surface ... we are just learning statistical correlations in text. But it turns out that ... to compress them really well, what the neural network learns is some representation of the process that produced the text. This text is actually a projection of the world. ... The neural network learns a compressed, abstract, usable representation of that.

Everything that intelligence does can be boiled down to prediction.

The brain is a biological computer, so why can't a digital computer, a digital brain, do the same things?

We should take seriously the idea that intelligence and awareness could emerge from the same principles.

Brain is More Than Computation



“Artificial intelligence” is a misnomer – it is not really intelligence. It is more like ‘artificial cleverness.’

Intelligence involves consciousness... Computer is a very specific mathematical structure. It means computational mathematics. It’s a very limited part of mathematics.

Understanding is not algorithmic. You can’t write a program that is going to understand the way a human being understands.

- Roger Penrose

Roger Penrose’s definition of “intelligence” is deeply tied to **human understanding**, **insight**, and **consciousness**—particularly our ability to grasp non-algorithmic truths, such as those revealed in Gödel’s incompleteness theorems.

Closing Remarks / My Opinion

1. Calling AI statistical models vastly underestimates what's happening inside these models and what they are capable of. AI gains generalizable understanding through efficient representations of abstract concepts and relations. They are not just curve fitting, they are discovering latent structures in the world.
2. Beyond the two views — that consciousness arises from/with intelligence, or that true intelligence requires consciousness — there's a third possibility:

Perhaps the two aren't deeply connected at all. AI may show us that superintelligence doesn't need consciousness.

A Bold Claim:

Intelligence is solved.

The only mystery is consciousness.