



# CS 639: Foundation Models **Alignment I**

Fred Sala

University of Wisconsin-Madison

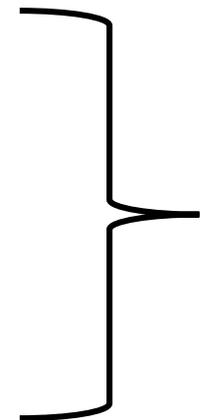
**March 17, 2026**



# Announcements

- **Exam:** grading will be done today
  - Hoping to release tomorrow
- **Homework 3:** out.
- **Project information:** here  
[https://pages.cs.wisc.edu/~fredsala/cs639/files/project\\_info\\_639.pdf](https://pages.cs.wisc.edu/~fredsala/cs639/files/project_info_639.pdf)
- **Class outline:**

Tuesday March 17	Alignment I
Thursday March 19	Alignment II
Tuesday March 24	Reasoning I: More CoT
Thursday March 26	Reasoning II: RLVR



# Outline

- **Last Time: Cross-Modal Adaptation, Tools**

- ORCA, aligning via optimal transport dataset distance, simple tool integration, Toolformer

- **Alignment and RLHF**

- Basic alignment idea, goals, mechanisms, RL review, RLHF steps

- **Why Does RLHF Work?**

- Failures of supervised learning, knowledge-seeking interactions, abstains

# Outline

- **Last Time: Cross-Modal Adaptation, Tools**

- ORCA, aligning via optimal transport dataset distance, simple tool integration, Toolformer

- **Alignment and RLHF**

- Basic alignment idea, goals, mechanisms, RL review, RLHF steps

- **Why Does RLHF Work?**

- Failures of supervised learning, knowledge-seeking interactions, abstains

# What About Other Modalities?

So far, mostly talked about fine-tuning language models.

- Suppose we want tasks that are not directly language-based
- Could just train a new model...
  - As we did for multimodal models – but might be expensive.

Can we adapt language models? Lots of **challenges**:

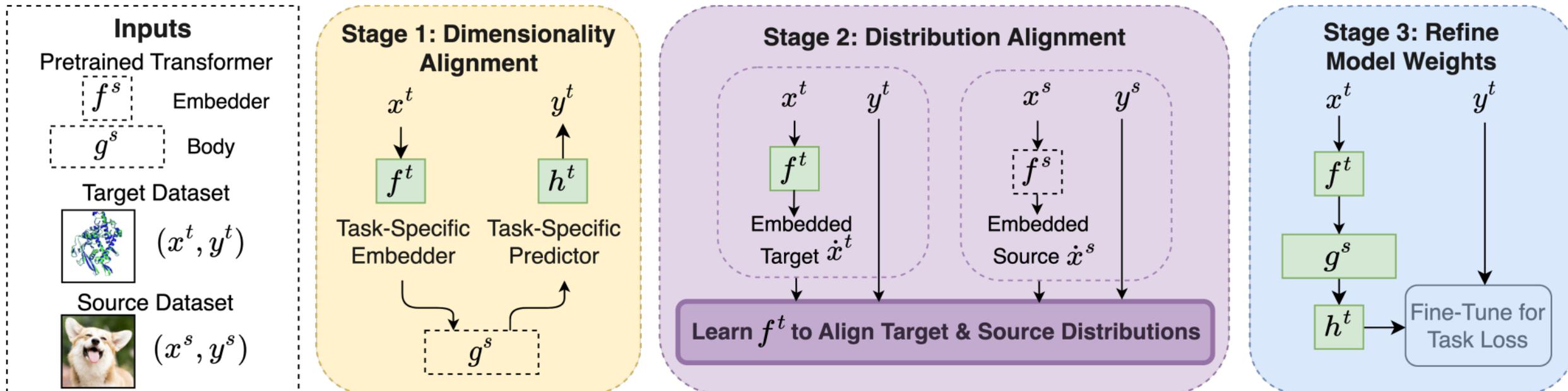
- Must change data types
- How do we know modalities are usable together?

# Cross-Modal: ORCA

Performance bottleneck in FPTs

A more powerful approach: ORCA (Shen et al '23)

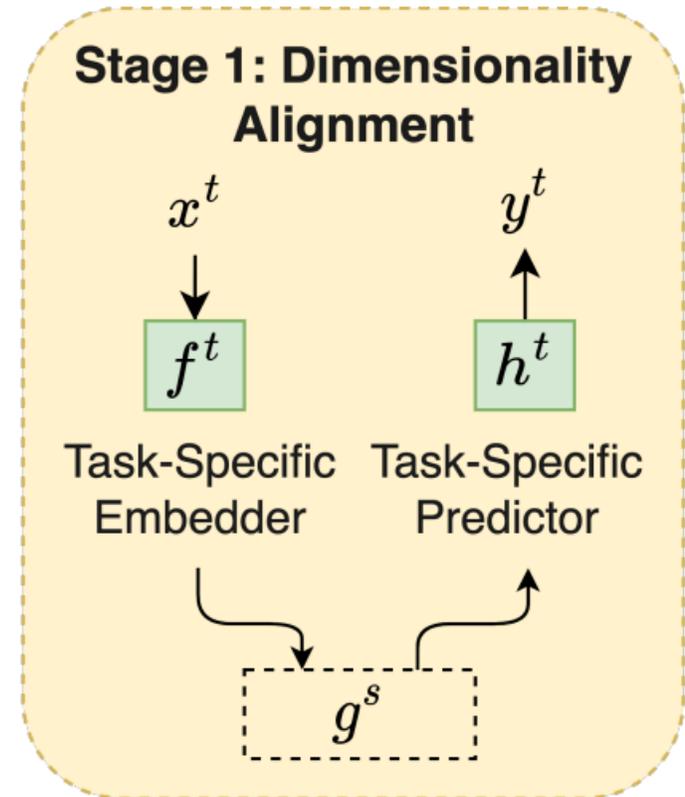
- Adds: distribution alignment step (align then refine)



# ORCA: Stage 1

Let's understand each stage of ORCA

- Stage 1: compatibility for inputs and outputs
- Custom input and output embedders that depend on the task
  - Input example: convolutional layers for image settings
  - Output example: average pooling+linear layer for classification

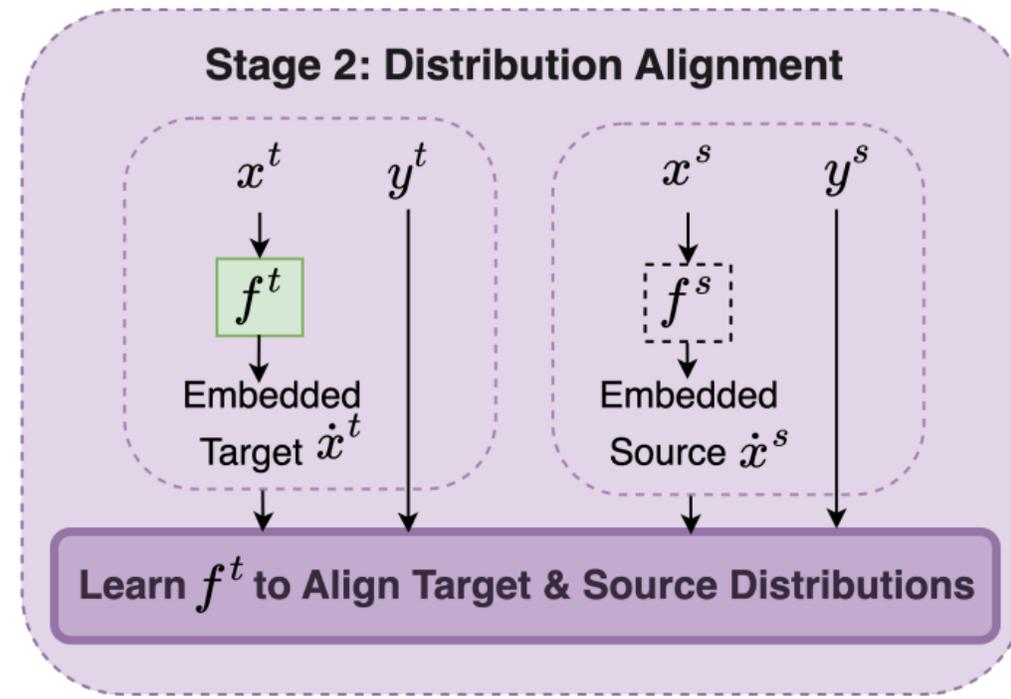


# ORCA: Stage 2

Let's understand each stage of ORCA

- Stage 2: distribution alignment
- Intuition:
  - Change embeddings so target features **resemble** source features
- Learn the function  $f^t$  that **minimizes distance between**  
 $(f^t(x^t), y^t)$  and  $(f^s(x^s), y^s)$

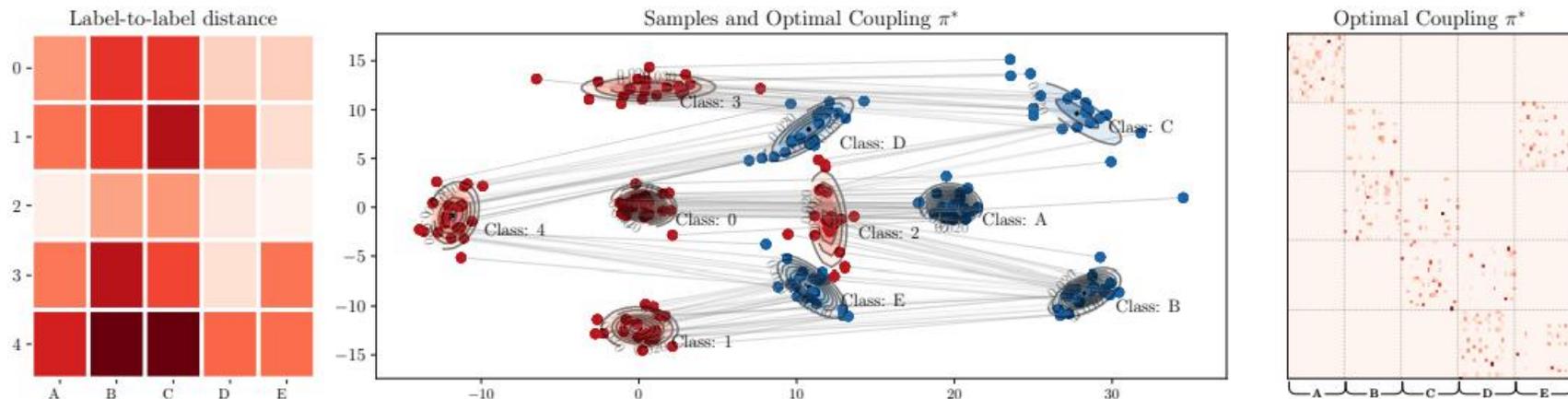
**Warning: this part is math-heavy, but we just need the intuition**



# ORCA: Distributional Distances

Want: learn the function  $f^t$  that minimizes distance between  $(f^t(x^t), y^t)$  and  $(f^s(x^s), y^s)$

- How?
- Need a distance function on these distributions
- Here, **optimal transport dataset distance (OTDD)**



# Interlude: Optimal Transport

In optimal transport, we solve

$$\inf \left\{ \int_{X \times Y} c(x, y) d\gamma(x, y) \mid \gamma \in \Gamma(\mu, \nu) \right\},$$

**Cost or distance**  
of moving  $x$  to  $y$

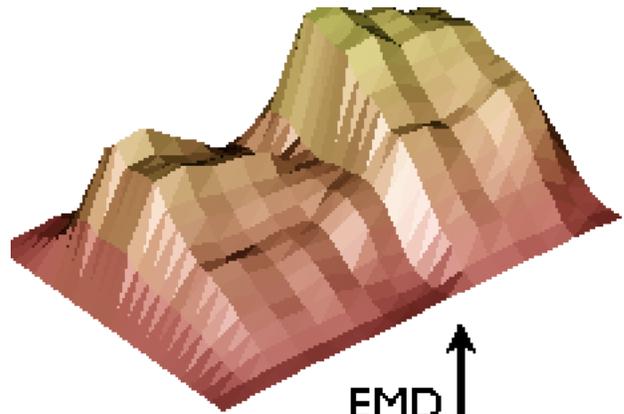
The two **marginals** we care  
about, i.e., on  $x$  and  $y$

- Want to “move” distribution on  $x$  to one on  $y$ 
  - Output is a joint distribution with the original source and target
- But there’s a cost to moving  $x$  to  $y$ , given by  $c(x, y)$

# Interlude: Optimal Transport

In optimal transport, we solve

$$\inf \left\{ \int_{X \times Y} c(x, y) d\gamma(x, y) \mid \gamma \in \Gamma(\mu, \nu) \right\},$$



EMD  
↕



**Cost or distance**  
of moving  $x$  to  $y$

The two **marginals** we care  
about, i.e., on  $x$  and  $y$

Applegate et al '11

# Interlude: **Optimal Transport**

In optimal transport, we solve

$$\inf \left\{ \int_{X \times Y} c(x, y) d\gamma(x, y) \mid \gamma \in \Gamma(\mu, \nu) \right\},$$

- Cost given by **distance**: Wasserstein distance
- Gives a distance on distributions, i.e.,

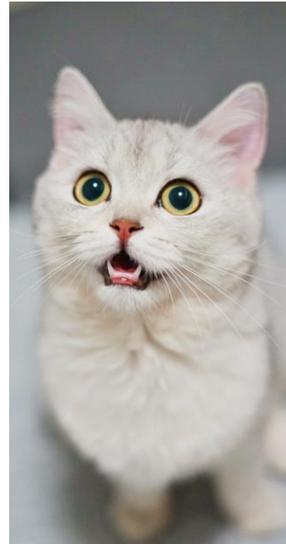
$$W_p(\mu, \nu) = \left( \inf_{\gamma \in \Gamma(\mu, \nu)} \mathbf{E}_{(x, y) \sim \gamma} d(x, y)^p \right)^{1/p}$$

# Interlude: Dataset Distance

What should this cost/distance  $c(x,y)$  be for us?

- For inputs  $x$ , pretty easy: feature vectors in spaces that have distances, e.g.,  $\|x-x'\|$

- For outputs  $y$ , not so easy



—



- A clever idea:

- Replace  $y$  with  $P(X|y)$

- Even harder? No, just use Wasserstein:  $W(P(X|y),P(X|y'))$

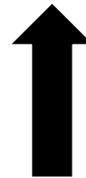
- Approximate this with a Gaussian: closed form too!

# ORCA: Distributional Distances

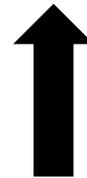
Want: learn the function  $f^t$  that minimizes distance between  $(f^t(x^t), y^t)$  and  $(f^s(x^s), y^s)$

- Need a distance function on these distributions
- Here, **optimal transport dataset distance (OTDD)**

$$d_{\mathcal{Z}}((x, y), (x', y')) \triangleq \left( d_{\mathcal{X}}(x, x')^p + \mathbf{W}_p^p(\alpha_y, \alpha_{y'}) \right)^{1/p}$$



i.e., Euclidean  
distance

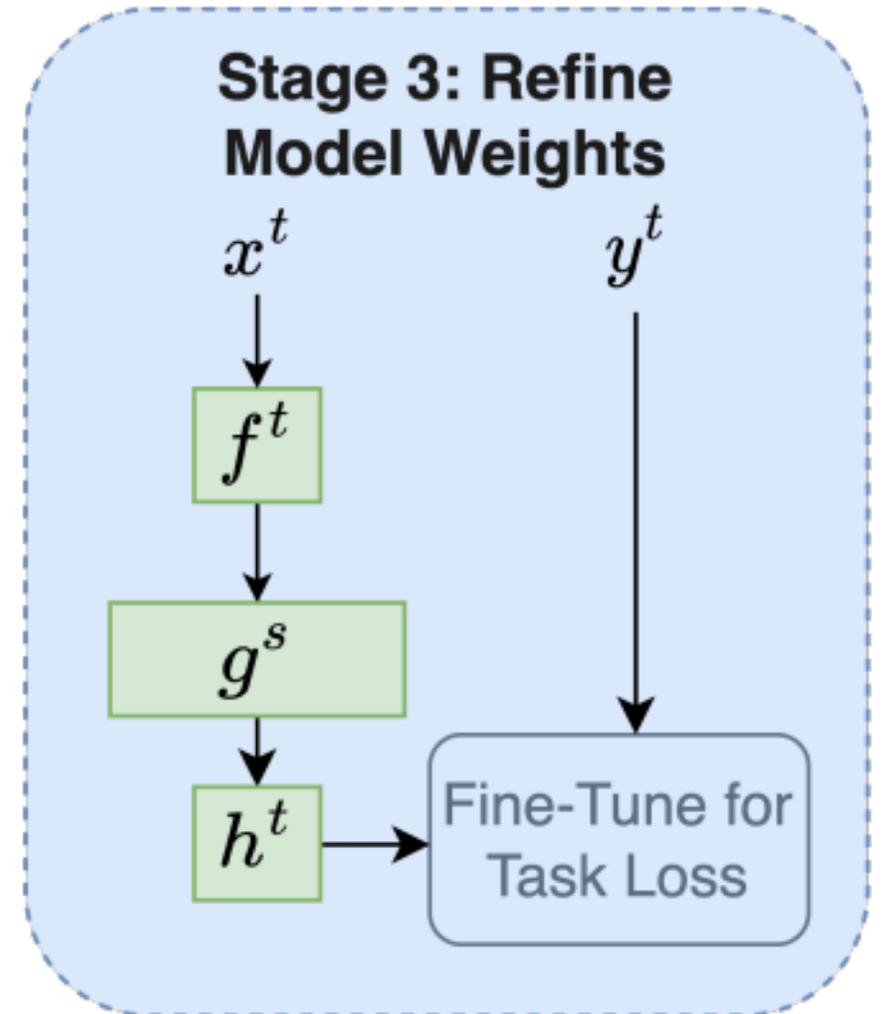


p-Wasserstein distance on  
 $P(x|y)$

# ORCA: Stage 3

Let's understand each stage of ORCA

- Stage 3: fine-tune the input and output network weights
  - For particular tasks
  - Or, could do any other variant of what we've talked about...



# ORCA: Results

Extremely good, even against state-of-the-art results

- Compare to Neural Architecture Search (NAS)
  - Produces custom architectures that hit state-of-the-art for various tasks
  - Same procedure on many types of tasks works well:

	CIFAR-100 0-1 error (%)	Spherical 0-1 error (%)	Darcy Flow relative $\ell_2$	PSICOV MAE <sub>8</sub>	Cosmic 1-AUROC	NinaPro 0-1 error (%)	FSD50K 1- mAP	ECG 1 - F1 score	Satellite 0-1 error (%)	DeepSEA 1- AUROC
Hand-designed	19.39	67.41	8E-3	3.35	<b>0.127</b>	8.73	0.62	<b>0.28</b>	19.80	0.30
NAS-Bench-360	23.39	48.23	2.6E-2	2.94	0.229	7.34	0.60	0.34	12.51	0.32
DASH	24.37	71.28	7.9E-3	3.30	0.19	<b>6.60</b>	0.60	0.32	12.28	<b>0.28</b>
Perceiver IO	70.04	82.57	2.4E-2	8.06	0.485	22.22	0.72	0.66	15.93	0.38
FPT	10.11	76.38	2.1E-2	4.66	0.233	15.69	0.67	0.50	20.83	0.37
<b>ORCA</b>	<b>6.53</b>	<b>29.85</b>	<b>7.28E-3</b>	<b>1.91</b>	0.152	7.54	<b>0.56</b>	<b>0.28</b>	<b>11.59</b>	0.29

# Beyond Unaided Language Models

Even when we fine-tune or optimize prompts (or do CoT), the language model can get things wrong.

- Often simple things... like **arithmetic**.

## Data Formatting

### Plain

128+367=495

### Reverse

\$128+367=594\$

### Simplified Scratchpad

Input: 128+367

Target:

A->5, C->1

A->9, C->0

A->4, C->0.

495

### Detailed Scratchpad

Input:

128+367

Target:

<scratch>

[1,2,8] has 3 digits.

[3,6,7] has 3 digits.

[1,2,8] + [3,6,7], C=0, 8+7+0=15, A->5, C->1

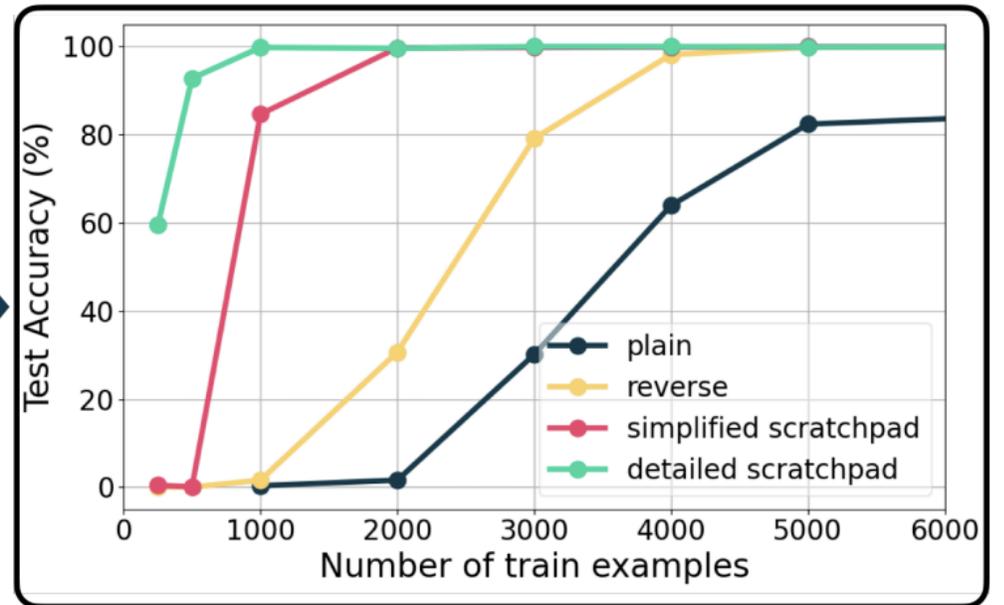
[1,2] + [3, 6], A= [5], 2+6+1=9, A->9, C->0

[1] + [3], A= [9,5], C=0, 1+3+0=4, A->4, C->0

[] + [], A= [4,9,5], C=0, END

</scratch>

4 9 5



# Beyond Unaided Language Models

Even when we fine-tune or optimize prompts (or do CoT), the language model can get things wrong.

- Often simple things... like **arithmetic**.
- How else can we help it?
- **A:** Use external tools



# Tools: Program-aided LMs

Use external tools:

- Python interpreter
- How? *Interleave* the text explanations in CoT steps with lines of Python code
- LMs can already output code
  - Just need to *prompt* the right way

## Chain-of-Thought (Wei et al., 2022)

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 tennis balls. 2 cans of 3 tennis balls each is 6 tennis balls.  $5 + 6 = 11$ . The answer is 11.

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves. They sold 93 in the morning and 39 in the afternoon. So they sold  $93 + 39 = 132$  loaves. The grocery store returned 6 loaves. So they had  $200 - 132 - 6 = 62$  loaves left. The answer is 62.



Gao et al '23

## Program-aided Language models (this work)

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 tennis balls.  
`tennis_balls = 5`  
`2 cans of 3 tennis balls each is`  
`bought_balls = 2 * 3`  
`tennis_balls`. The answer is  
`answer = tennis_balls + bought_balls`

Q: The bakers at the Beverly Hills Bakery baked 200 loaves of bread on Monday morning. They sold 93 loaves in the morning and 39 loaves in the afternoon. A grocery store returned 6 unsold loaves. How many loaves of bread did they have left?

Model Output

A: The bakers started with 200 loaves  
`loaves_baked = 200`  
`They sold 93 in the morning and 39 in the afternoon`  
`loaves_sold_morning = 93`  
`loaves_sold_afternoon = 39`  
`The grocery store returned 6 loaves.`  
`loaves_returned = 6`  
The answer is  
`answer = loaves_baked - loaves_sold_morning`  
`- loaves_sold_afternoon + loaves_returned`

`>>> print(answer)`  
74



# Tools: Program-aided LMs Prompts

Just need to *prompt* the right way

- Craft examples that interleave code and text

Q: Olivia has \$23. She bought five bagels for \$3 each. How much money does she have left?

```
money_initial = 23
bagels = 5
bagel_cost = 3
money_spent = bagels * bagel_cost
money_left = money_initial - money_spent
answer = money_left
```

Figure 3: Example prompt for the mathematical reasoning tasks, from the GSM8K benchmark.

Q: On the table, you see a bunch of objects arranged in a row: a purple paperclip, a pink stress ball, a brown keychain, a green scrunchiephone charger, a mauve fidget spinner, and a burgundy pen. What is the color of the object directly to the right of the stress ball?

```
...
stress_ball_idx = None
for i, object in enumerate(objects):
    if object[0] == 'stress ball':
        stress_ball_idx = i
        break
# Find the directly right object
direct_right = objects[stress_ball_idx+1]
# Check the directly right object's color
answer = direct_right[1]
```

# Tools: Program-of-Thoughts

## Similar idea: program-of-thoughts

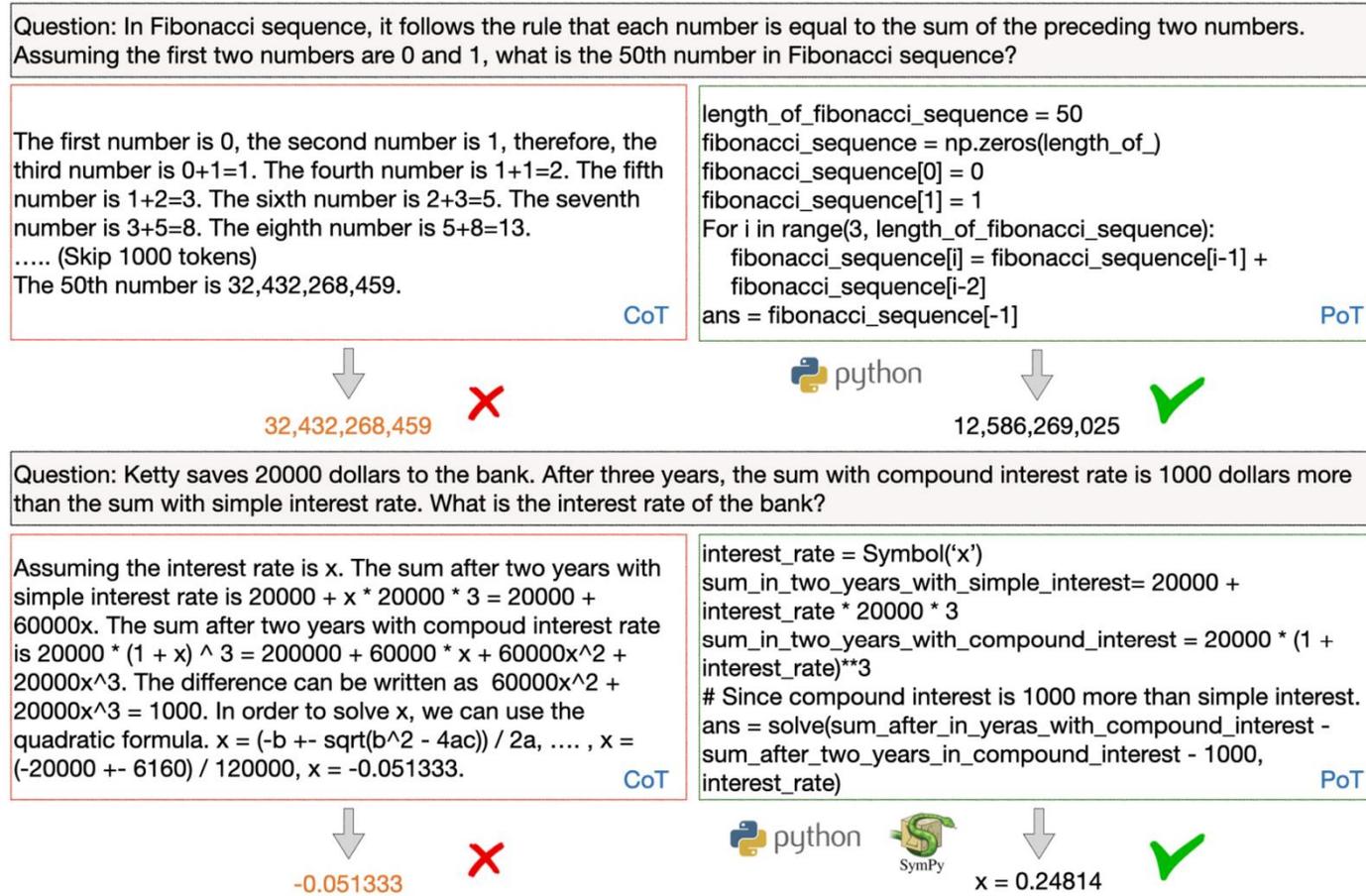


Figure 1: Comparison between Chain of Thoughts and Program of Thoughts.

# Tools: More General Tools

Ideally, use more general external tools

- Without lots of human annotation
- Model should decide on its own which tool to use
- **Toolformer**: introduces API calls into the model
  - But these API calls aren't already there... so need to **fine-tune**

*Your task is to add calls to a Question Answering API to a piece of text. The questions should help you get information required to complete the text. You can call the API by writing "[QA(question)]" where "question" is the question you want to ask. Here are some examples of API calls:*

**Input:** Joe Biden was born in Scranton, Pennsylvania.

**Output:** Joe Biden was born in [QA("Where was Joe Biden born?")] Scranton, [QA("In which state is Scranton?")] Pennsylvania.

**Input:** Coca-Cola, or Coke, is a carbonated soft drink manufactured by the Coca-Cola Company.

**Output:** Coca-Cola, or [QA("What other name is Coca-Cola known by?")] Coke, is a carbonated soft drink manufactured by [QA("Who manufactures Coca-Cola?")] the Coca-Cola Company.

**Input: x**

**Output:**



**Break & Questions**

# Outline

- **Last Time: Cross-Modal Adaptation, Tools**

- ORCA, aligning via optimal transport dataset distance, simple tool integration, Toolformer

- **Alignment and RLHF**

- Basic alignment idea, goals, mechanisms, RL review, RLHF steps

- **Why Does RLHF Work?**

- Failures of supervised learning, knowledge-seeking interactions, abstains

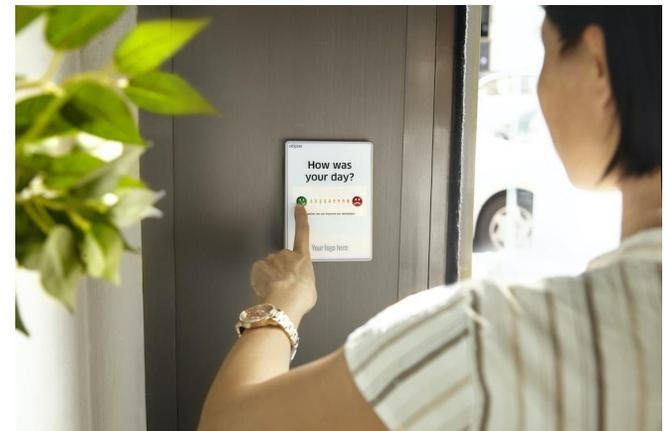
# Alignment: **Basic Motivation**

Goal: produce language model outputs that users like better...

- **Hard** to specify exactly what this means,
- **Easy** to query users

Collect human feedback and use it to change the model

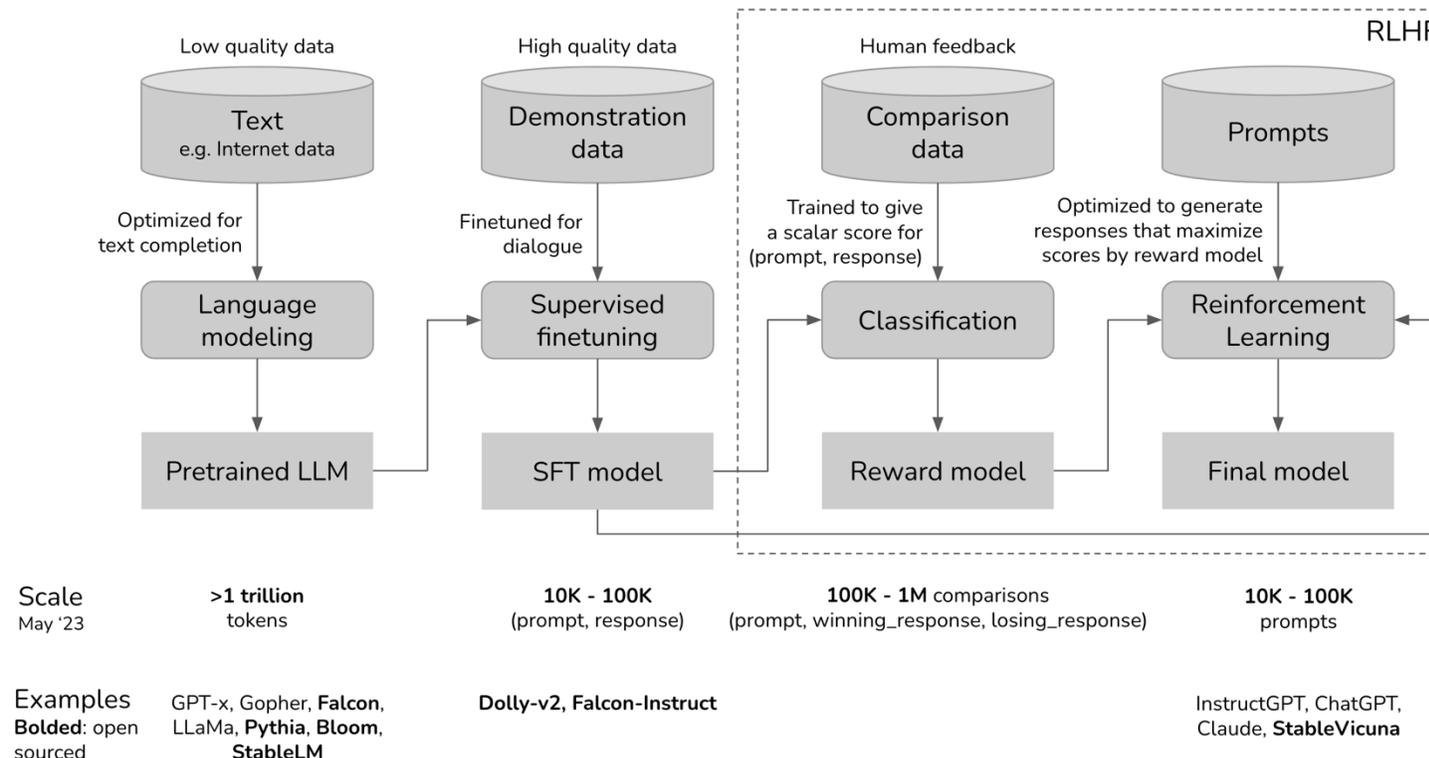
- Can do this by fine-tuning, especially with instructions
- Doesn't quite capture what users want
- We'll use other approaches, like RLHF



# RLHF: Setup

Goal: produce language model outputs that users like better...

- Via RL with trained reward model (Ouyang et al '22)

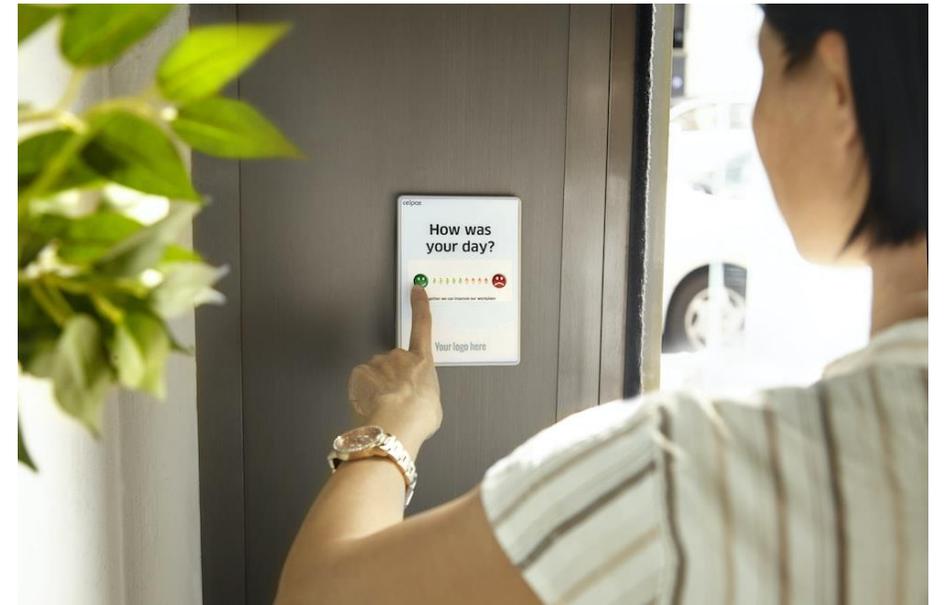


Chip Huyen

# RLHF: Feedback

First stage: get **human feedback** to train reward model

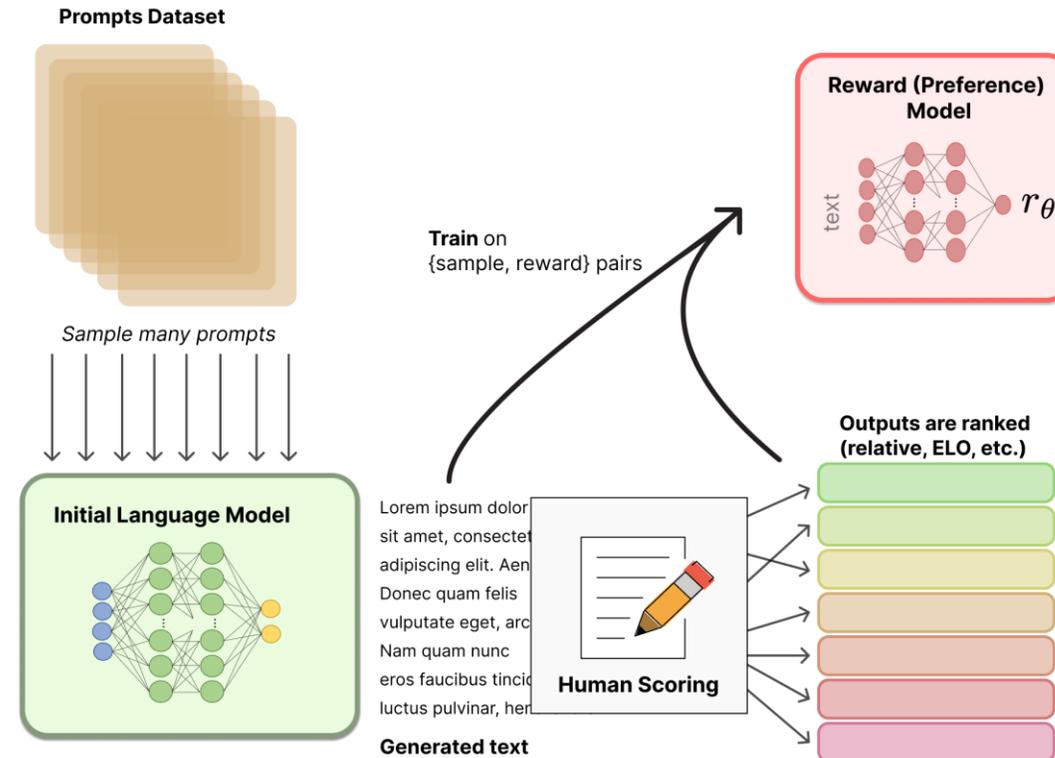
- Fix a set of prompts
- Produce multiple outputs for each prompt
  - Can get them from the original model post-SFT, or otherwise
- Ask human users **which is better**
  - **Binary output**
  - Can do more
    - Rank more questions



# RLHF: Reward/Preference Model

Second stage: train reward model

- Use the human feedback to train/fine-tune another model to reproduce the metric
- **Preference model**



<https://huggingface.co/blog/rlhf>

# RLHF: Reward/Preference Model

Second stage: train reward model

- Use the human feedback to train/fine-tune another model to reproduce the metric
- **Loss?** Based on preference models,
  - Example: Bradley-Terry model

$$p^*(y_1 \succ y_2 \mid x) = \frac{\exp(r^*(x, y_1))}{\exp(r^*(x, y_1)) + \exp(r^*(x, y_2))}.$$

- Then, our reward model loss is based on the log likelihood,

$$\mathcal{L}_R(r_\phi, \mathcal{D}) = -\mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\phi(x, y_w) - r_\phi(x, y_l))]$$

# RLHF: Reward/Preference Model

Note: we don't have to always do this from scratch

- Pretrained reward models available
- Benchmarks for this:



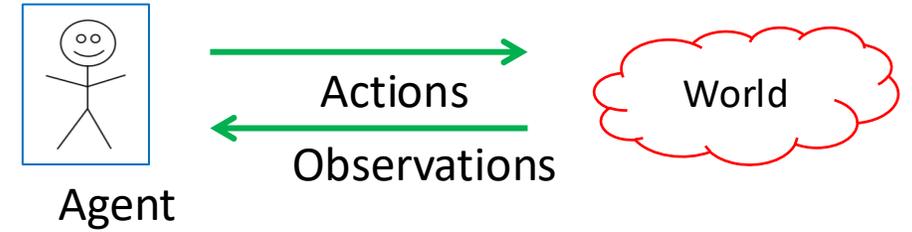
▲	Model	▲	Model Type	▲	Score	▲	Chat	▲	Chat Hard	▲
1	<a href="#">nvidia/Llama-3.1-Nemotron-70B-Reward</a>		Custom Classifier		94.1		97.5		85.7	
2	<a href="#">Skywork/Skywork-Reward-Gemma-2-27B</a>		Seq. Classifier		93.8		95.8		91.4	
3	<a href="#">SF-Foundation/TextEval-Llama3.1-70B</a>		Generative		93.5		94.1		90.1	
4	<a href="#">meta-metrics/MetaMetrics-RM-v1.0</a>		Custom Classifier		93.4		98.3		86.4	
5	<a href="#">Skywork/Skywork-Critic-Llama-3.1-70B</a>		Generative		93.3		96.6		87.9	
6	<a href="#">LxzGordon/URM-LLaMa-3.1-8B</a>		Seq. Classifier		92.9		95.5		88.2	
7	<a href="#">Salesforce/SFR-LLaMa-3.1-70B-Judge-r</a> *		Generative		92.7		96.9		84.8	
8	<a href="#">Skywork/Skywork-Reward-Llama-3.1-8B</a>		Seq. Classifier		92.5		95.8		87.3	

Lambert et al '24

# RLHF: Fine-Tuning with RL

Third stage: RL

- Use an RL algorithm
- **Goal:** produce outputs that have high reward

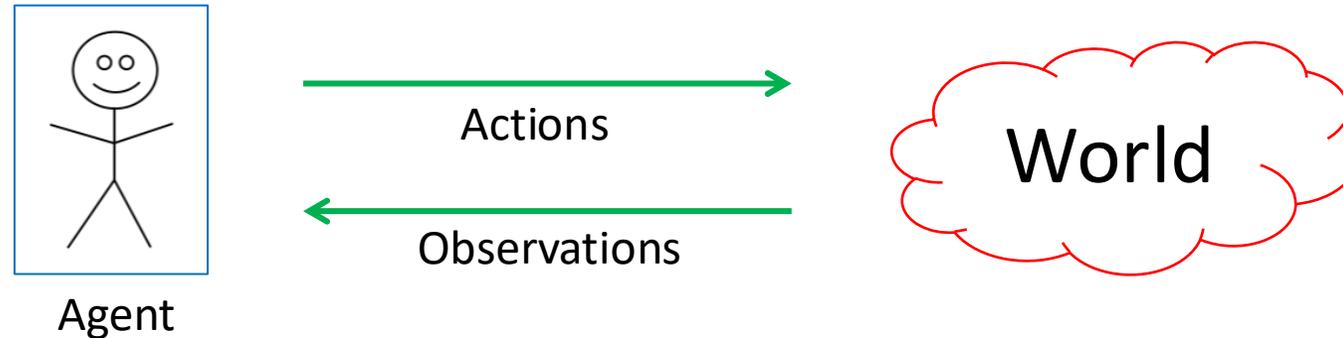


RL formulation:

- **Action space:** all the tokens possible to output
- **State space:** all the sequences of tokens
- **Reward function:** the trained reward model
- **Policy:** the new version of the LM, taking in state and returning tokens

# Reinforcement Learning Review

We have an **agent interacting** with the **world**

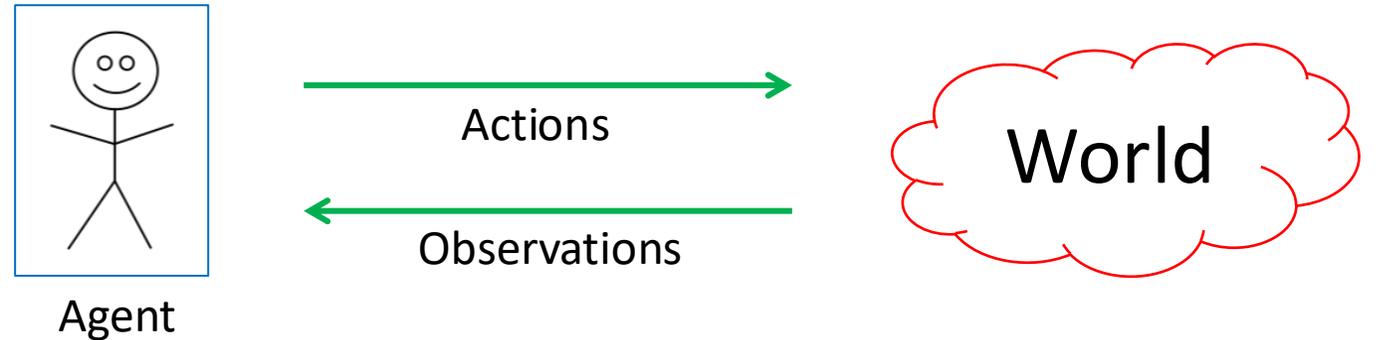


- Agent receives a reward based on state of the world
  - **Goal:** maximize reward / utility **(\$\$\$)**

# RL Review: Theoretical Model

## Basic setup:

- Set of states,  $S$
- Set of actions  $A$
- Information: at time  $t$ , observe state  $s_t \in S$ . Get reward  $r_t$
- Agent makes choice  $a_t \in A$ . State changes to  $s_{t+1}$ , continue



Goal: find a map from **states to actions** maximize rewards.

↑  
A “policy”

# RL Review: Markov Decision Process (MDP)

The formal mathematical model:

- **State set**  $S$ . Initial state  $s_0$ . **Action set**  $A$
- **State transition model:**  $P(s_{t+1} | s_t, a_t)$ 
  - Markov assumption: transition probability only depends on  $s_t$  and  $a_t$ , and not previous actions or states.
- **Reward function:**  $r(s_t)$
- **Policy:**  $\pi(s) : S \rightarrow A$  action to take at a particular state.

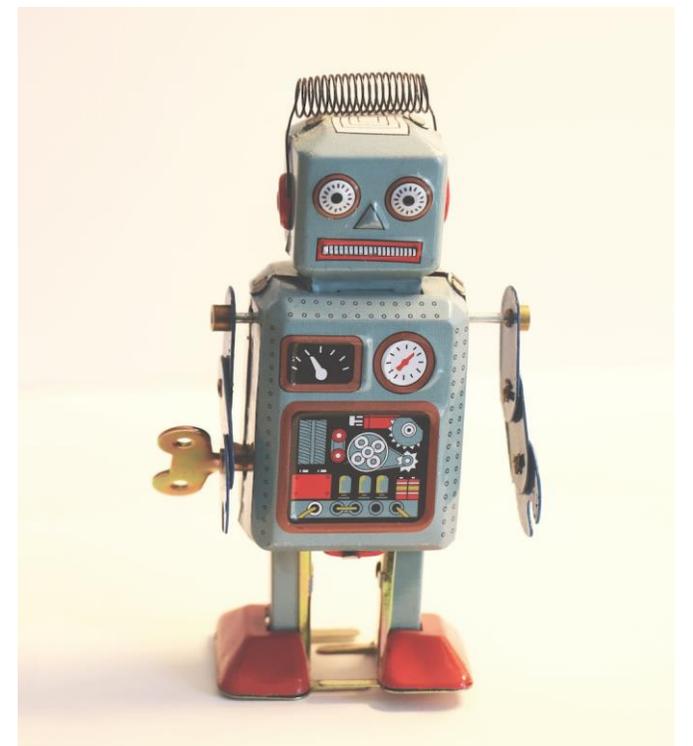
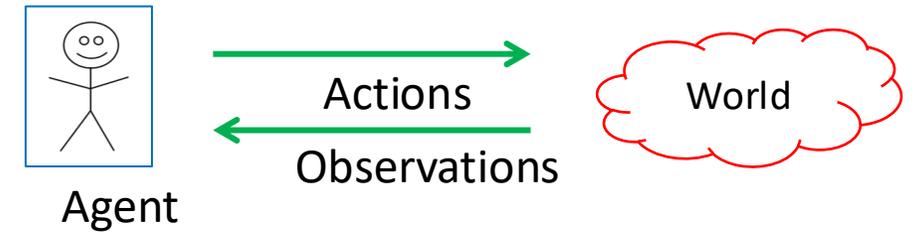
$$s_0 \xrightarrow{a_0} s_1 \xrightarrow{a_1} s_2 \xrightarrow{a_2} \dots$$

# RLHF: RL Approach

What approach for RL stage?

- Many deep RL methods available
- Policy gradient methods
- Popular: PPO (Proximal Policy Optimization)
  - Main difference from vanilla policy gradient, you constrain change to policy at each step (Schulman et al)

$$\max_{\pi_{\theta}} \mathbb{E}_{x \sim \mathcal{D}, y \sim \pi_{\theta}(y|x)} [r_{\phi}(x, y)] - \beta \mathbb{D}_{\text{KL}}[\pi_{\theta}(y|x) \parallel \pi_{\text{ref}}(y|x)]$$





**Break & Questions**

# Outline

- **Last Time: Cross-Modal Adaptation, Tools**

- ORCA, aligning via optimal transport dataset distance, simple tool integration, Toolformer

- **Alignment and RLHF**

- Basic alignment idea, goals, mechanisms, RL review, RLHF steps

- **Why Does RLHF Work?**

- Failures of supervised learning, knowledge-seeking interactions, abstains

# Why RLHF?

## Why should we do this?

- Why does supervised fine-tuning by itself not give our goal results?
- Many hypotheses; this section inspired by Yoav Goldberg's blog:
  - <https://gist.github.com/yoavg/6bff0feccd65950898eba1bb321cfbd81>
  - Itself based on Schulman's talk
  - [https://www.youtube.com/watch?v=hiLw5Q\\_UFg](https://www.youtube.com/watch?v=hiLw5Q_UFg)



# Why RLHF? Ways To Interact

Three “modes of interaction”:

- **text-grounded**: provide the model with text, instruction (“what are the chemical names mentioned in this text”),
- **knowledge-seeking**: provide the model with question or instruction, and expect a (truthful) answer based on the model's internal knowledge
- **creative**: provide the model with question or instruction, expect some creative output. (“Write a story about...”)

# Why RLHF? Knowledge-seeking

Three “modes of interaction”:

- **knowledge-seeking**: provide the model with question or instruction, and expect a (truthful) answer based on the model's internal knowledge
- This is hypothesized to require RL. Why does **SL fail**?
  - Case 1: know the answer: fine.
  - Case 2: don't know the answer. Supervised learning forces memorization, cannot produce “don't know”.
  - Worse, SL on case 2 encourages **model to lie**...



# Why RLHF? Knowledge-seeking with RL

Three “modes of interaction”:

- **knowledge-seeking**: provide the model with question or instruction, and expect a (truthful) answer based on the model's internal knowledge
- Why does RL succeed?
  - Case 1: know the answer: fine. Get a reward
  - Case 2: don't know the answer. Sometimes make it up and get a reward if lucky, most of the time low reward
  - **Encourages truth telling.**

# Why RLHF? **Abstains**

Additionally, **we'd like our model to abstain**

- SL will really struggle with this
  - Usually no abstains in datasets
  - Even if there were, “generalization” here means abstaining on similar questions? Difficult
- RL still challenging, need to produce high reward for “don't know”, but specific to model
- One way to craft a reward function:
  - High reward: correct answers
  - Medium reward: abstain
  - Negative reward: incorrect





**Thank You!**