

# **Chain of Thought and RL**

# Scratchpads

- Teaching the model to provide rationale (e.g., computations) instead of just predicting the final answer.

Input:

2 9 + 5 7

Target:

<scratch>

2 9 + 5 7 , C: 0

2 + 5 , 6 C: 1 # added 9 + 7 = 6 carry 1

, 8 6 C: 0 # added 2 + 5 + 1 = 8 carry 0

0 8 6

</scratch>

8 6

Figure 2: Example of input and target for addition with a scratchpad. The carry is recorded in the digit following “C:”. Comments (marked by #) are added for clarity and are not part of the target.

SHOW YOUR WORK: SCRATCHPADS FOR INTERMEDIATE COMPUTATION WITH LANGUAGE MODELS

Maxwell Nye<sup>1,2\*</sup> Anders Johan Andreassen<sup>3</sup> Guy Gur-Ari<sup>3</sup> Henryk Michalewski<sup>2</sup>

Jacob Austin<sup>2</sup> David Bieber<sup>2</sup> David Dohan<sup>2</sup> Aitor Lewkowycz<sup>3</sup> Maarten Bosma<sup>2</sup>

David Luan<sup>2</sup> Charles Sutton<sup>2</sup> Augustus Odena<sup>2</sup>

<sup>1</sup>MIT

<sup>2</sup>Google Research, Brain Team

<sup>3</sup>Google Research, Blueshift Team

# Chain of Thought

Takeshi Kojima  
The University of Tokyo  
t.kojima@weblab.t.u-tokyo.ac.jp

Shixiang Shane Gu  
Google Research, Brain Team

Machel Reid  
Google Research\*

Yutaka Matsuo  
The University of Tokyo

Yusuke Iwasawa  
The University of Tokyo

- LLMs can generate rationale without explicit training:
  - “Think step by step”
  - Few-shot CoT

(a) Few-shot

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The answer is 8. ✗

(b) Few-shot-CoT

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: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) The juggler can juggle 16 balls. Half of the balls are golf balls. So there are  $16 / 2 = 8$  golf balls. Half of the golf balls are blue. So there are  $8 / 2 = 4$  blue golf balls. The answer is 4. ✓

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) 8 ✗

(d) Zero-shot-CoT (Ours)

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: **Let's think step by step.**

(Output) There are 16 balls in total. Half of the balls are golf balls. That means that there are 8 golf balls. Half of the golf balls are blue. That means that there are 4 blue golf balls. ✓

# Review

- Scratchpad/CoT improves both the **expressivity** and **learning**.
  - With a granular scratchpad, any poly-time algorithm can be taught to a LM.
  - We should keep the **globality** of our scratchpads low.
- Length generalization and its relation to positional embeddings.
  - Relative positional embeddings are often better than absolute ones.
  - However, they're not enough and other ideas such as **attention masking** should be used.

# Tool calling

Timo Schick Jane Dwivedi-Yu Roberto Dessì<sup>†</sup> Roberta Raileanu  
 Maria Lomeli Luke Zettlemoyer Nicola Cancedda Thomas Scialom  
 Meta AI Research <sup>†</sup>Universitat Pompeu Fabra

- Addition can be hard to learn, it's also not very efficient
- Let's just call a calculator!
- At training, we only need to exclude the output of the tool from the next-token prediction loss.
- At inference, we call the tool.
- Can increase **efficiency** and **reliability**.

---

The Nile has an approximate length of `<API> QA(What is the approximate length of the Nile?)`  
`→ 6,853 km </API>` 6,853 kilometers, the White Nile being its main source.

---

If Venus had an atmosphere similar to Earth's then you would expect Venus' mean temperature to be 499 K (1.74 x 287) rather than 735 K which is `<API> Calculator(735 / 499) → 1.47 </API>` 1.47 (735 / 499) times hotter than it should be.

---

API Name	Example Input	Example Output
Question Answering	Where was the Knights of Columbus founded?	New Haven, Connecticut
Wikipedia Search	Fishing Reel Types	Spin fishing > Spin fishing is distinguished between fly fishing and bait cast fishing by the type of rod and reel used. There are two types of reels used when spin fishing, the open faced reel and the closed faced reel.
Calculator	27 + 4 * 2	35
Calendar	ε	Today is Monday, January 30, 2023.
Machine Translation	sûreté nucléaire	nuclear safety

# Vision Transformers (ViT)

- Transformers can be also used on images.
- Each image is broken into *patches*, and patches are treated as tokens.
  - Often 14x14 or 16x16 patches.
- We often use an encoder-only model.
  - Having *bidirectional attention*.
- Using a classification token: <CLS> token

AN IMAGE IS WORTH 16X16 WORDS:  
TRANSFORMERS FOR IMAGE RECOGNITION AT SCALE

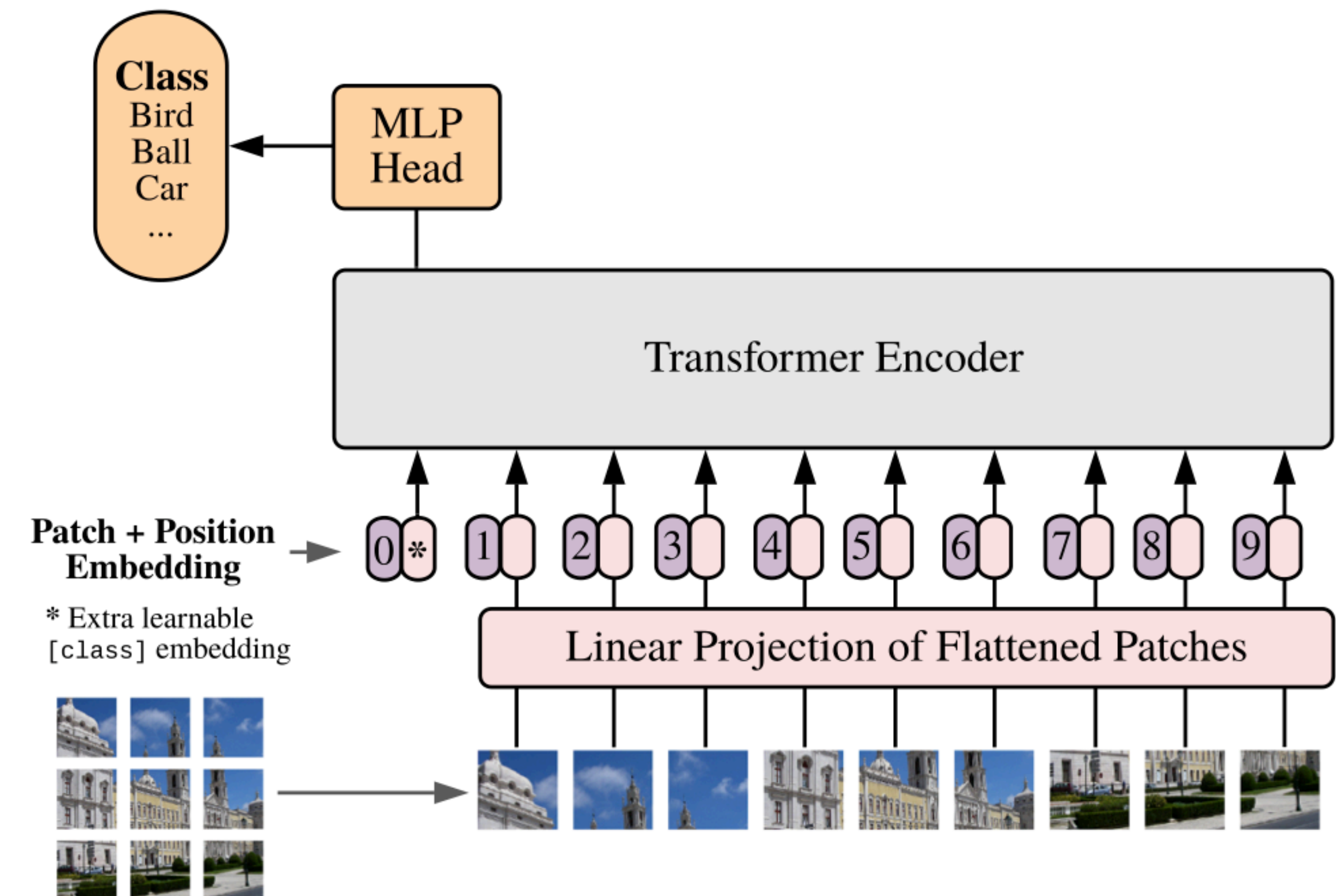
Alexey Dosovitskiy<sup>\*,†</sup>, Lucas Beyer<sup>\*</sup>, Alexander Kolesnikov<sup>\*</sup>, Dirk Weissenborn<sup>\*</sup>,  
Xiaohua Zhai<sup>\*</sup>, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer,  
Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, Neil Houlsby<sup>\*,†</sup>

<sup>\*</sup>equal technical contribution, <sup>†</sup>equal advising

Google Research, Brain Team

{adosovitskiy, neilhousby}@google.com

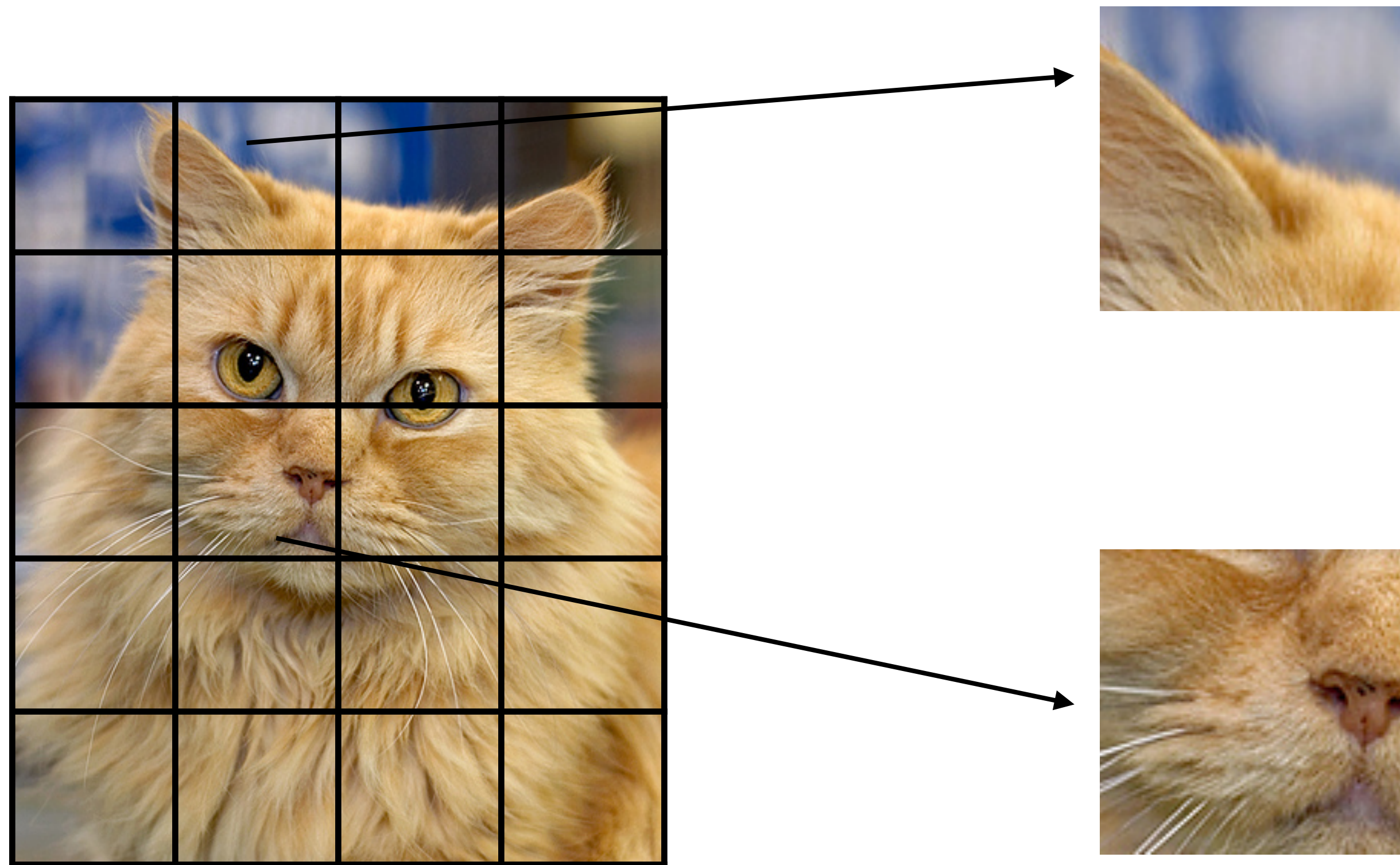
## Vision Transformer (ViT)



# Most vision tasks are local

A few local features (e.g., patches) provide significant information on the output.

=> Common image classification is a *local* task.

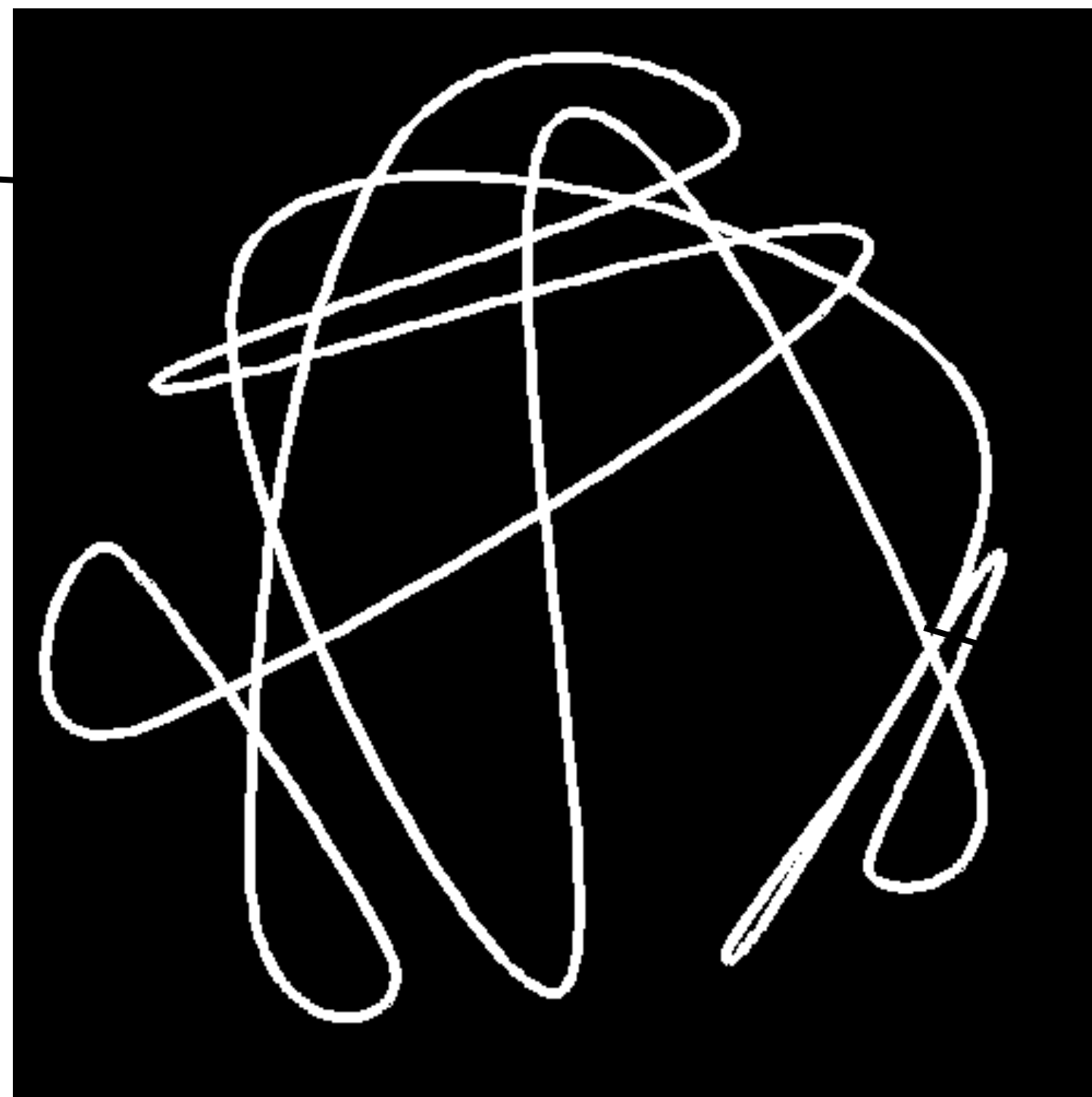


*Probabilities given the two patches*

*Persian cat: 0.38*  
*Egyptian cat: 0.15*  
*tiger cat: 0.08*  
*Siamese cat: 0.02*  
*...*

# Not all visual tasks are local

How many strings are in the image? One or two?



Humans can't solve this task instantaneously either.

=> System II reasoning for vision.

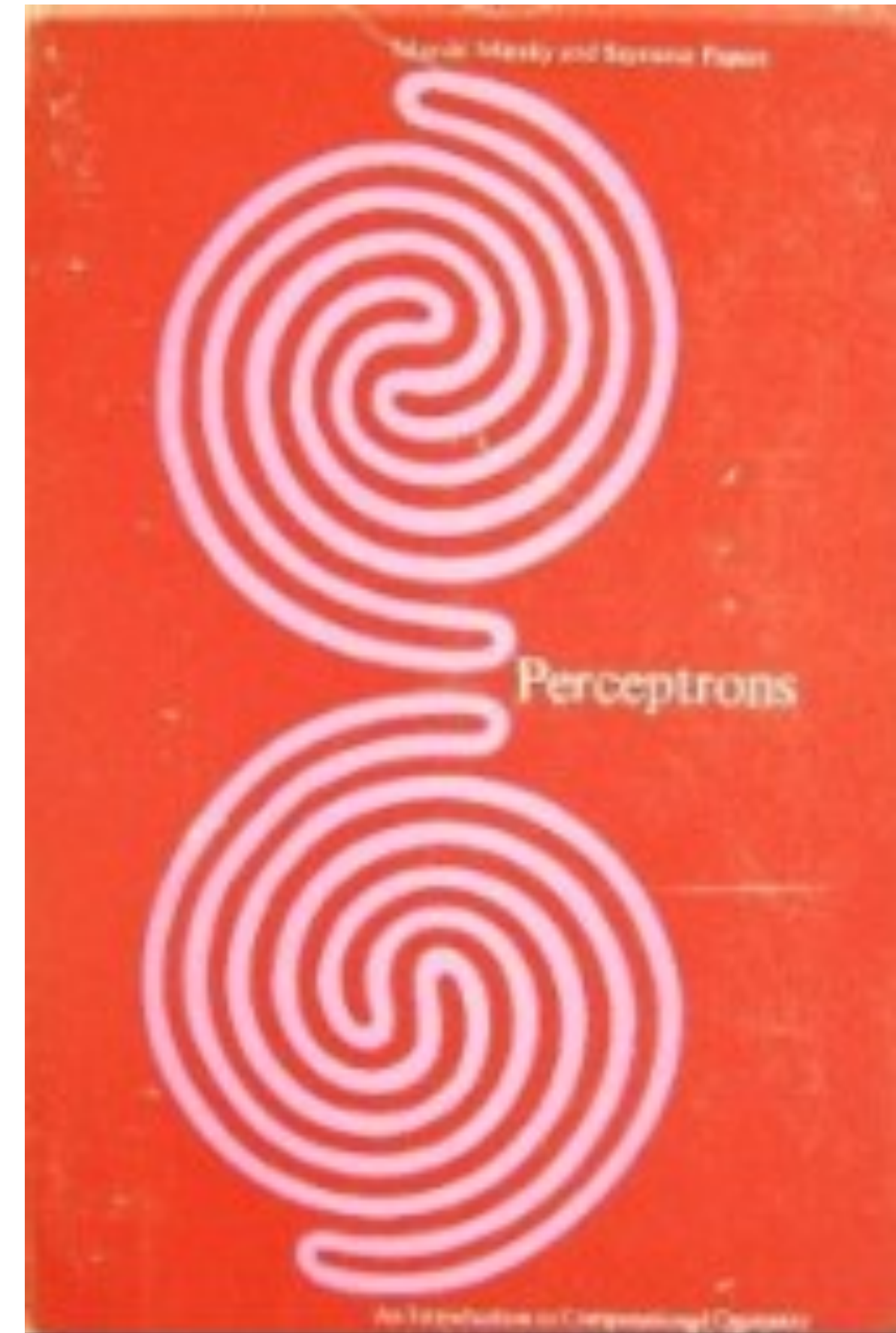
Single patches are uninformative on the output —> global reasoning is needed.



# Not all tasks are local

*“Perceptrons can’t distinguish these figures.”*

- *Current neural networks don’t have expressivity issues.  
But, can they learn to distinguish such figures?*

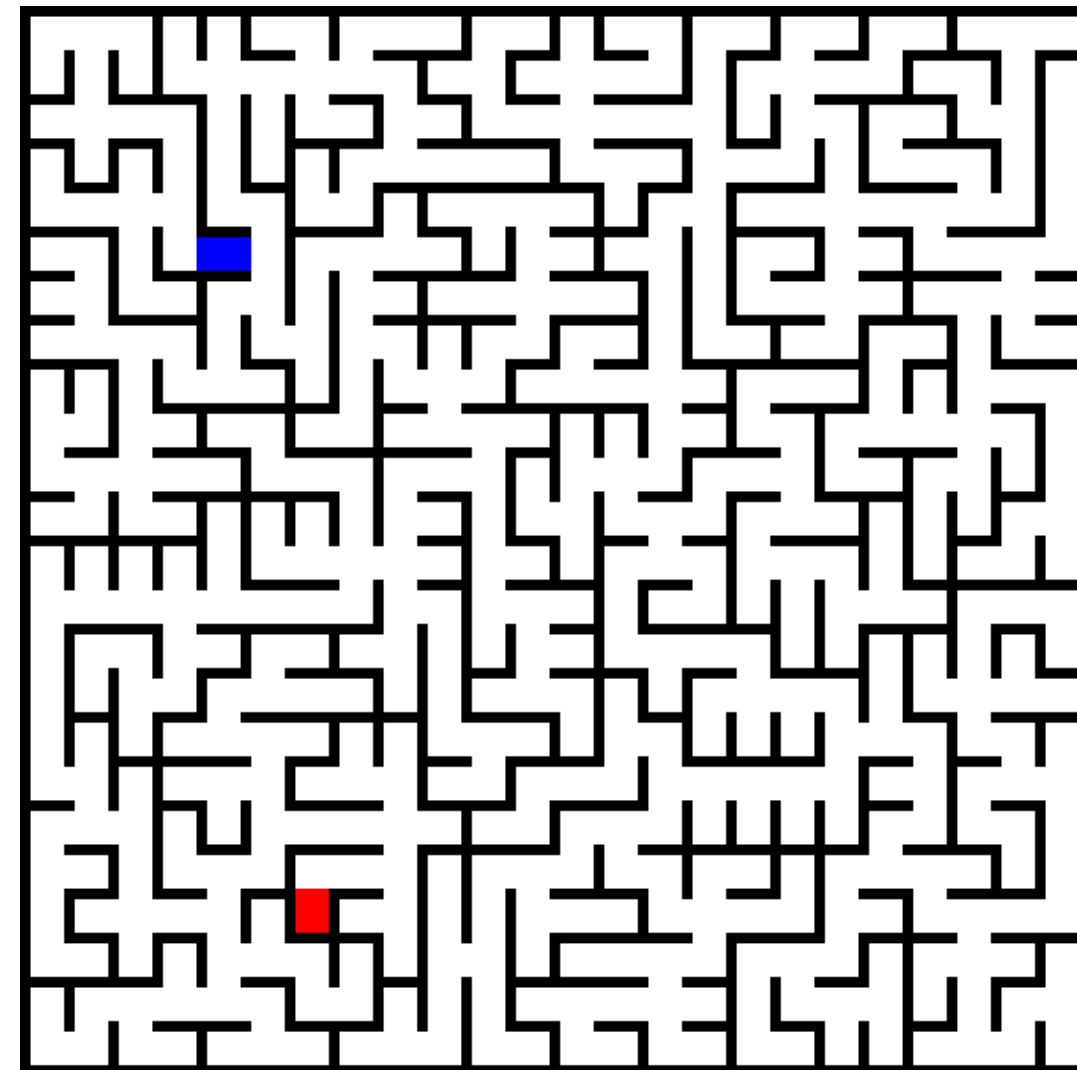


*Perceptrons* by M. Minsky & S. Papert, 1972

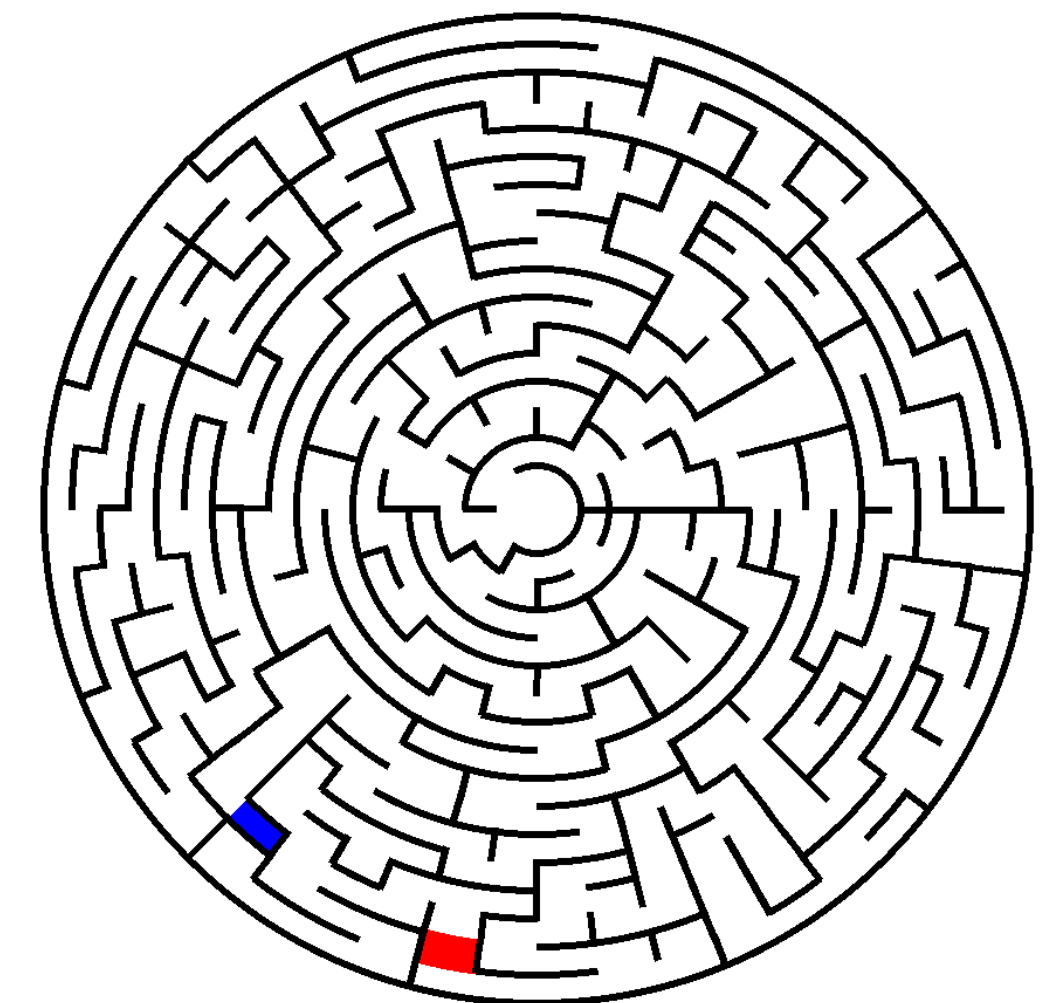
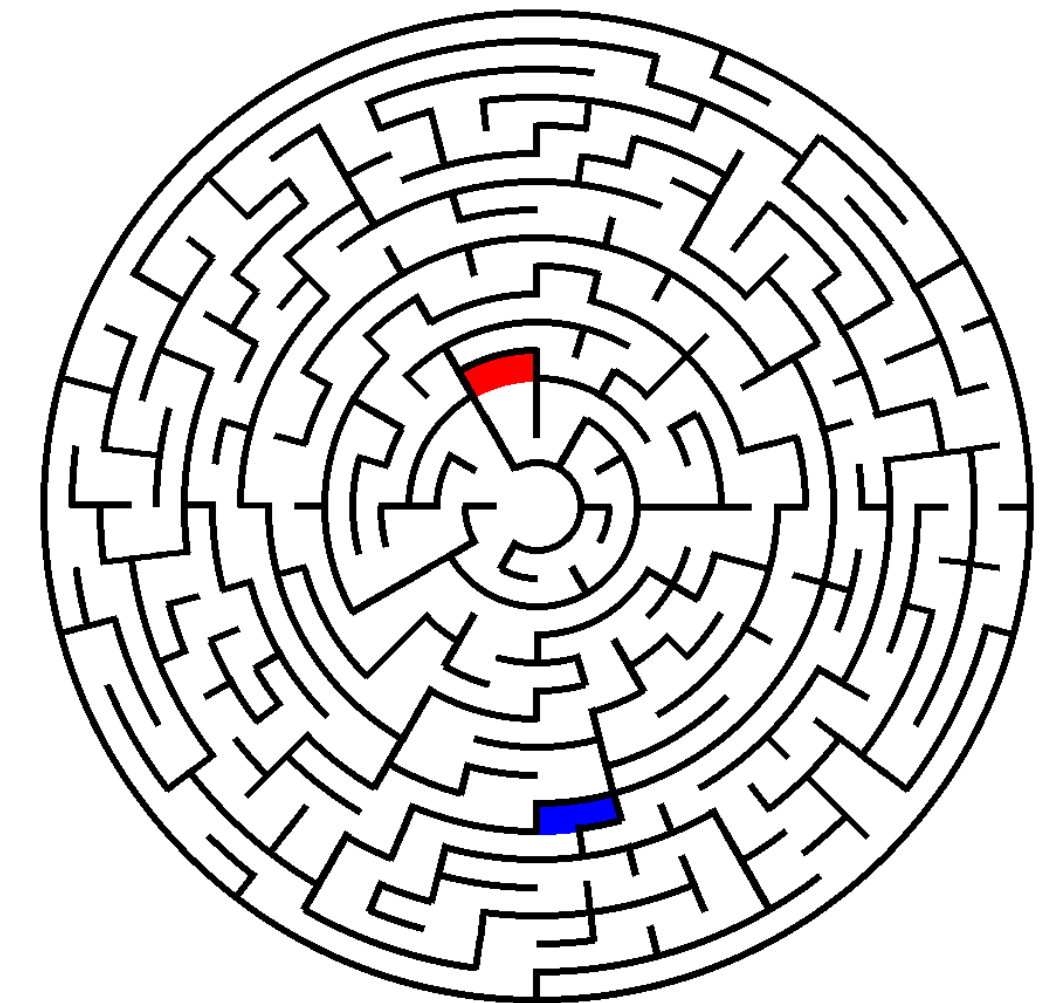
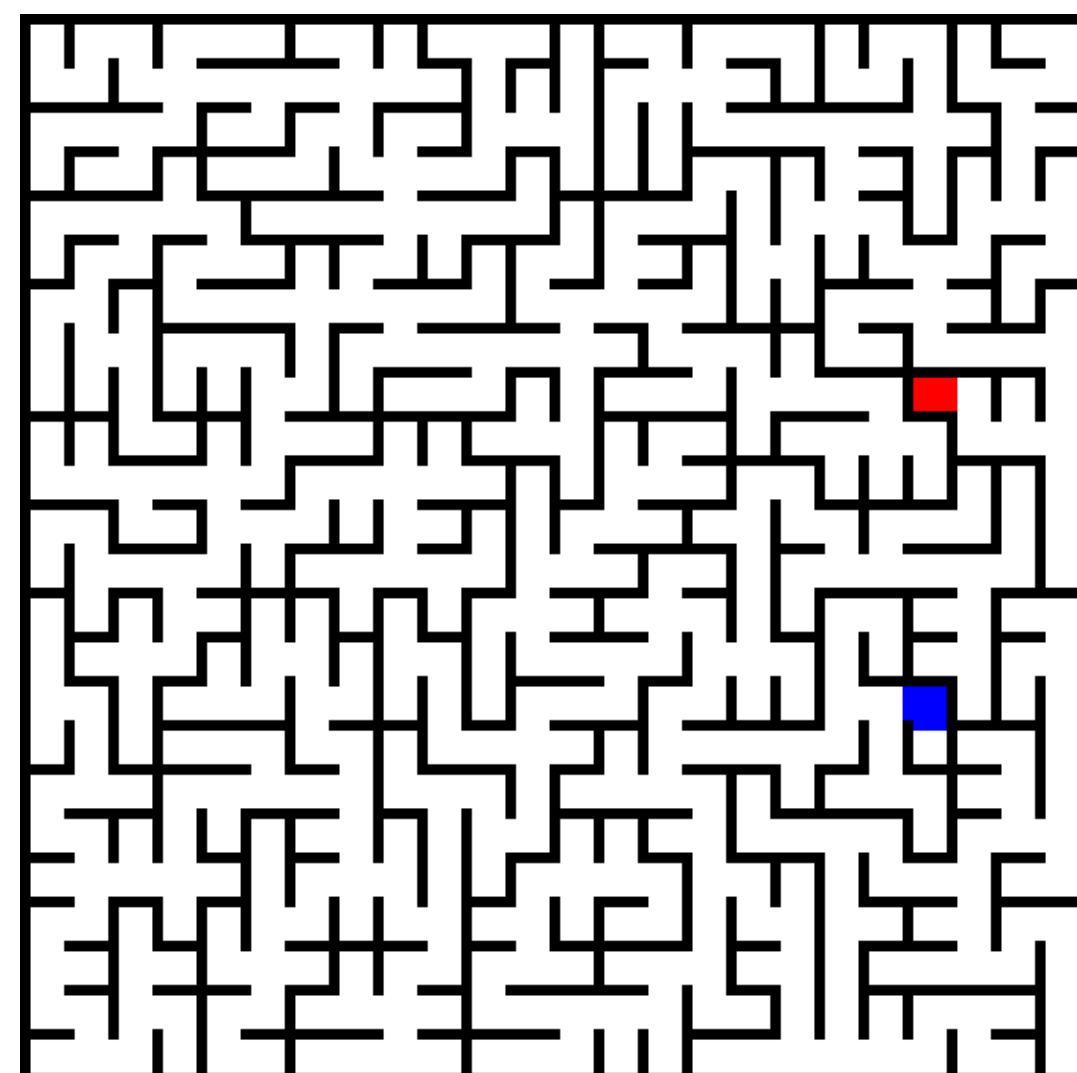
# More examples

Are the two points connected?

Yes



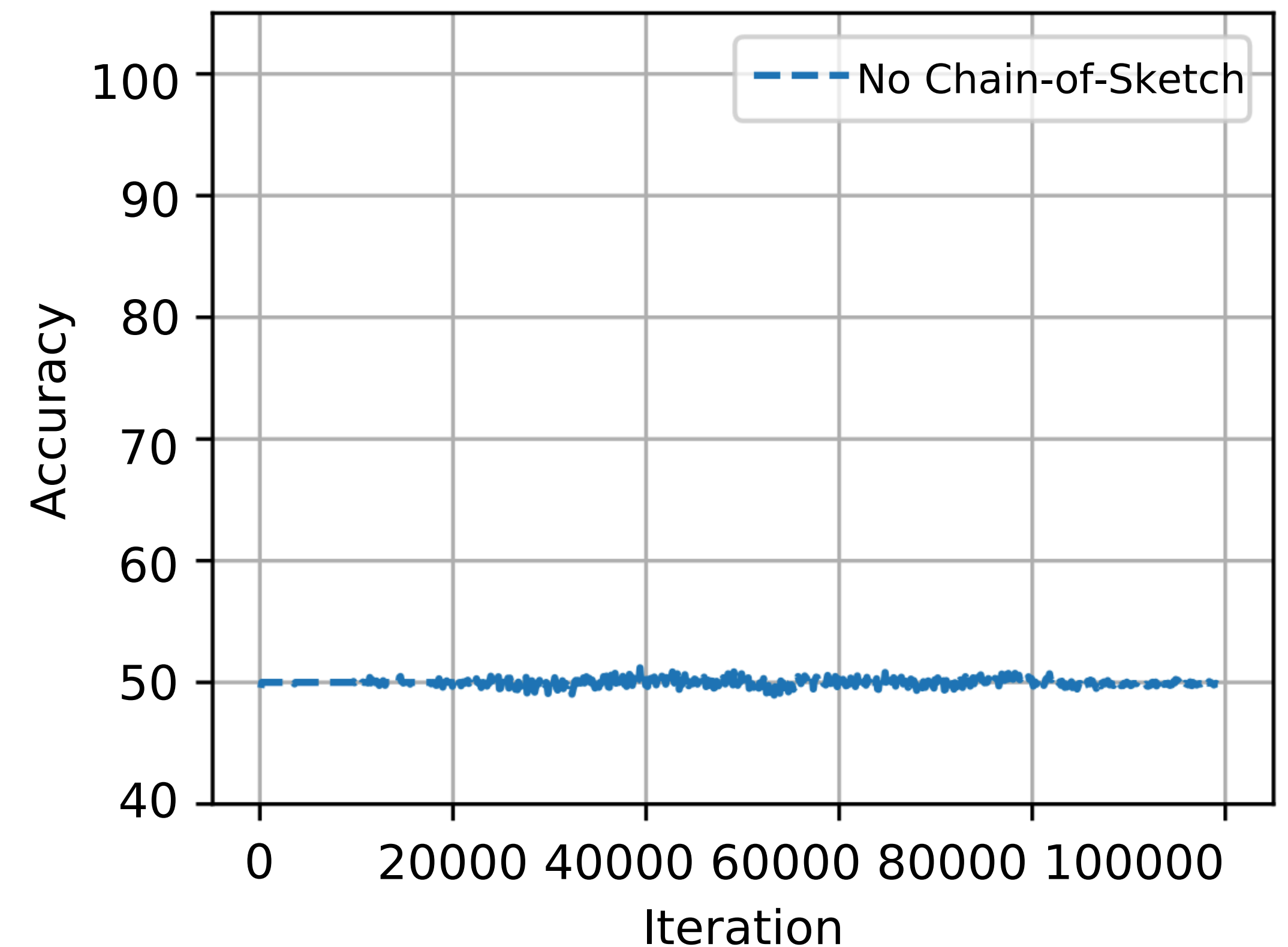
No



# These tasks are not (efficiently) learnable

- E.g., training a ViT with 1M samples of strings task with 16 nodes:

- Using pre-trained models is not helpful.
- Issue: **high globality**

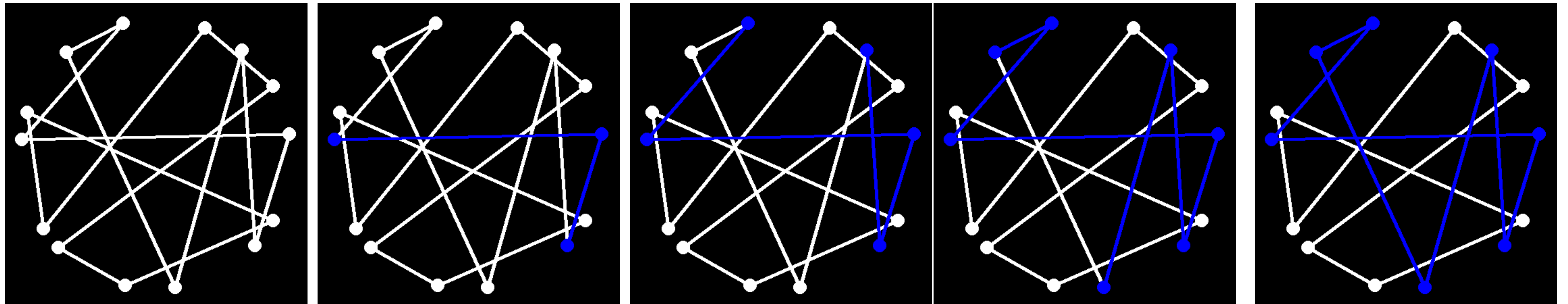


# These tasks are also hard for SOTA multi-modal LLMs.

Task	Size	GPT-4o		o4-mini		o3	
		Prompt	ICL	Prompt	ICL	Prompt	ICL
Cycles	6 nodes	61.9%	58.9%	68.9%	62.2%	63.4%	62.1%
	8 nodes	51.6%	50.4%	55.4%	52.9%	54.7%	51.4%
	10 nodes	49.9%	48.8%	50.5%	48.4%	49.8%	51.3%
Maze (rect.)	4 × 4	50.6%	55.4%	55.3%	52.4%	52.4%	51.9%
	6 × 6	48.4%	51.8%	49.4%	51.3%	51.2%	48.6%

# Solution: what would humans do?

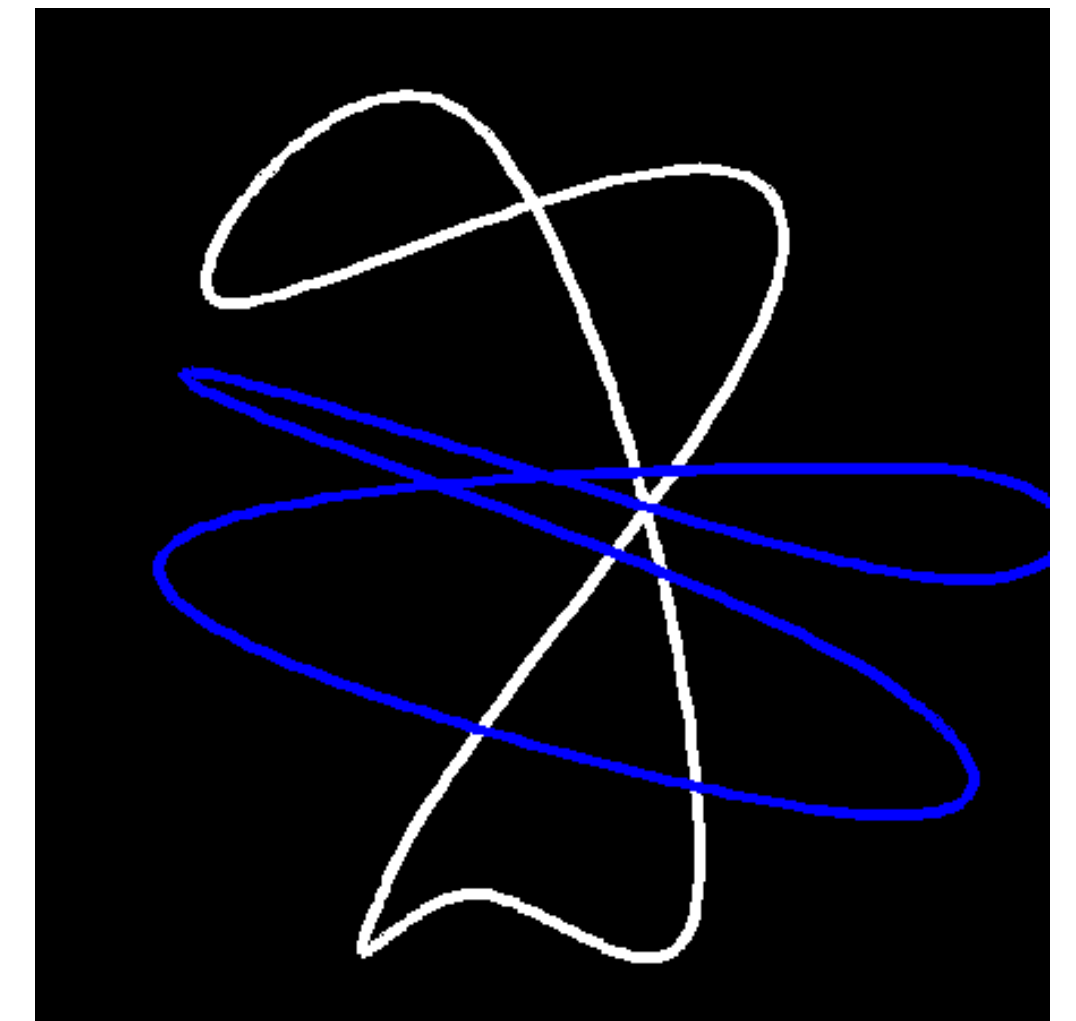
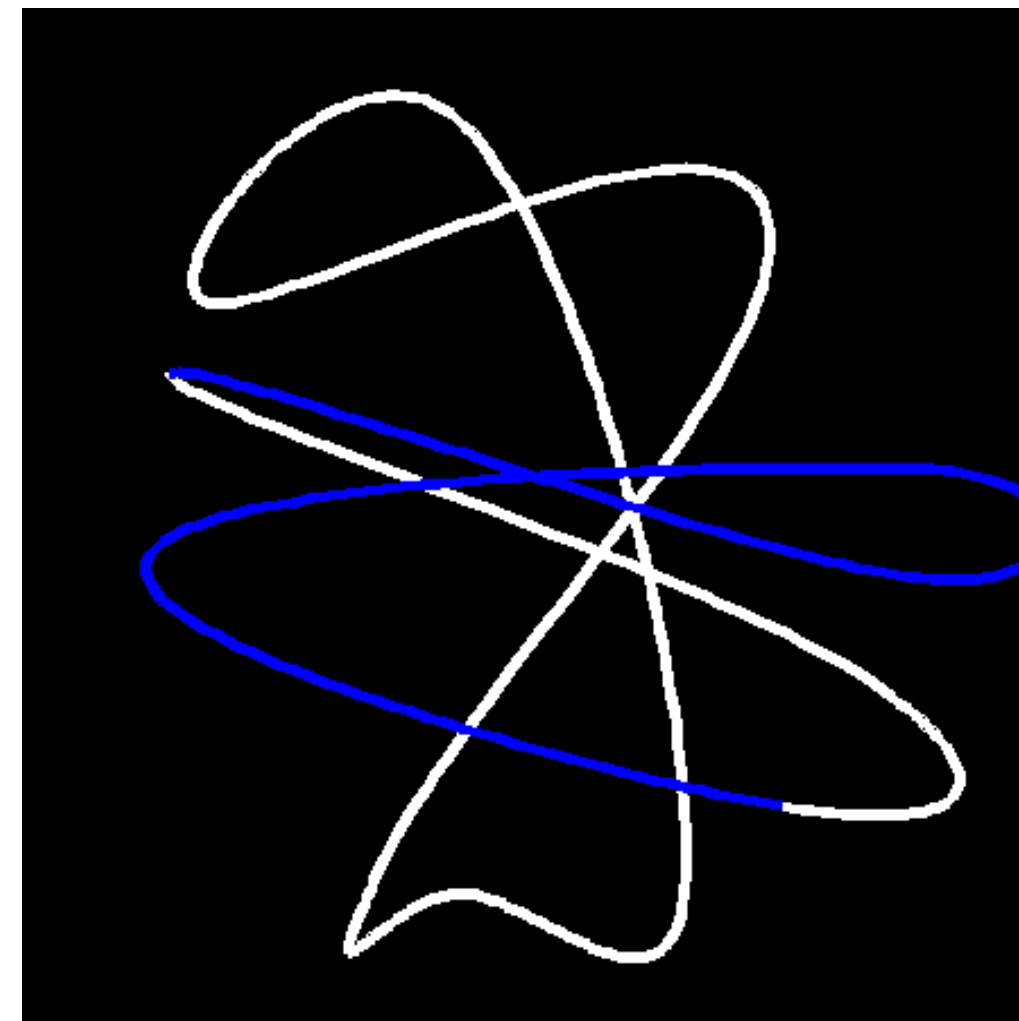
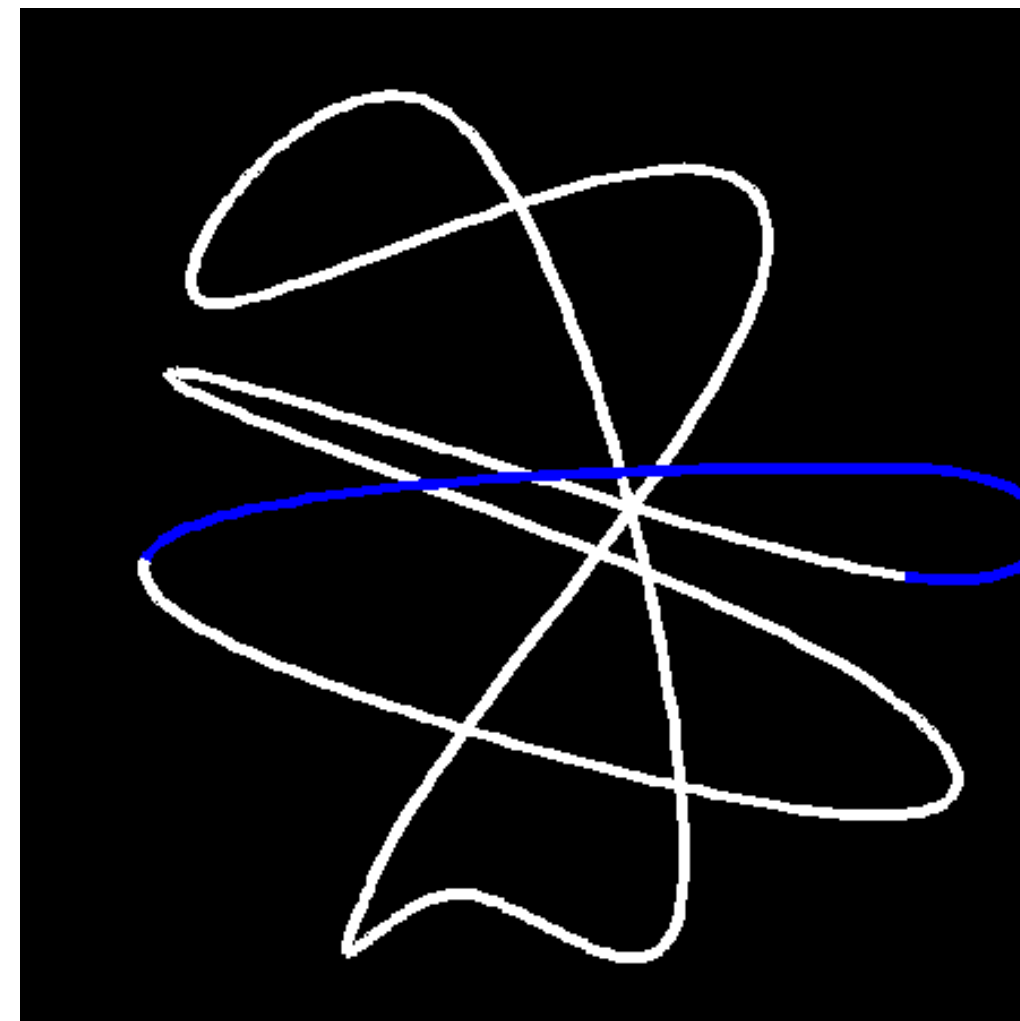
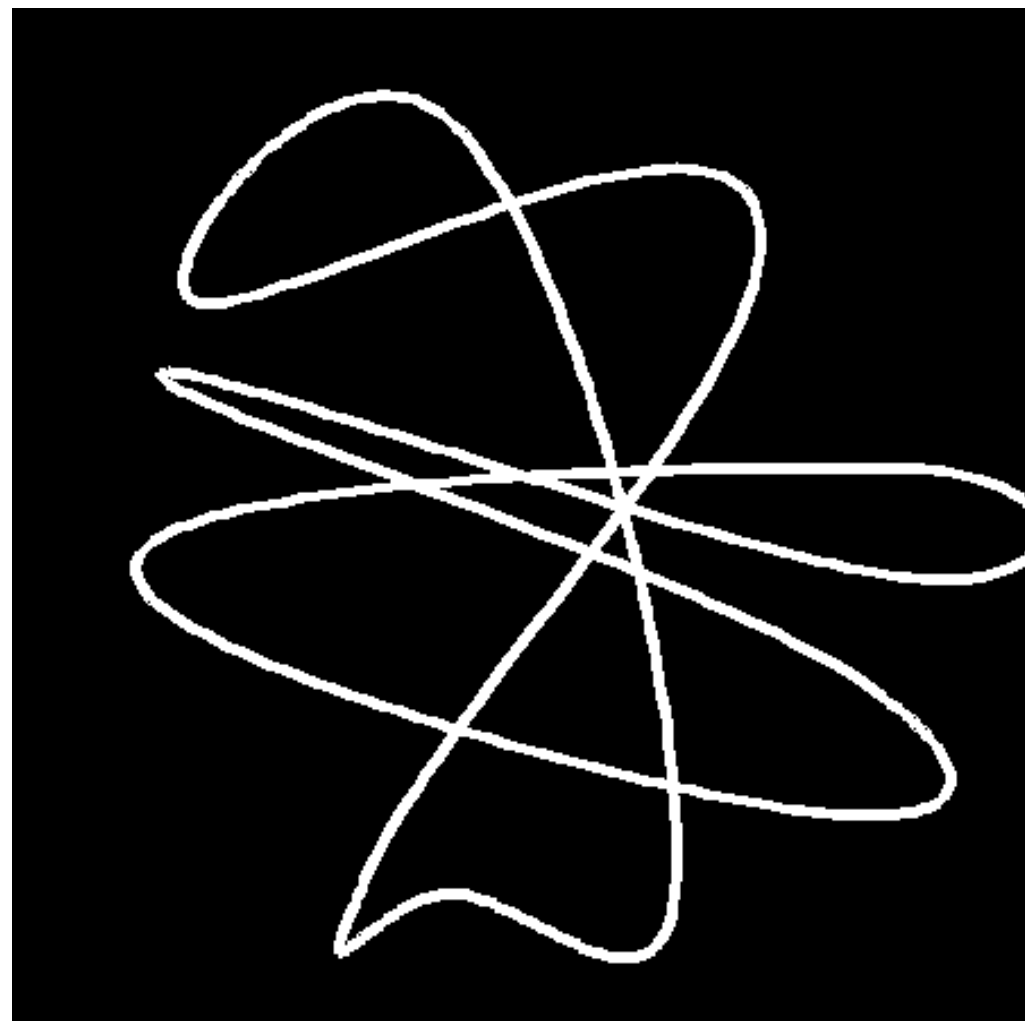
- Classical (and local) vision tasks are often solved by humans instantly and intrinsically.
- Global tasks require thinking  $\rightarrow$  we may need to write something down or annotate the image, etc.



**Chain-of-Sketch**

# Solution: what would humans do?

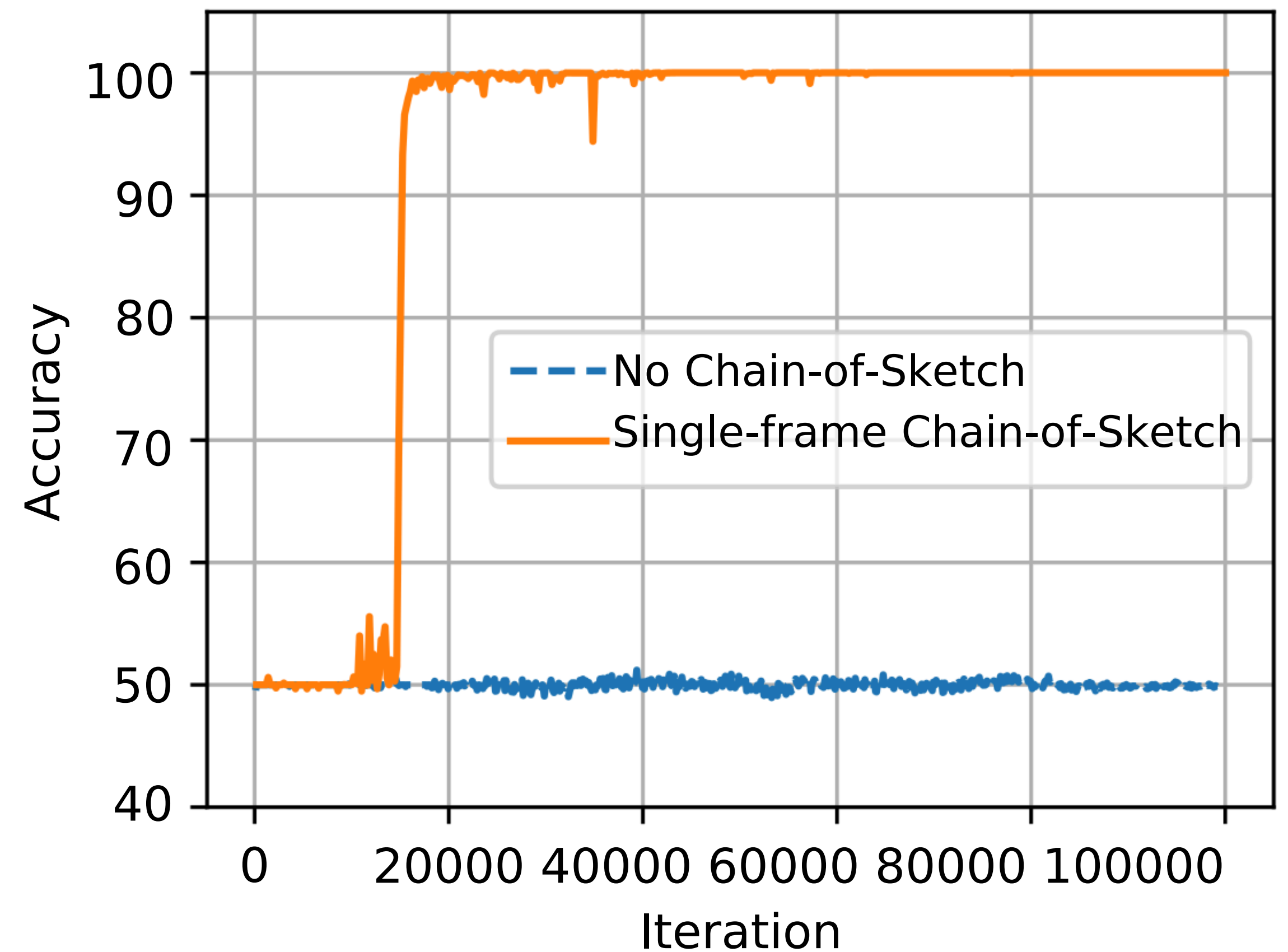
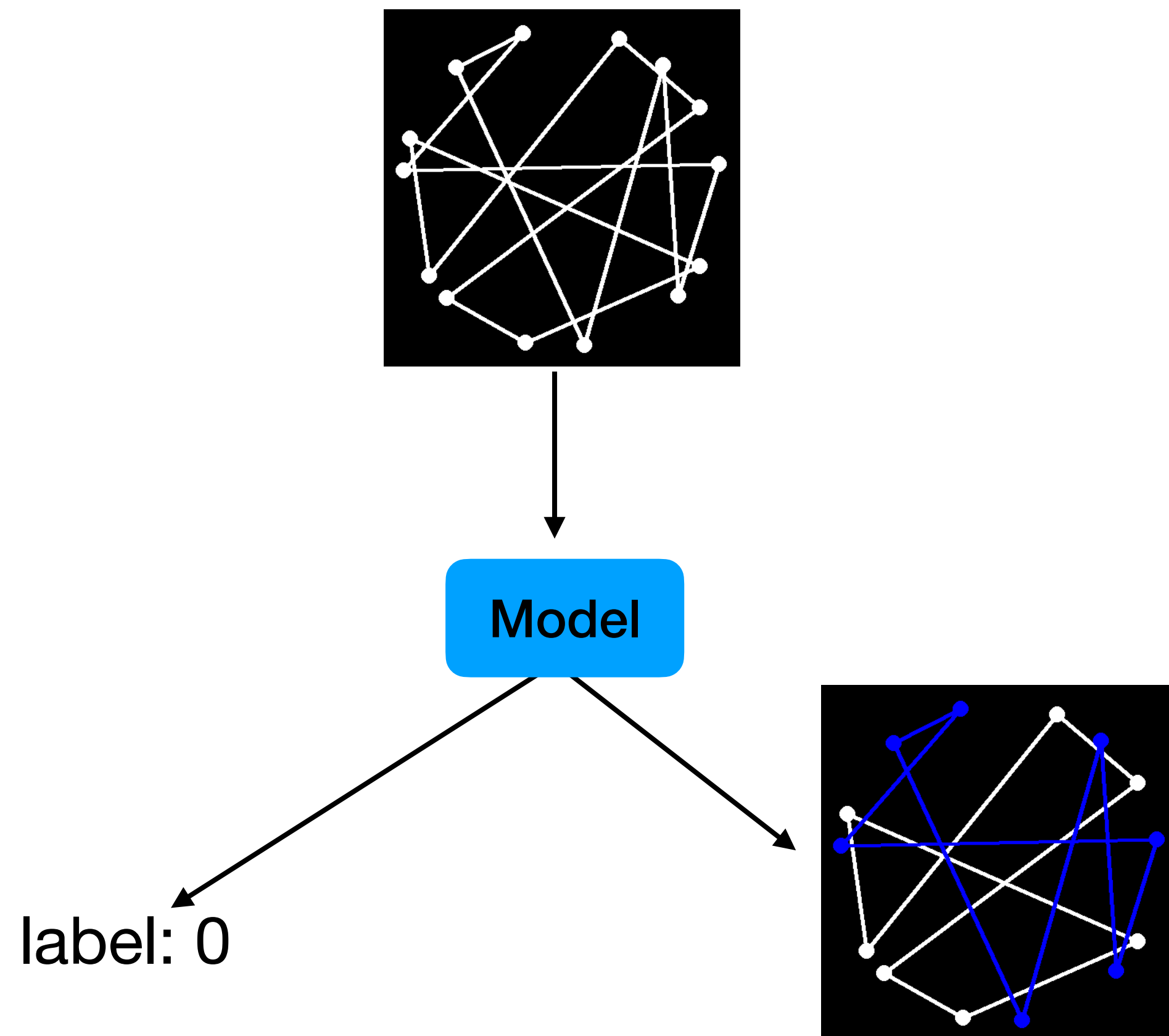
- Classical (and local) vision tasks are often solved by humans instantly and intrinsically.
- Global tasks require thinking —> we may need to write something down or annotate the image, etc.



## Chain-of-Sketch

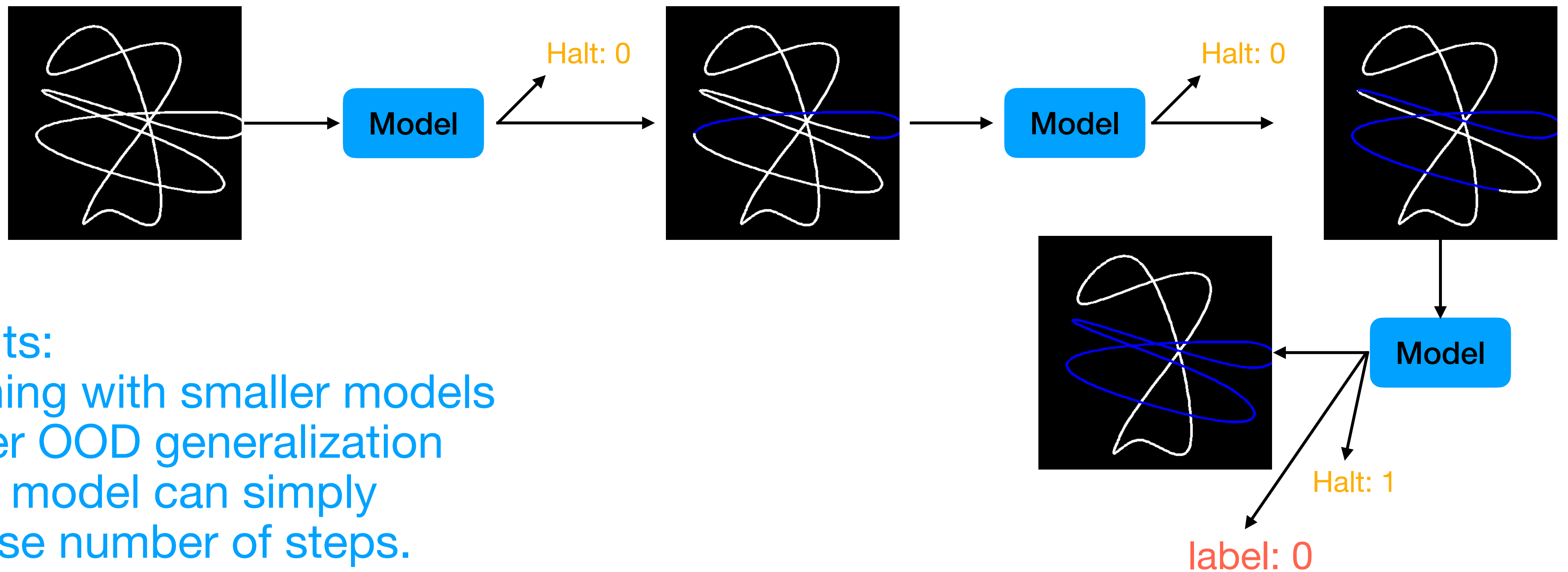
# Single-frame Chain-of-Sketch

Training the model to predict the underlying solution is helpful.



# Inductive Chain-of-sketch

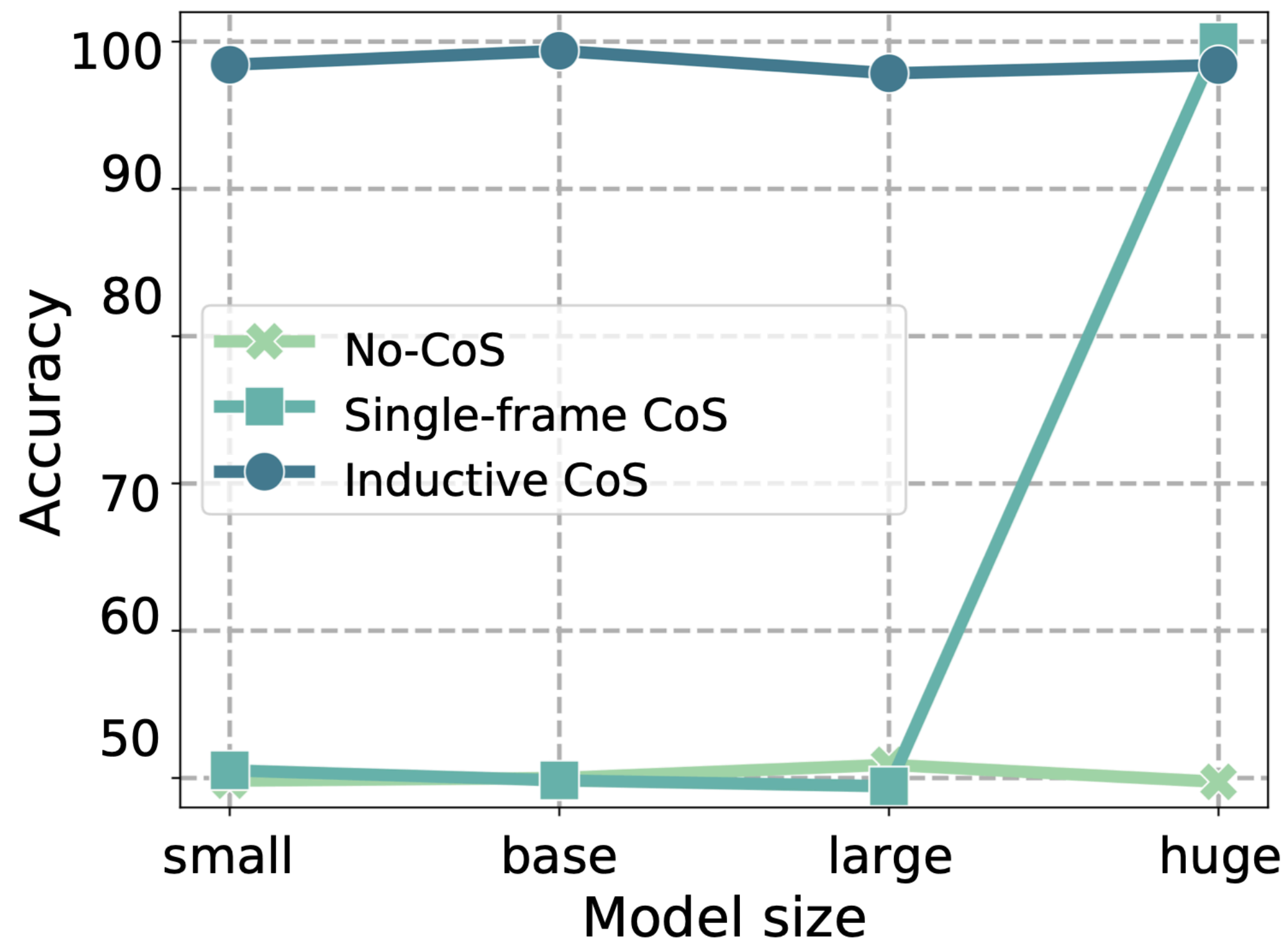
- Training the model to predict the underlying solution in multiple frames.
- Each time the model predicts the next frame using only the previous frame.



## Benefits:

- learning with smaller models
- better OOD generalization as the model can simply increase number of steps.

# Inductive Chain-of-sketch — benefits



Cycles 24

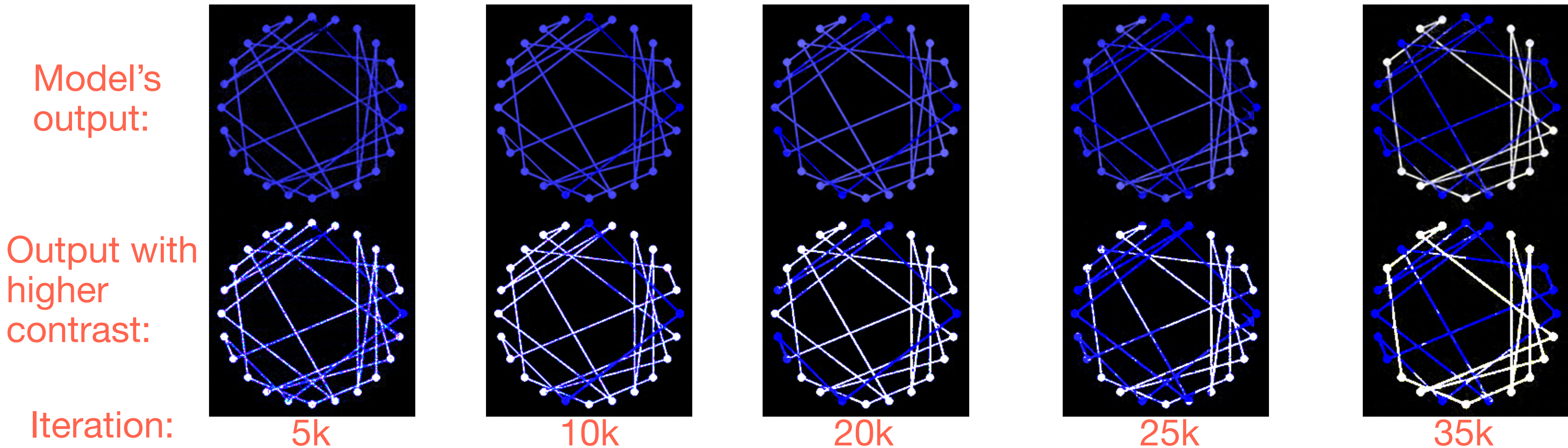
+ Learning with smaller models

Dataset	Method	Accuracy (%)	
		ID	OOD
Maze 24 (Rect.)	single-frame CoS	<b>100.0</b>	54.4
	inductive CoS	99.8	<b>99.8</b>
PVR Grid 7x7	single-frame CoS	<b>100.0</b>	48.6
	inductive CoS	99.5	<b>82.2</b>

+ Better OOD generalization as the model can simply increase number of steps.

# Hierarchical learning of the single-frame CoS

- We observed that the learning process of the scratchpad is hierarchical.



- This is connected to the staircase phenomenon observed/proved when learning Boolean functions.

[E Abbe, E Boix, T Misiakiewicz '22 & '23]

***Takeaway: We need to train models with high quality CoT data to improve their reasoning.***

*How to obtain such data?*

# Focus: Mathematical Reasoning

- There are datasets of math problems and solutions (with steps).

<p><b>Problem:</b> Beth bakes 4, 2 dozen batches of cookies in a week. If these cookies are shared amongst 16 people equally, how many cookies does each person consume?</p> <p><b>Solution:</b> Beth bakes 4 2 dozen batches of cookies for a total of <math>4 \times 2 = \ll 4 \times 2 = 8 \gg</math> 8 dozen cookies There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of <math>12 \times 8 = \ll 12 \times 8 = 96 \gg</math> 96 cookies She splits the 96 cookies equally amongst 16 people so they each eat <math>96/16 = \ll 96/16 = 6 \gg</math> 6 cookies <b>Final Answer:</b> 6</p>
<p><b>Problem:</b> Mrs. Lim milks her cows twice a day. Yesterday morning, she got 68 gallons of milk and in the evening, she got 82 gallons. This morning, she got 18 gallons fewer than she had yesterday morning. After selling some gallons of milk in the afternoon, Mrs. Lim has only 24 gallons left. How much was her revenue for the milk if each gallon costs \$3.50?</p> <p>Mrs. Lim got 68 gallons - 18 gallons = <math>\ll 68 - 18 = 50 \gg</math> 50 gallons this morning. So she was able to get a total of 68 gallons + 82 gallons + 50 gallons = <math>\ll 68 + 82 + 50 = 200 \gg</math> 200 gallons. She was able to sell 200 gallons - 24 gallons = <math>\ll 200 - 24 = 176 \gg</math> 176 gallons. Thus, her total revenue for the milk is <math>\\$3.50/\text{gallon} \times 176 \text{ gallons} = \ll 3.50 \times 176 = 616 \gg</math> 616. <b>Final Answer:</b> 616</p>
<p><b>Problem:</b> Tina buys 3 12-packs of soda for a party. Including Tina, 6 people are at the party. Half of the people at the party have 3 sodas each, 2 of the people have 4, and 1 person has 5. How many sodas are left over when the party is over?</p> <p><b>Solution:</b> Tina buys 3 12-packs of soda, for <math>3 \times 12 = \ll 3 \times 12 = 36 \gg</math> 36 sodas 6 people attend the party, so half of them is <math>6/2 = \ll 6/2 = 3 \gg</math> 3 people Each of those people drinks 3 sodas, so they drink <math>3 \times 3 = \ll 3 \times 3 = 9 \gg</math> 9 sodas Two people drink 4 sodas, which means they drink <math>2 \times 4 = \ll 2 \times 4 = 8 \gg</math> 8 sodas With one person drinking 5, that brings the total drank to <math>5 + 9 + 8 + 3 = \ll 5 + 9 + 8 + 3 = 25 \gg</math> 25 sodas As Tina started off with 36 sodas, that means there are <math>36 - 25 = \ll 36 - 25 = 11 \gg</math> 11 sodas left <b>Final Answer:</b> 11</p>

Figure 1: Three example problems from GSM8K. Calculation annotations are highlighted in red.

**GSM8k** [Cobbe et al. 2021]

- Issue: Limited size.
- Can we get more data?

## MATH Dataset (Ours)

**Problem:** Tom has a red marble, a green marble, a blue marble, and three identical yellow marbles. How many different groups of two marbles can Tom choose?

**Solution:** There are two cases here: either Tom chooses two yellow marbles (1 result), or he chooses two marbles of different colors ( $\binom{4}{2} = 6$  results). The total number of distinct pairs of marbles Tom can choose is  $1 + 6 = \boxed{7}$ .

**Problem:** If  $\sum_{n=0}^{\infty} \cos^{2n} \theta = 5$ , what is  $\cos 2\theta$ ?

**Solution:** This geometric series is  $1 + \cos^2 \theta + \cos^4 \theta + \dots = \frac{1}{1 - \cos^2 \theta} = 5$ . Hence,

$$\cos^2 \theta = \frac{4}{5}. \text{ Then } \cos 2\theta = 2 \cos^2 \theta - 1 = \boxed{\frac{3}{5}}.$$

**Problem:** The equation  $x^2 + 2x = i$  has two complex solutions. Determine the product of their real parts.

**Solution:** Complete the square by adding 1 to each side.

Then  $(x + 1)^2 = 1 + i = e^{\frac{i\pi}{4}} \sqrt{2}$ , so  $x + 1 = \pm e^{\frac{i\pi}{8}} \sqrt[4]{2}$ .

The desired product is then

$$\begin{aligned} & (-1 + \cos(\frac{\pi}{8}) \sqrt[4]{2}) (-1 - \cos(\frac{\pi}{8}) \sqrt[4]{2}) = \\ & 1 - \cos^2(\frac{\pi}{8}) \sqrt{2} = 1 - \frac{(1 + \cos(\frac{\pi}{4}))}{2} \sqrt{2} = \boxed{\frac{1 - \sqrt{2}}{2}}. \end{aligned}$$

**MATH** [Hendrycks et al. 2021]

# Generating more data

## Using the model itself

- Feature of (some) math datasets: having a final answer —> **easy verification**
  - Not true for theorems —> using formal math (e.g., Lean)
- How to generate data using the model?
  - Give questions to the model and generate solutions.
  - Check the final answers and keep the correct ones. **(slightly error prone)**
  - Fine tune the model on the correct solutions.

# Generating more data

## STaR

---

### Algorithm 1 STaR

---

**Input**  $M$ : a pretrained LLM; dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^D$  (w/ few-shot prompts)

- 1:  $M_0 \leftarrow M$  # Copy the original model
- 2: **for**  $n$  **in**  $1 \dots N$  **do** # Outer loop
- 3:    $(\hat{r}_i, \hat{y}_i) \leftarrow M_{n-1}(x_i) \quad \forall i \in [1, D]$  # Perform rationale generation
- 4:    $(\hat{r}_i^{\text{rat}}, \hat{y}_i^{\text{rat}}) \leftarrow M_{n-1}(\text{add\_hint}(x_i, y_i)) \quad \forall i \in [1, D]$  # Perform rationalization
- 5:    $\mathcal{D}_n \leftarrow \{(x_i, \hat{r}_i, y_i) \mid i \in [1, D] \wedge \hat{y}_i = y_i\}$  # Filter rationales using ground truth answers
- 6:    $\mathcal{D}_n^{\text{rat}} \leftarrow \{(x_i, \hat{r}_i^{\text{rat}}, y_i) \mid i \in [1, D] \wedge \hat{y}_i \neq y_i \wedge \hat{y}_i^{\text{rat}} = y_i\}$  # Filter rationalized rationales
- 7:    $M_n \leftarrow \text{train}(M, \mathcal{D}_n \cup \mathcal{D}_n^{\text{rat}})$  # Finetune the original model on correct solutions - inner loop
- 8: **end for**

---

Wasting incorrect solutions :(

	GSM8K Test Accuracy (%)	Train Data Used (%)
Few-shot Direct GPT-J	3.0	$\sim 0$
Few-shot CoT GPT-J	3.1	$\sim 0$
GPT-J Direct Finetuned	5.8	100
STaR without rationalization	10.1	25.0
STaR with rationalization	<b>10.7</b>	28.7

---

# Reinforcement Learning

# Sources



## Alignment and RL for LLMs

Caglar Gulcehre and Volkan Cevher

EPFL

1

## Lecture 1: Introduction to RL

Professor Emma Brunskill

CS234 RL

Winter 2025

- Today the 3rd part of the lecture includes some slides from David Silver's introduction to RL slides or modifications of those slides

- With many slides from or derived from David Silver and John Schulman and Pieter Abbeel

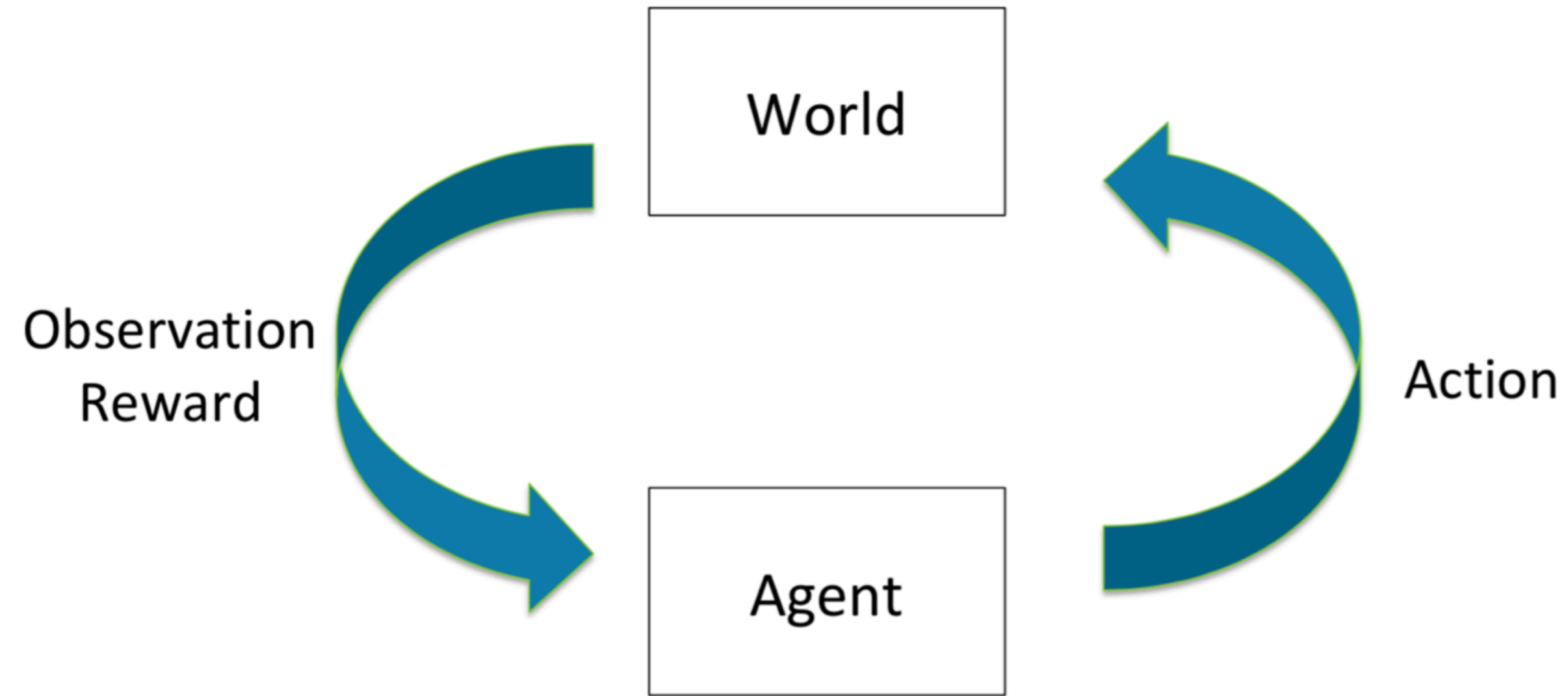
# Reinforcement learning (RL)

- **Reinforcement learning:** Learning through experience/data to make good decisions under uncertainty.
- Builds strongly from theory and ideas starting in the 1950s with Richard Bellman.
- Commonly used for robotics.
- Impressive success cases:
  - math, physics (nuclear fusion), games, finance, ...



The image is a screenshot of a BBC News article. At the top, the BBC logo is visible on the left, and navigation links for 'Home', 'News', 'Sport', 'Business', and 'Innovation' are on the right. Below the navigation is a red banner with the word 'NEWS' in white. Underneath the banner, there are more navigation links: 'Home', 'Israel-Gaza war', 'War in Ukraine', 'Climate', 'Video', 'World', 'UK', 'Business', 'Tech', and 'Science'. The article title is 'Artificial intelligence: Google's AlphaGo beats Go master Lee Se-dol' in a large, bold, black font. Below the title, the date '12 March 2016' is shown. There is a red share icon to the left of the main image. The main image is a Go board with black and white stones. In the bottom left corner of the image, there is a logo for 'AlphaGo Google DeepMind' and a timer showing '00:08:32'. In the bottom right corner, there is a small inset image of Lee Se-dol and a timer showing '00:00:27'. The word 'OTHER' is visible in the bottom left corner of the image area.

A computer program has beaten a master Go player 3-0 in a best-of-five competition, in what is seen as a landmark moment for artificial intelligence.



Goal: Select actions to maximize total expected future reward

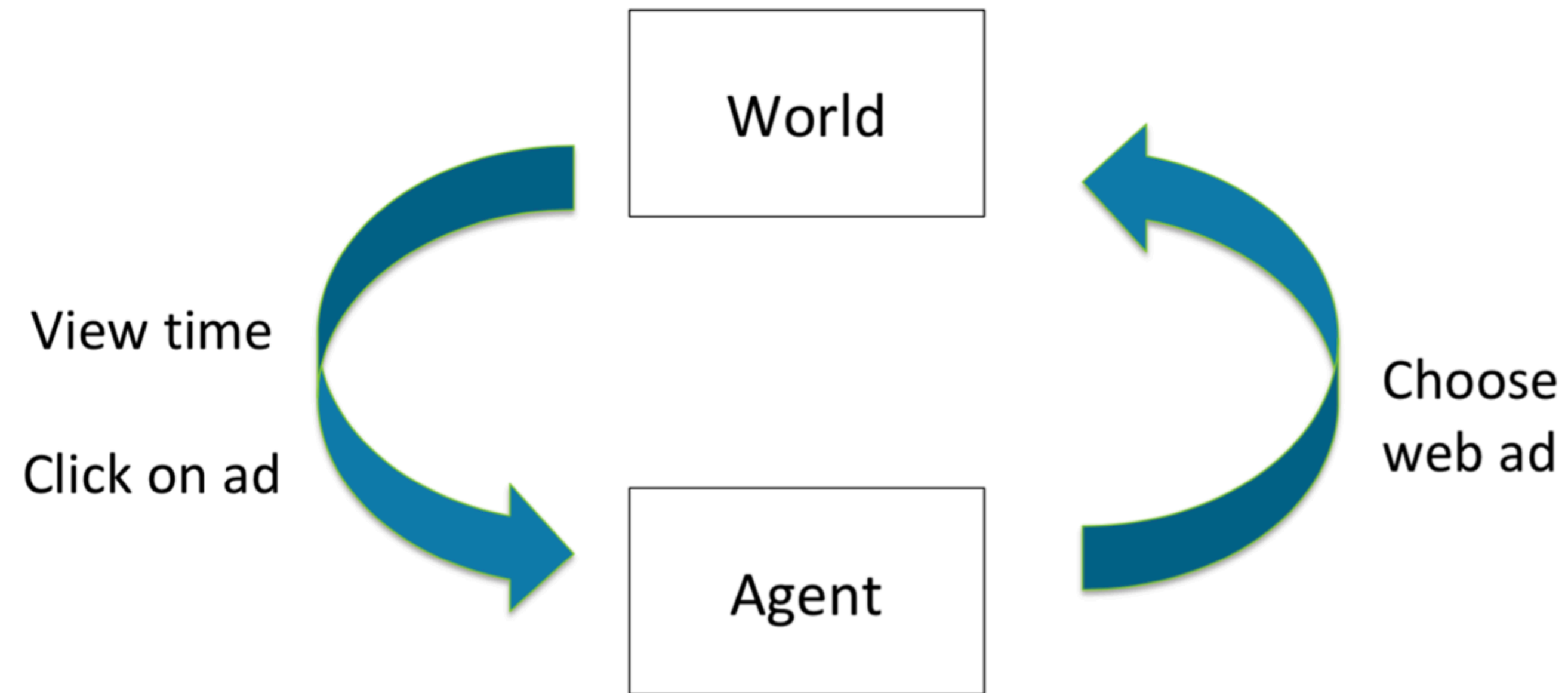
May require balancing immediate & long term rewards

# RL properties

- Delayed consequences
  - Decisions now can impact things much later.
  - E.g., Saving for retirement, actions in a game, or actions of a robot.
- Challenges:
  - When planning: decisions involve reasoning about not just immediate benefit of a decision but also its longer term consequences.
  - When learning: temporal credit assignment is hard (what caused later high or low rewards?)

# RL properties

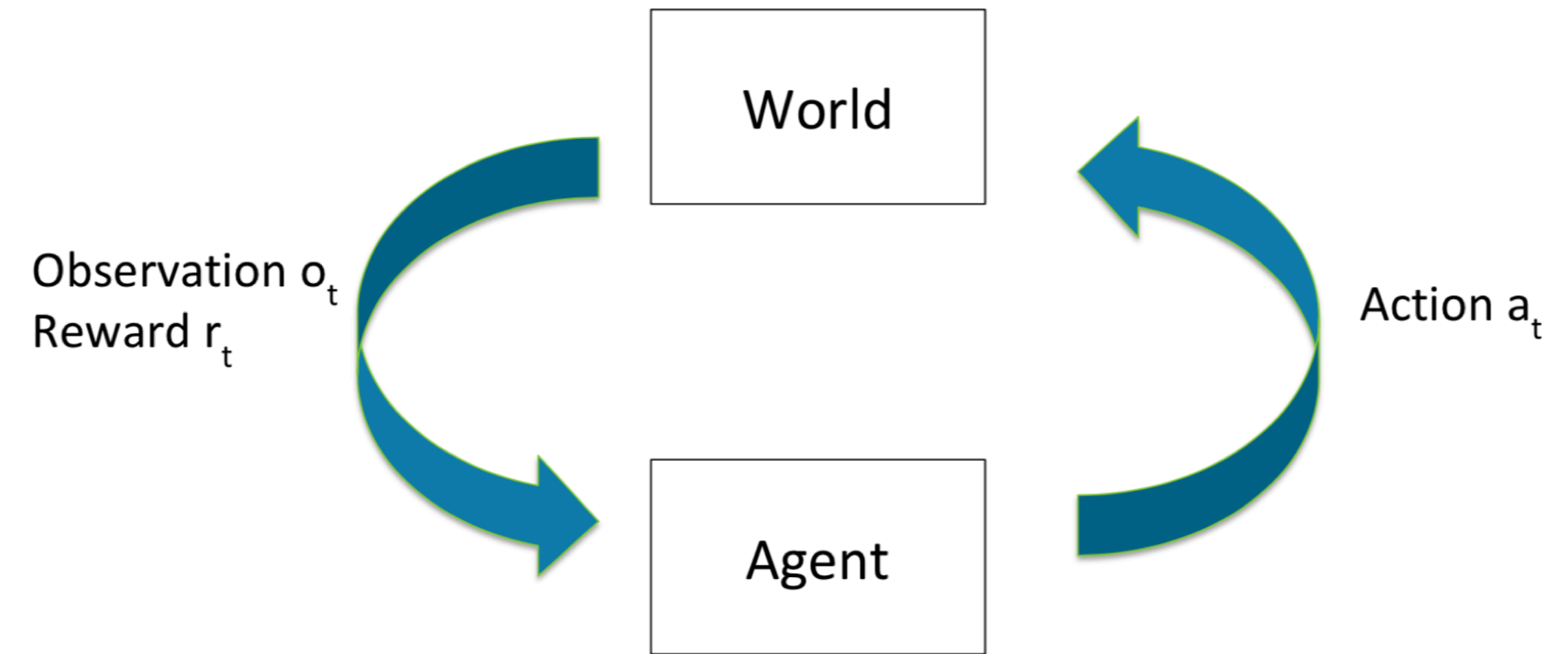
## Web advertising



# RL properties

- RL: Learning about the world by making decisions
  - Learn to ride a bike by trying (and failing)
- We want a mapping from current state to the action
  - However, the space of states is usually very large
    - E.g., in (video) games or robotics
    - A simple video game frame, may contain 200x300 resolution, and each pixel can have  $256^3$  values.

# RL properties



- History  $h_t = (a_1, o_1, r_1, \dots, a_t, o_t, r_t)$
- State is information assumed to determine what happens next Function of history:  $s_t = (h_t)$
- **Markov assumption**
  - Information state: sufficient statistic of the history
  - $\mathbb{P}(s_{t+1} | s_t, a_t) = \mathbb{P}(s_{t+1} | h_t, a_t)$  (next state might not be deterministic)
  - In practice we may assume that  $o_t = s_t$

# Markov Decision Process (MDP)

- $S$  is a set of Markov **states**.
- We have a set of **actions**  $A$ .
- We have a **transition function** specifying  $\mathbb{P}(s_{t+1} = s' \mid s_t = s, a_t = a)$ .
- We have **reward** function  $R(s_t = s, a_t = a) = \mathbb{E}[r_t \mid s_t = s, a_t = a]$ .
- **Discount factor**  $\gamma \in [0,1]$ .
  - Total reward is computed by  $r_1 + \gamma r_2 + \gamma^2 r_3 + \dots$
  - Preferring earlier rewards (like humans)
  - If horizon is finite, we can use  $\gamma = 1$ .

# Markov Decision Process (MDP)

- Policy:  $\pi(a | s) = \mathbb{P}(a_t = a | s_t = s)$

$$R^\pi(s) = \sum_{a \in A} \pi(a|s) R(s, a)$$

$$P^\pi(s'|s) = \sum_{a \in A} \pi(a|s) P(s'|s, a)$$

$$V^\pi(s) = \underbrace{R^\pi(s)}_{\text{Immediate reward}} + \underbrace{\gamma \sum_{s' \in S} P^\pi(s'|s) V^\pi(s')}_{\text{Discounted sum of future rewards}}$$

$$V^\pi = R^\pi + \gamma P^\pi V^\pi$$

# Markov Decision Process (MDP)

**Policy:**  $\pi(a | s) = \mathbb{P}(a_t = a | s_t = s)$

**State-action value:**  $Q^\pi(s, a) = R(s, a) + \gamma \sum_{s'} P(s' | s, a) V^\pi(s')$

How to find the optimal policy? (Assuming finite states)

- **Policy iteration**
  - Compute values for a policy and use it to get a better policy.
- **Value iteration**
  - Compute optimal value for horizon  $k$  and then  $k + 1$ .

# RL for LMs

Agent (policy) = language model

Initial state

**Problem:** Suppose  $a$  and  $b$  are positive real numbers with  $a > b$  and  $ab = 8$ . Find the minimum value of  $\frac{a^2+b^2}{a-b}$ .

**Ground truth solution:** We can write  $\frac{a^2+b^2}{a-b} = \frac{a^2+b^2-2ab+16}{a-b} = \frac{(a-b)^2+16}{a-b} = a - b + \frac{16}{a-b}$ . By AM-GM,  $a - b + \frac{16}{a-b} \geq 2\sqrt{(a-b) \cdot \frac{16}{a-b}} = 8$ . Equality occurs when  $a - b = 4$  and  $ab = 8$ . We can solve these equations to find  $a = 2\sqrt{3} + 2$  and  $b = 2\sqrt{3} - 2$ . Thus, the minimum value is  $\boxed{8}$ .

Steps = actions

reward = 1 if  
answer is correct

**Intuition:** Sparse-reward MDP with deterministic dynamics

# LMs and RL

Policy is parametrized:  $\pi_{\theta}(s, a) = \mathbb{P}(a | s; \theta)$

Goal: maximize  $V(\theta) := V(s_0, \theta)$

We need to compute the **policy gradient**  $\nabla_{\theta} V(s_0, \theta)$

- assuming  $\pi_{\theta}$  is differentiable

$$V(\theta) = \sum_{\tau} \mathbb{P}(\tau; \theta) R(\tau) \text{ where } \tau \text{ presents a trajectory } (s_0, a_0, r_0, \dots)$$

# Policy Gradient

$$\begin{aligned}\nabla_{\theta} V(\theta) &= \nabla_{\theta} \sum_{\tau} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} \frac{P(\tau; \theta)}{P(\tau; \theta)} \nabla_{\theta} P(\tau; \theta) R(\tau) \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \underbrace{\frac{\nabla_{\theta} P(\tau; \theta)}{P(\tau; \theta)}}_{\text{likelihood ratio}} \\ &= \sum_{\tau} P(\tau; \theta) R(\tau) \nabla_{\theta} \log P(\tau; \theta)\end{aligned}$$

Approximation using  $m$  samples under  $\pi_{\theta}$ :

$$\nabla_{\theta} V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_{\theta} \log P(\tau^{(i)}; \theta)$$

# Policy Gradient

Approximation using  $m$  samples under  $\pi_\theta$ :

$$\nabla_\theta V(\theta) \approx \hat{g} = (1/m) \sum_{i=1}^m R(\tau^{(i)}) \nabla_\theta \log P(\tau^{(i)}; \theta)$$

$$\begin{aligned} \nabla_\theta \log P(\tau^{(i)}; \theta) &= \nabla_\theta \log \left[ \underbrace{\mu(s_0)}_{\text{Initial state distrib.}} \prod_{t=0}^{T-1} \underbrace{\pi_\theta(a_t|s_t)}_{\text{policy}} \underbrace{P(s_{t+1}|s_t, a_t)}_{\text{dynamics model}} \right] \\ &= \nabla_\theta \left[ \log \mu(s_0) + \sum_{t=0}^{T-1} \log \pi_\theta(a_t|s_t) + \log P(s_{t+1}|s_t, a_t) \right] \\ &= \sum_{t=0}^{T-1} \underbrace{\nabla_\theta \log \pi_\theta(a_t|s_t)}_{\text{no dynamics model required!}} \quad \text{score function} \end{aligned}$$

## REINFORCE:

Initialize policy parameters  $\theta$  arbitrarily

**for** each episode  $\{s_1, a_1, r_2, \dots, s_{T-1}, a_{T-1}, r_T\} \sim \pi_\theta$  **do**

**for**  $t = 1$  to  $T - 1$  **do**

$$\theta \leftarrow \theta + \alpha \nabla_\theta \log \pi_\theta(s_t, a_t) G_t$$

**endfor**

**endfor**

**return**  $\theta$

where  $G_t = \sum_{k=t}^{T-1} \gamma^{k-t} R_{k+1}$  is the [return](#).

# Generating more data

## STaR

---

### Algorithm 1 STaR

---

**Input**  $M$ : a pretrained LLM; dataset  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^D$  (w/ few-shot prompts)

- 1:  $M_0 \leftarrow M$  # Copy the original model
- 2: **for**  $n$  **in**  $1 \dots N$  **do** # Outer loop
- 3:    $(\hat{r}_i, \hat{y}_i) \leftarrow M_{n-1}(x_i) \quad \forall i \in [1, D]$  # Perform rationale generation
- 4:    $(\hat{r}_i^{\text{rat}}, \hat{y}_i^{\text{rat}}) \leftarrow M_{n-1}(\text{add\_hint}(x_i, y_i)) \quad \forall i \in [1, D]$  # Perform rationalization
- 5:    $\mathcal{D}_n \leftarrow \{(x_i, \hat{r}_i, y_i) \mid i \in [1, D] \wedge \hat{y}_i = y_i\}$  # Filter rationales using ground truth answers
- 6:    $\mathcal{D}_n^{\text{rat}} \leftarrow \{(x_i, \hat{r}_i^{\text{rat}}, y_i) \mid i \in [1, D] \wedge \hat{y}_i \neq y_i \wedge \hat{y}_i^{\text{rat}} = y_i\}$  # Filter rationalized rationales
- 7:    $M_n \leftarrow \text{train}(M, \mathcal{D}_n \cup \mathcal{D}_n^{\text{rat}})$  # Finetune the original model on correct solutions - inner loop
- 8: **end for**

---

*STaR is very similar to REINFORCE*

$$J(M, X, Y) = \sum_i \mathbb{E}_{\hat{r}_i, \hat{y}_i \sim p_M(\cdot | x_i)} \mathbb{1}(\hat{y}_i = y_i), \quad (1)$$

$$\nabla J(M, X, Y) = \sum_i \mathbb{E}_{\hat{r}_i, \hat{y}_i \sim p_M(\cdot | x_i)} [\mathbb{1}(\hat{y}_i = y_i) \cdot \nabla \log p_M(\hat{y}_i, \hat{r}_i | x_i)], \quad (2)$$

# Policy Gradient with baseline

We can deduct a baseline value from the return:

$$\nabla_{\theta} \mathbb{E}_{\tau} [R] = \mathbb{E}_{\tau} \left[ \sum_{t=0}^{T-1} \nabla_{\theta} \log \pi(a_t | s_t; \theta) \left( \sum_{t'=t}^{T-1} r_{t'} - b(s_t) \right) \right] \quad (\text{wlog } \gamma = 1)$$

+ The estimator remains **unbiased**.

+ We can **reduce the variance**.

- In particular the expected return or the state-value  $V^{\pi}(s)$  is a great baseline.
- We'd like to assess how much the return is better/worse than the average.
- We call  $G_t - b(s_t) = G_t - V^{\pi}(s_t) = A^{\pi}(s_t, a_t)$  the **advantage**.

# Challenges of policy gradient

- Distance in parameter space vs. distance in policy space
  - Choosing the right learning rate can be hard. —> **stability** issues.
- Sample efficiency is poor.
  - PG is **on-policy** => we need to get rid of the data after each gradient step.
  - Ideally we'd like to **use old data for multiple gradient steps**.
- Solution: We can use samples from the old policy, as long as:
  - Policies are close, in particular, **low KL divergence**.
  - **Importance sampling weight** is included.

# Challenges of policy gradient

- Assuming  $D_{KL}(\pi' || \pi)$  is small, we can equivalently find  $\pi'$  maximizing:

$$\mathbb{E}_{\tau \sim \pi} \left[ \sum_{t=0}^{\infty} \gamma^t \frac{\pi'(a_t | s_t)}{\pi(a_t | s_t)} A^{\pi}(s_t, a_t) \right]$$

where  $\pi$  is the old policy.



*stability challenge*

# Proximal Policy Optimization (PPO)

• Using **KL penalty**:  $\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}(\theta) - \beta_k \bar{D}_{KL}(\theta || \theta_k)$

• Using **clipped objective**:  $\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$

where  $r_t(\theta) = \pi_{\theta}(a_t | s_t) / \pi_{\theta_k}(a_t | s_t)$ .

# Proximal Policy Optimization (PPO)

---

## Algorithm PPO with Clipped Objective

---

Input: initial policy parameters  $\theta_0$ , clipping threshold  $\epsilon$

**for**  $k = 0, 1, 2, \dots$  **do**

Collect set of partial trajectories  $\mathcal{D}_k$  on policy  $\pi_k = \pi(\theta_k)$

Estimate advantages  $\hat{A}_t^{\pi_k}$  using any advantage estimation algorithm

Compute policy update

$$\theta_{k+1} = \arg \max_{\theta} \mathcal{L}_{\theta_k}^{CLIP}(\theta)$$

by taking  $K$  steps of minibatch SGD (via Adam), where

$$\mathcal{L}_{\theta_k}^{CLIP}(\theta) = \mathbb{E}_{\tau \sim \pi_k} \left[ \sum_{t=0}^T \left[ \min(r_t(\theta) \hat{A}_t^{\pi_k}, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t^{\pi_k}) \right] \right]$$

**end for**

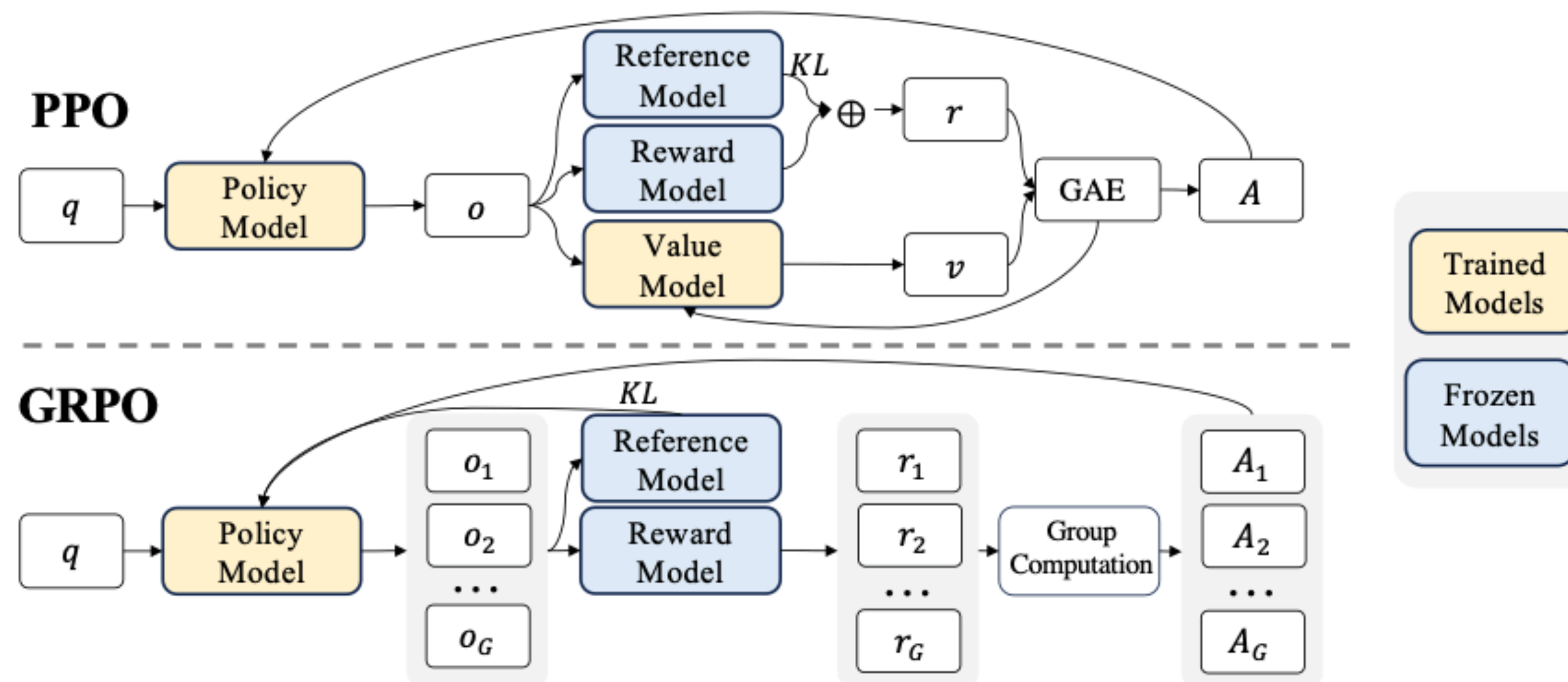
---

Advantage is often estimated by learning a separate value model, e.g.,

$$A_t = G_t - V_{\phi}(s_t)$$

# Group Relative Policy Optimization (GRPO)

- GRPO generates multiple trajectories per question, and uses the averaged (normalized) reward as advantage. —> **no need for a value model!**



DeepSeekMath: Pushing the Limits of Mathematical Reasoning in Open Language Models

Zhihong Shao<sup>1,2\*†</sup>, Peiyi Wang<sup>1,3\*†</sup>, Qihao Zhu<sup>1,3\*†</sup>, Runxin Xu<sup>1</sup>, Junxiao Song<sup>1</sup>,  
Xiao Bi<sup>1</sup>, Haowei Zhang<sup>1</sup>, Mingchuan Zhang<sup>1</sup>, Y.K. Li<sup>1</sup>, Y. Wu<sup>1</sup>, Daya Guo<sup>1\*</sup>

$$A_i = \frac{r_i - \text{mean}(r)}{\text{std}(r)}$$

*GRPO is now commonly used for training reasoning models.*

# Sparse rewards problem

- Hard questions may not provide any reward.
- Solutions:
  - Using a **curriculum**: training on simpler questions first.
  - Helping the model by **making the question easier**.
  - Process Reward Model (PRM): **rewarding individual steps** of the solution.