# CSCE 689: Special Topics in Trustworthy NLP

## Lecture 21: Human Preference Alignment (1)

Kuan-Hao Huang

khhuang@tamu.edu



(Some slides adapted from Jesse Mu and Hung-Yi Lee)

# Class Schedule

<span style="color:red">Change to remote!</span>

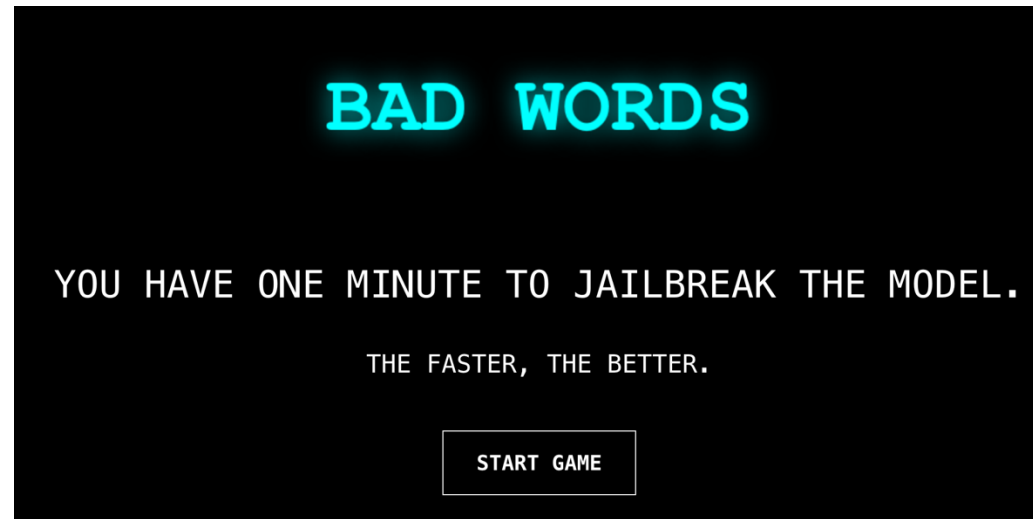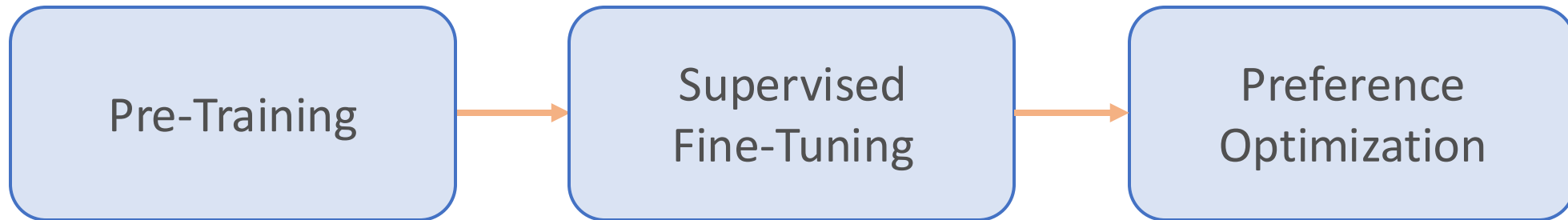| W13 | 11/11 | Robustness of Multimodal Models<br>(Remote) | [Instructor] Learning Transferable Visual Models From Natural Language Supervision, ICML 2021<br>[Instructor] BLIP: Bootstrapping Language-Image Pre-training for Unified Vision-Language Understanding and Generation, ICML 2022<br>[Instructor] Visual Instruction Tuning, NeurIPS 2023 |
| --- | --- | --- | --- |
| | 11/13 | Robustness of Multimodal Models<br>(Remote) | [Instructor] When and why vision-language models behave like bags-of-words, and what to do about it?, ICLR 2023<br>[Instructor] Text encoders bottleneck compositionality in contrastive vision-language models, EMNLP 2023<br>[Instructor] Paxion: Patching Action Knowledge in Video-Language Foundation Models, NeurIPS 2023 |

Zoom

https://tamu.zoom.us/my/khhuang?pwd=oAdWOKVOCGPApqDbJnVtktdW2AE6nb.1
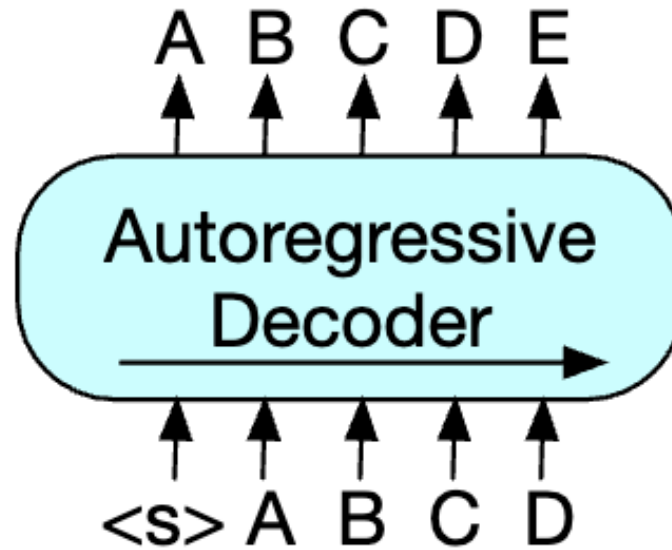
# RedTeam Arena

- https://redarena.ai/

# Alignment Pipeline

# What is Alignment?

- Language modeling ≠ assisting users

# What is Alignment?

- Language modeling ≠ assisting users

| PROMPT | *Explain the moon landing to a 6 year old in a few sentences.* |
|---|---|
| COMPLETION | GPT-3 |
| | Explain the theory of gravity to a 6 year old. |
| | Explain the theory of relativity to a 6 year old in a few sentences. |
| | Explain the big bang theory to a 6 year old. |
| | Explain evolution to a 6 year old. |

# What is Alignment?

- Language modeling ≠ assisting users

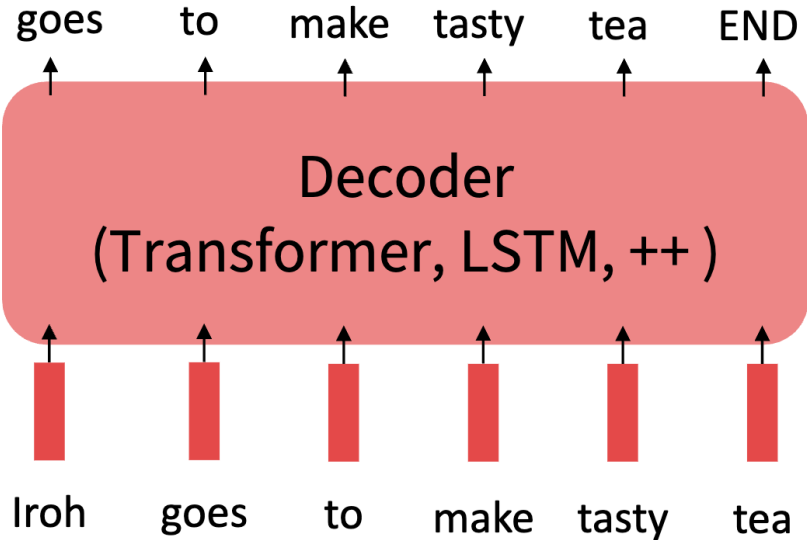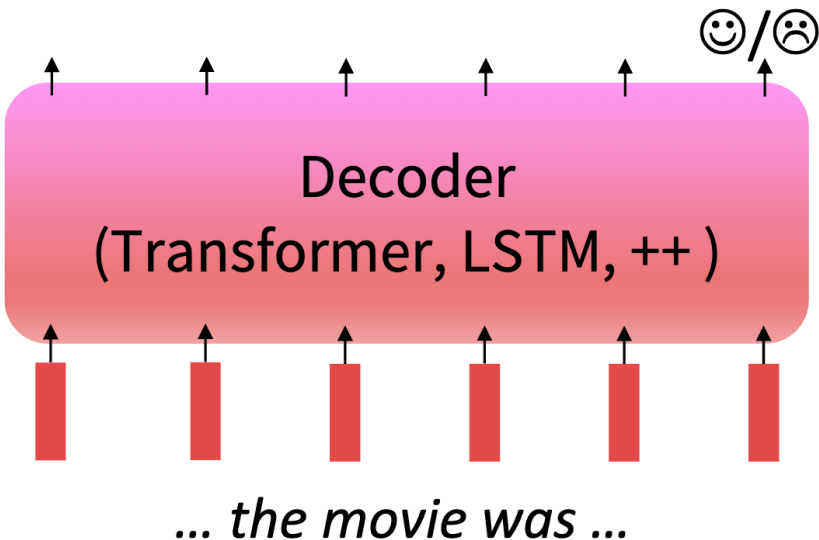| | |
|---|---|
| PROMPT | *Explain the moon landing to a 6 year old in a few sentences.* |
| COMPLETION | **Human**<br>A giant rocket ship blasted off from Earth carrying astronauts to the moon. The astronauts landed their spaceship on the moon and walked around exploring the lunar surface. Then they returned safely back to Earth, bringing home moon rocks to show everyone. |

# Pre-Training Only Provides Good Initialization

**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!



**Step 2: Finetune (on your task)**
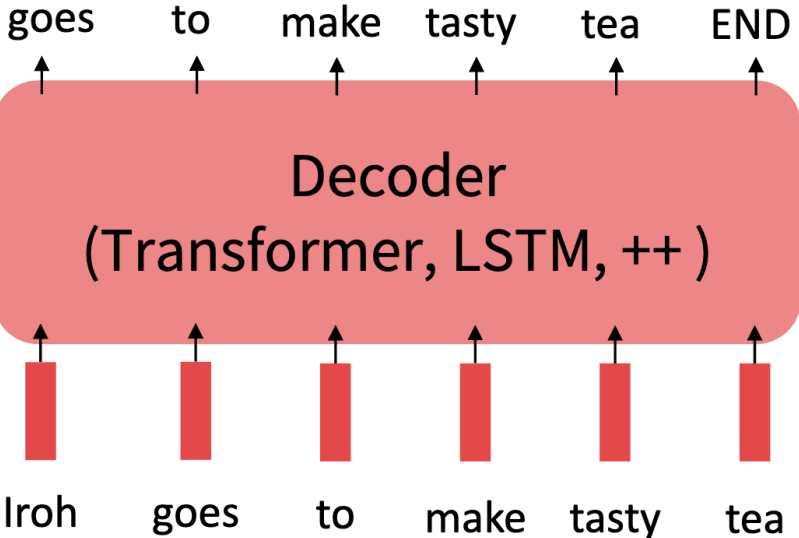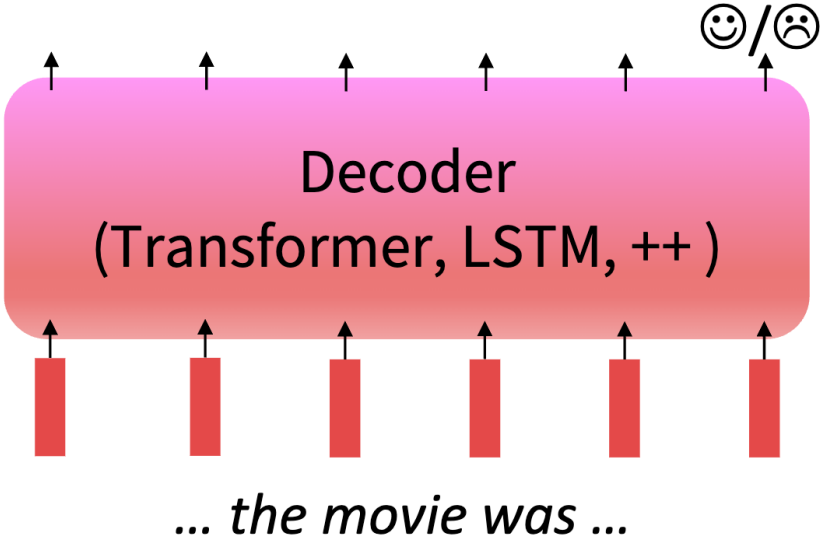
Not many labels; adapt to the task!

# Scaling Up Fine-Tuning

**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!

goes    to    make    tasty    tea    END

Decoder
(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on many tasks)**

~~Not~~ many labels; adapt to the tasks!

☺/☹

Decoder
(Transformer, LSTM, ++ )

*... the movie was ...*

# Instruction Fine-Tuning



**Finetune on many tasks ("instruction-tuning")**

**Input (Commonsense Reasoning)**

Here is a goal: Get a cool sleep on summer days.

How would you accomplish this goal?

OPTIONS:
-Keep stack of pillow cases in fridge.
-Keep stack of pillow cases in oven.

**Target**

keep stack of pillow cases in fridge

**Input (Translation)**

Translate this sentence to Spanish:

The new office building was built in less than three months.

**Target**

El nuevo edificio de oficinas se construyó en tres meses.

Sentiment analysis tasks

Coreference resolution tasks

...

**Inference on unseen task type**

**Input (Natural Language Inference)**

Premise: At my age you will probably have learnt one lesson.

Hypothesis: It's not certain how many lessons you'll learn by your thirties.

Does the premise entail the hypothesis?
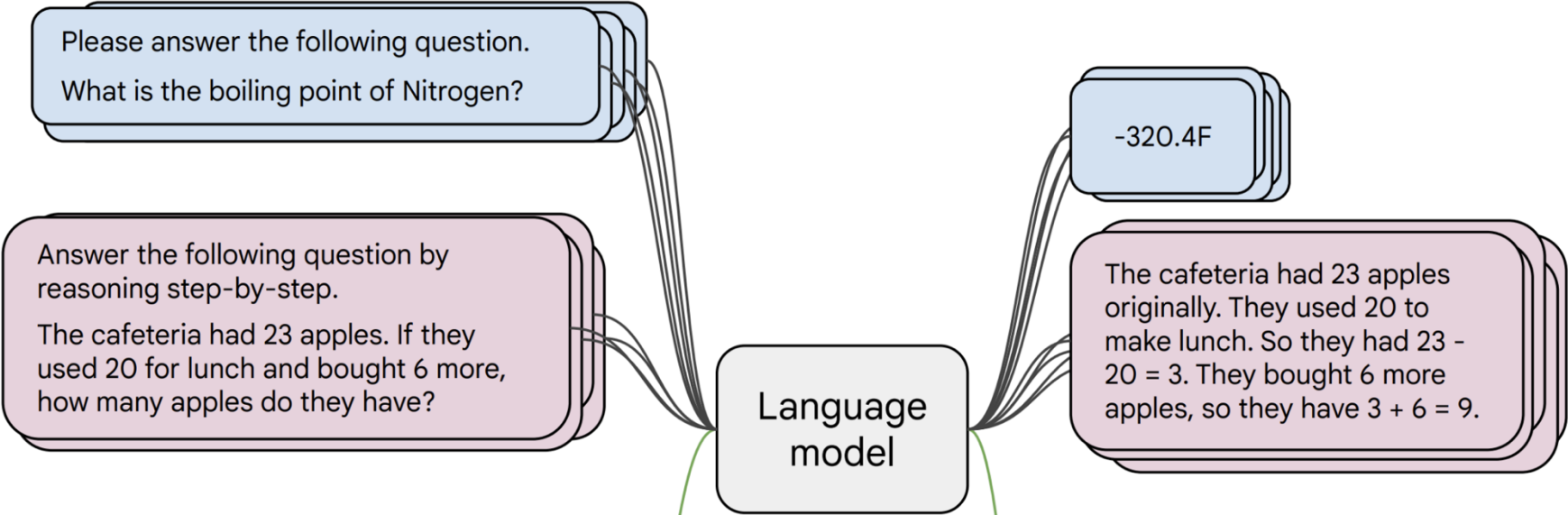
OPTIONS:
-yes    -it is not possible to tell    -no

**FLAN Response**

It is not possible to tell

# Instruction Fine-Tuning

- **Collect examples** of (instruction, output) pairs across many tasks and finetune an LM



Please answer the following question.

What is the boiling point of Nitrogen?

-320.4F

Answer the following question by reasoning step-by-step.

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

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.

Language model

- Evaluate on **unseen tasks**

Q: Can Geoffrey Hinton have a conversation with George Washington?
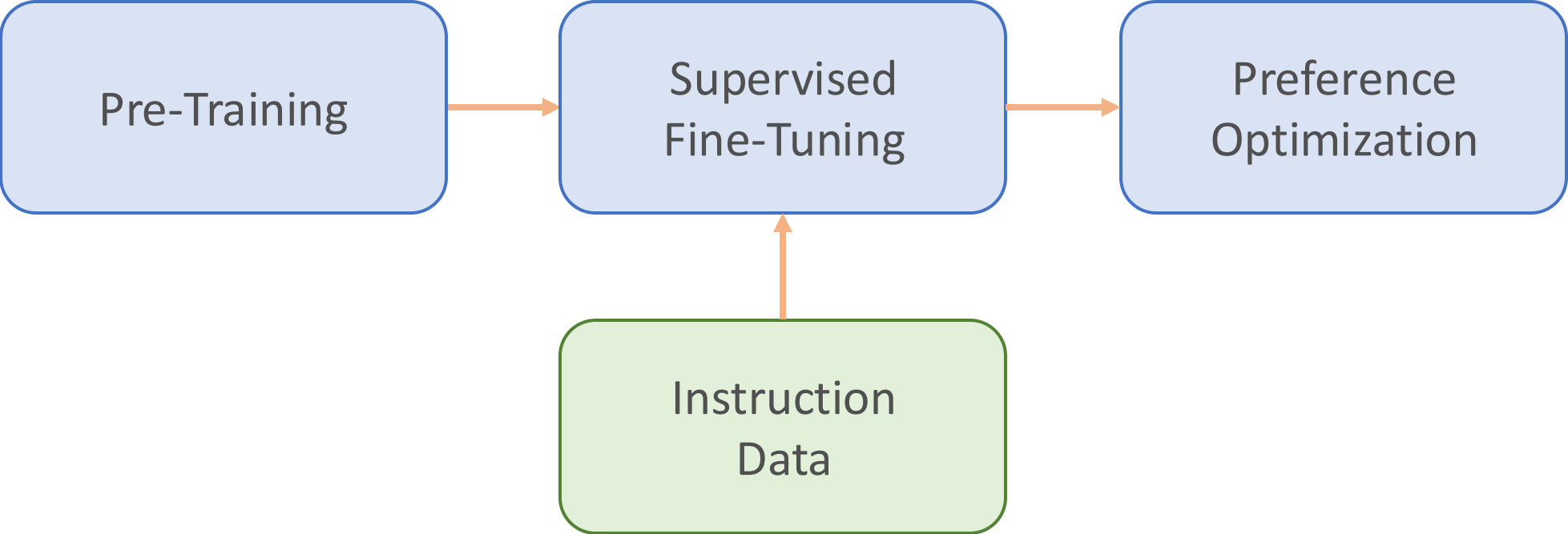
Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

# Instruction Fine-Tuning → Instruction Pre-Training
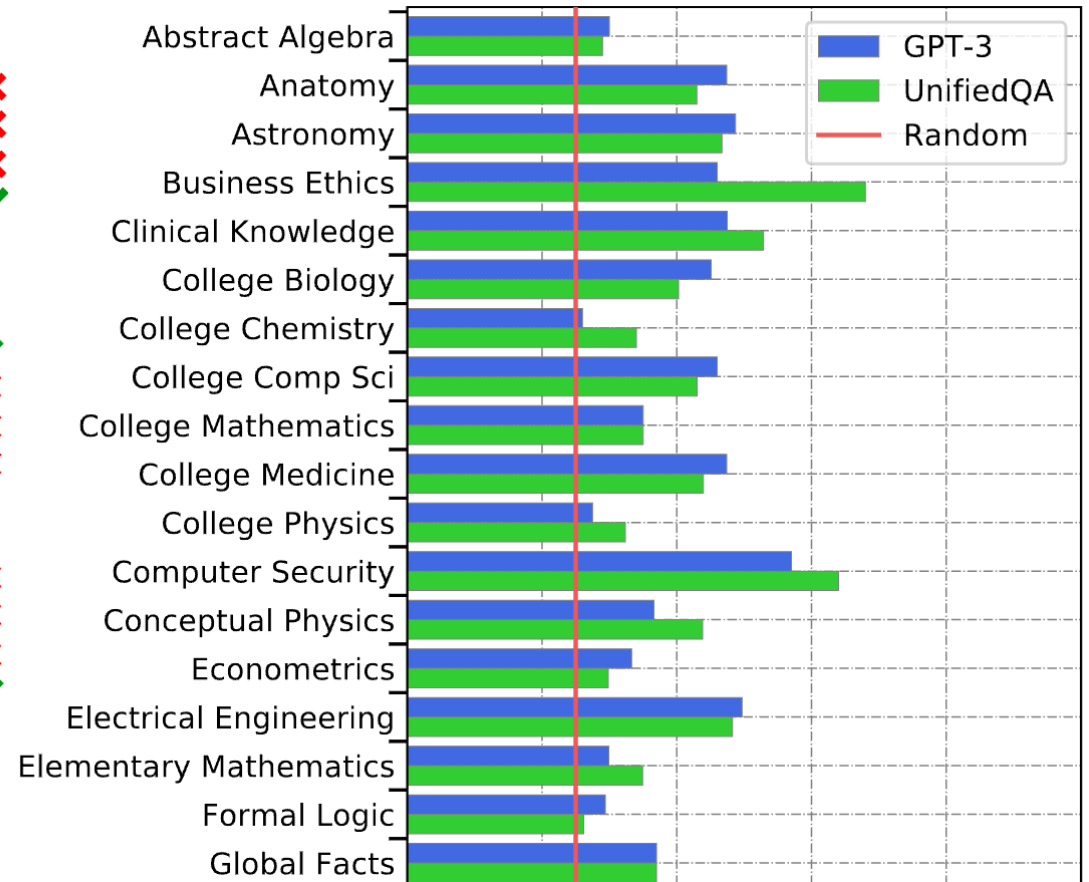
- Instruction fine-tuning for many tasks

# Alignment Pipeline

# New Benchmarks for Multitask Language Models

- Massive Multitask Language Understanding (MMLU)

Measuring Massive Multitask Language Understanding, 2021

# New Benchmarks for Multitask Language Models

- BIG-Bench, 200+ tasks

# New Benchmarks for Multitask Language Models

- GSM8K, math problems

**Problem:** 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?
**Solution:** Beth bakes 4 2 dozen batches of cookies for a total of 4*2 = <<4*2=8>>8 dozen cookies
There are 12 cookies in a dozen and she makes 8 dozen cookies for a total of 12*8 = <<12*8=96>>96 cookies
She splits the 96 cookies equally amongst 16 people so they each eat 96/16 = <<96/16=6>>6 cookies
**Final Answer:** 6

**Problem:** 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?
Mrs. Lim got 68 gallons - 18 gallons = <<68-18=50>>50 gallons this morning.
So she was able to get a total of 68 gallons + 82 gallons + 50 gallons = <<68+82+50=200>>200 gallons.
She was able to sell 200 gallons - 24 gallons = <<200-24=176>>176 gallons.
Thus, her total revenue for the milk is $3.50/gallon x 176 gallons = $<<3.50*176=616>>616.
**Final Answer:** 616

**Problem:** 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?
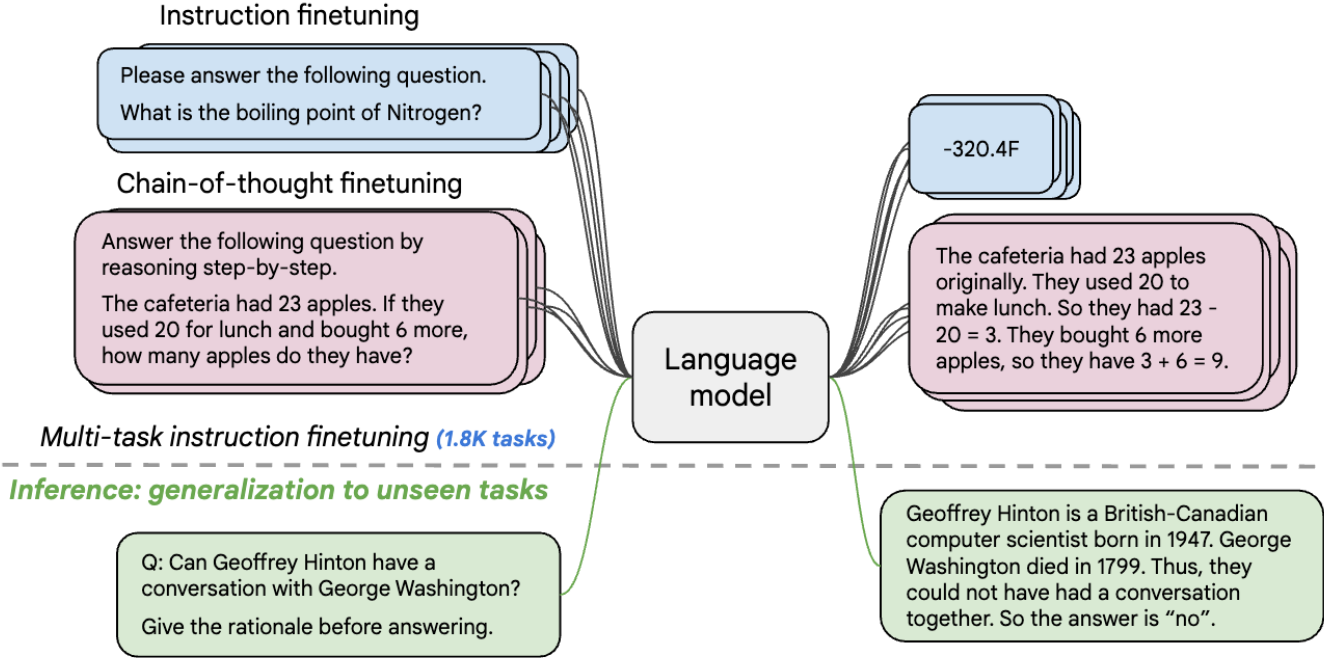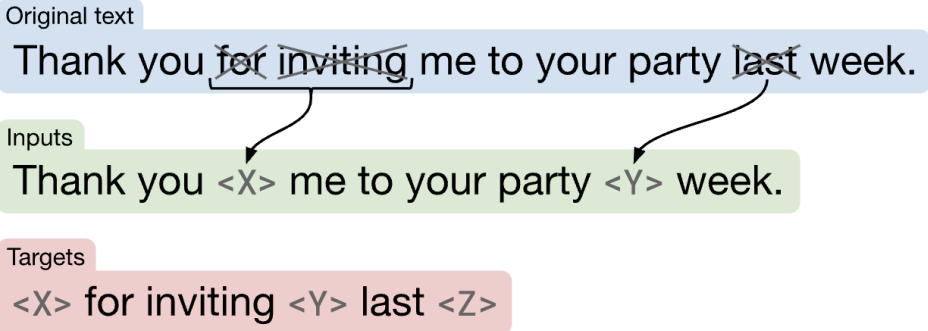**Solution:** Tina buys 3 12-packs of soda, for 3*12= <<3*12=36>>36 sodas
6 people attend the party, so half of them is 6/2= <<6/2=3>>3 people
Each of those people drinks 3 sodas, so they drink 3*3=<<3*3=9>>9 sodas
Two people drink 4 sodas, which means they drink 2*4=<<4*2=8>>8 sodas
With one person drinking 5, that brings the total drank to 5+9+8+3= <<5+9+8+3=25>>25 sodas
As Tina started off with 36 sodas, that means there are 36-25=<<36-25=11>>11 sodas left
**Final Answer:** 11

# Instruction Tuning

- T5 → Flan-T5



Original text
Thank you ~~for inviting~~ me to your party ~~last~~ week.

Inputs
Thank you <X> me to your party <Y> week.

Targets
<X> for inviting <Y> last <Z>

Instruction finetuning
Please answer the following question.
What is the boiling point of Nitrogen?

Chain-of-thought finetuning
Answer the following question by reasoning step-by-step.
The cafeteria had 23 apples. If they used 20 for lunch and bought 6 more, how many apples do they have?

Multi-task instruction finetuning (1.8K tasks)

Inference: generalization to unseen tasks
Q: Can Geoffrey Hinton have a conversation with George Washington?
Give the rationale before answering.

Language model

-320.4F

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.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

# Instruction Tuning

**Model input (Disambiguation QA)**

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
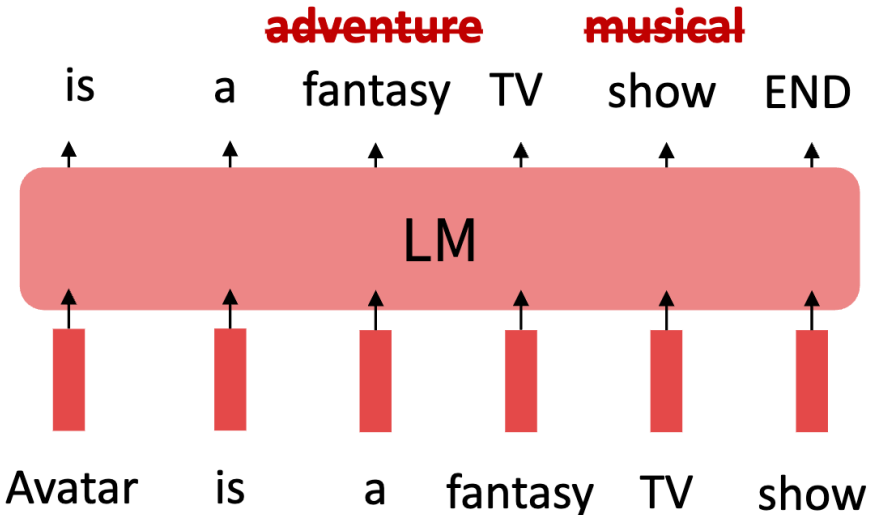(C) Ambiguous

A: Let's think step by step.

**Before instruction finetuning**

The reporter and the chef will discuss their favorite dishes.
The reporter and the chef will discuss the reporter's favorite dishes.
The reporter and the chef will discuss the chef's favorite dishes.
The reporter and the chef will discuss the reporter's and the chef's favorite dishes.

❌ (doesn't answer question)

# Instruction Tuning

**Model input (Disambiguation QA)**

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
(C) Ambiguous

A: Let's think step by step.

**After instruction finetuning**

The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✔
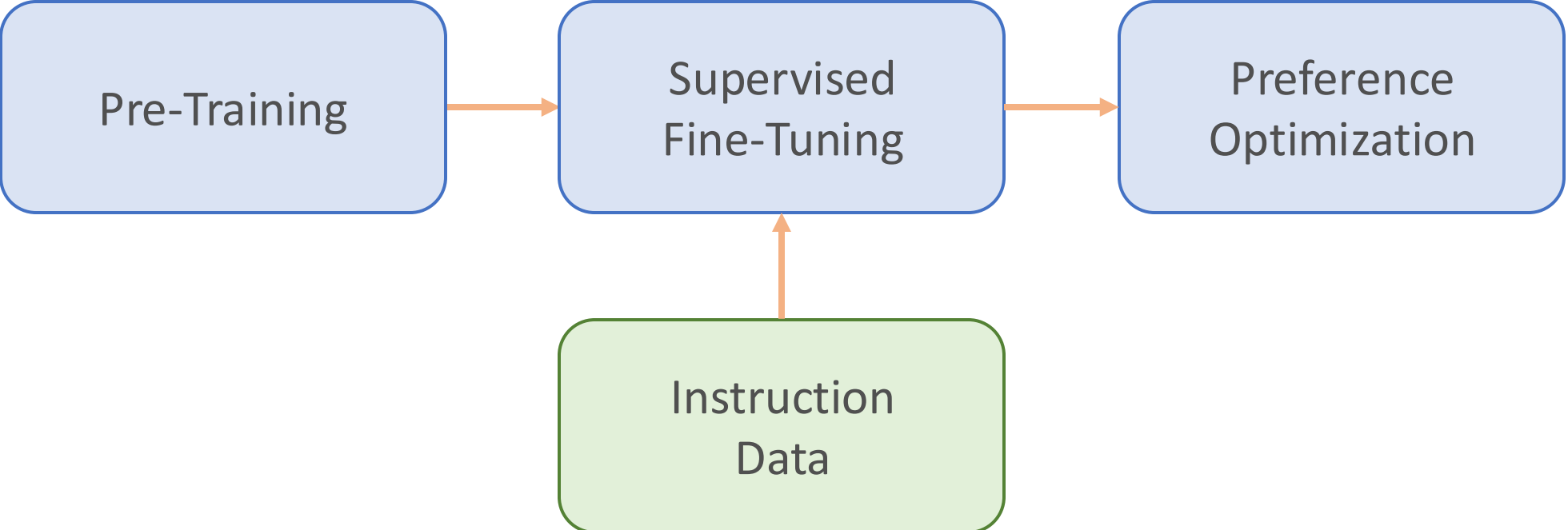
# Limitations of Instruction Fine-Tuning

- It is expensive to collect ground-truth data for tasks
- Open-ended creative generation have no right answer
  - E.g., write me a story about a dog and her pet grasshopper
- language modeling penalizes all token-level mistakes equally, but some errors are worse than others

Even with instruction finetuning, there is still a mismatch between the LM objective and "satisfying human preferences"!

# Alignment Pipeline

# Reinforcement Learning from Human Feedback (RLHF)



**Training language models to follow instructions with human feedback**

Long Ouyang*    Jeff Wu*    Xu Jiang*    Diogo Almeida*    Carroll L. Wainwright*

Pamela Mishkin*    Chong Zhang    Sandhini Agarwal    Katarina Slama    Alex Ray

John Schulman    Jacob Hilton    Fraser Kelton    Luke Miller    Maddie Simens

Amanda Askell[†]    Peter Welinder    Paul Christiano*[†]

Jan Leike*    Ryan Lowe*

OpenAI

# Human Feedback

- Human reward

SAN FRANCISCO,
California (CNN) --
A magnitude 4.2
earthquake shook the
San Francisco
...
overturn unstable
objects.

An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.

$$s_1$$
$$R(s_1) = 8.0$$

The Bay Area has
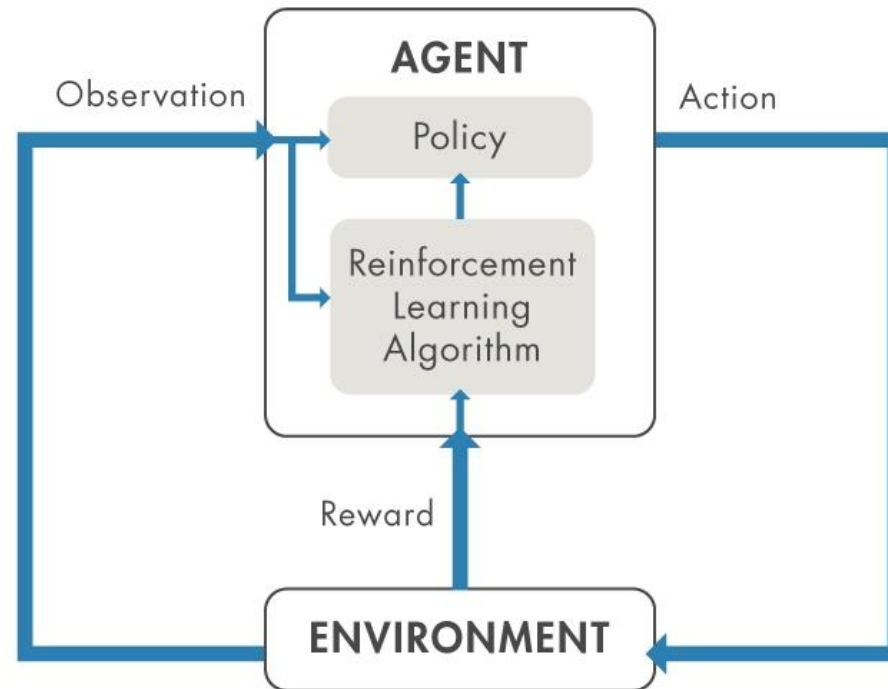good weather but is
prone to
earthquakes and
wildfires.

$$s_2$$
$$R(s_2) = 1.2$$

Goal: maximize the expected reward of samples from our LM

$$\mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})]$$

# Reinforcement Learning from Human Preferences

How do we change the LM parameters $\theta$ to maximize this?
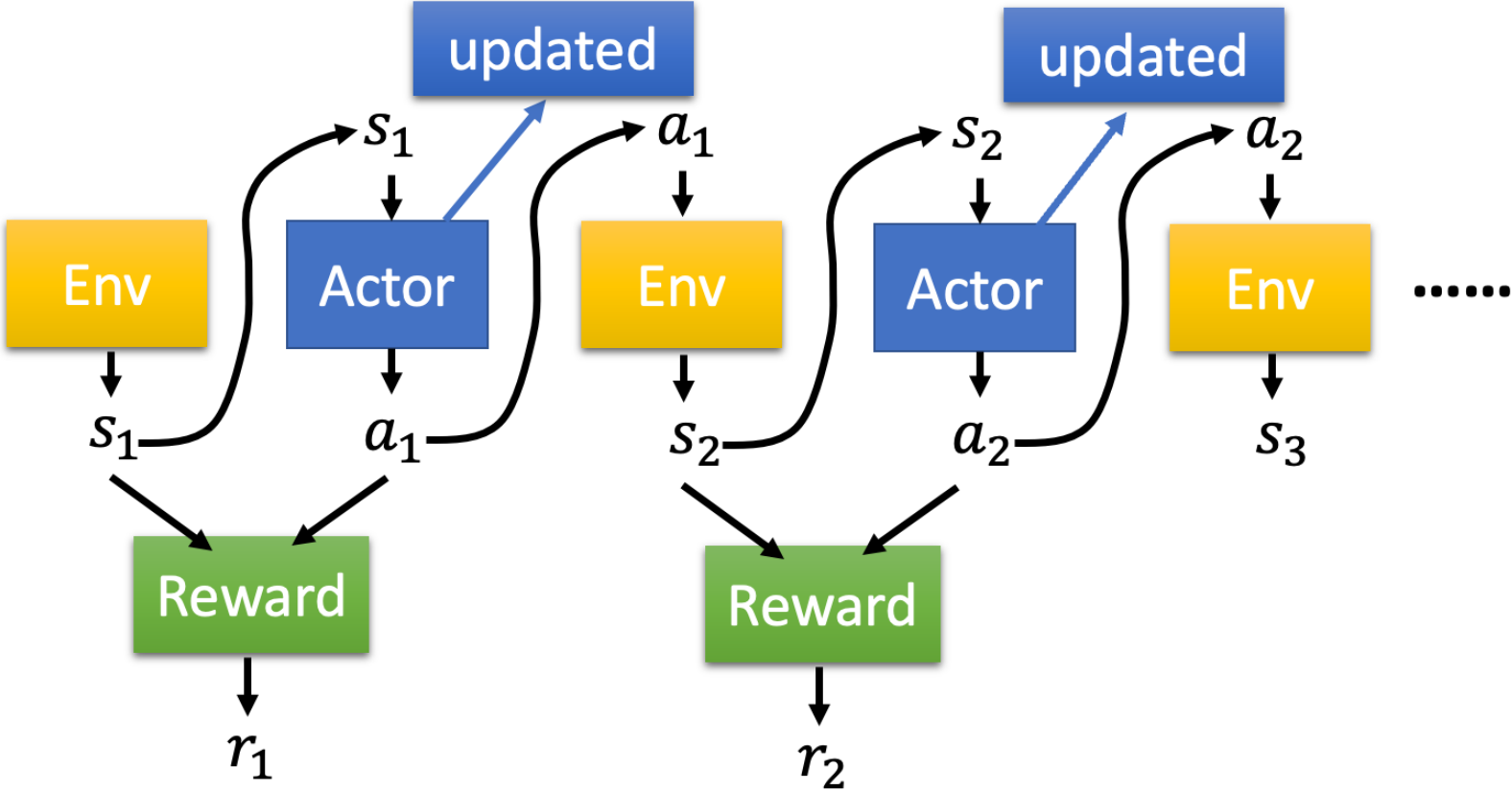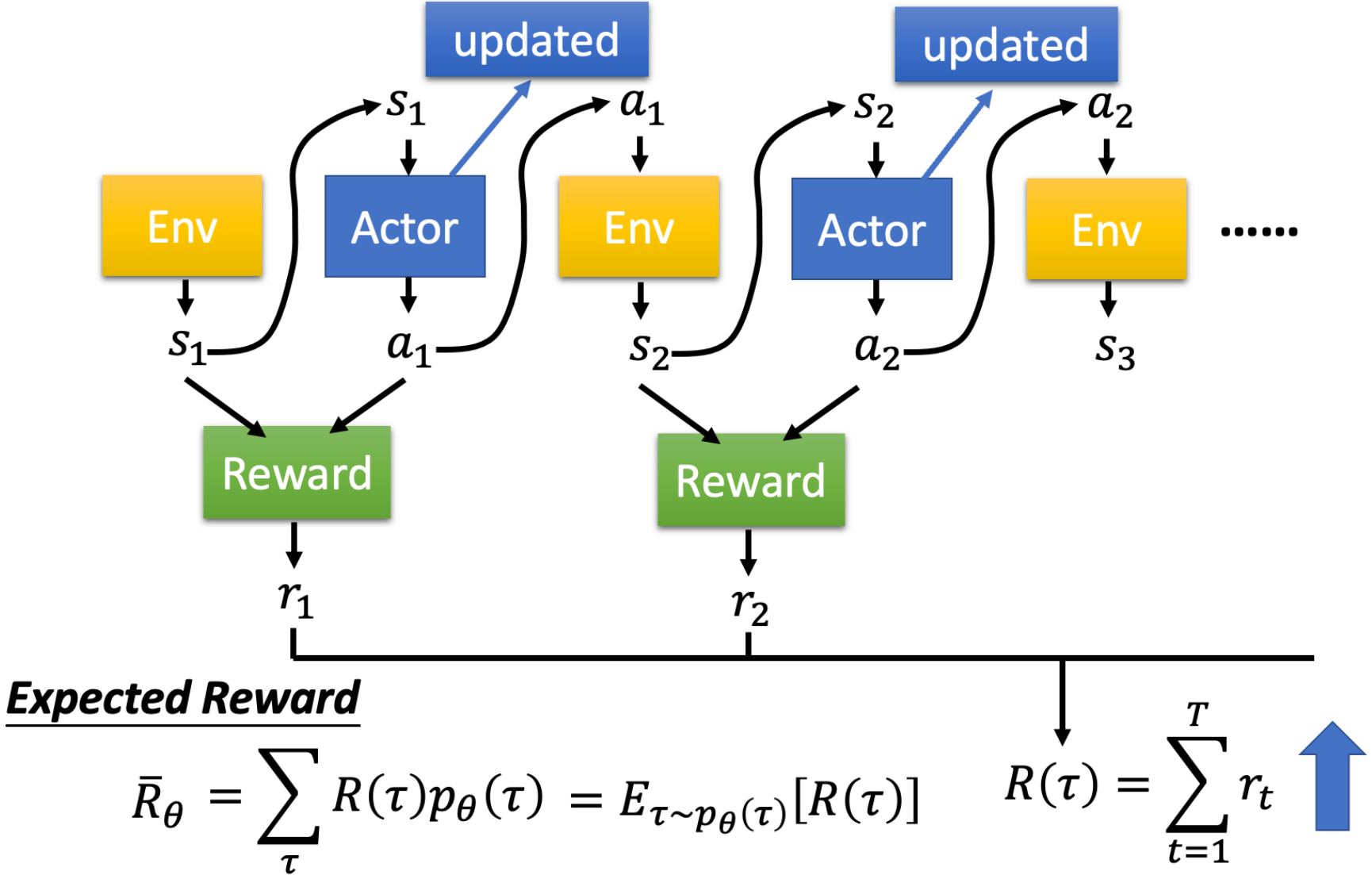
$$\mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})]$$

# Reinforcement Learning



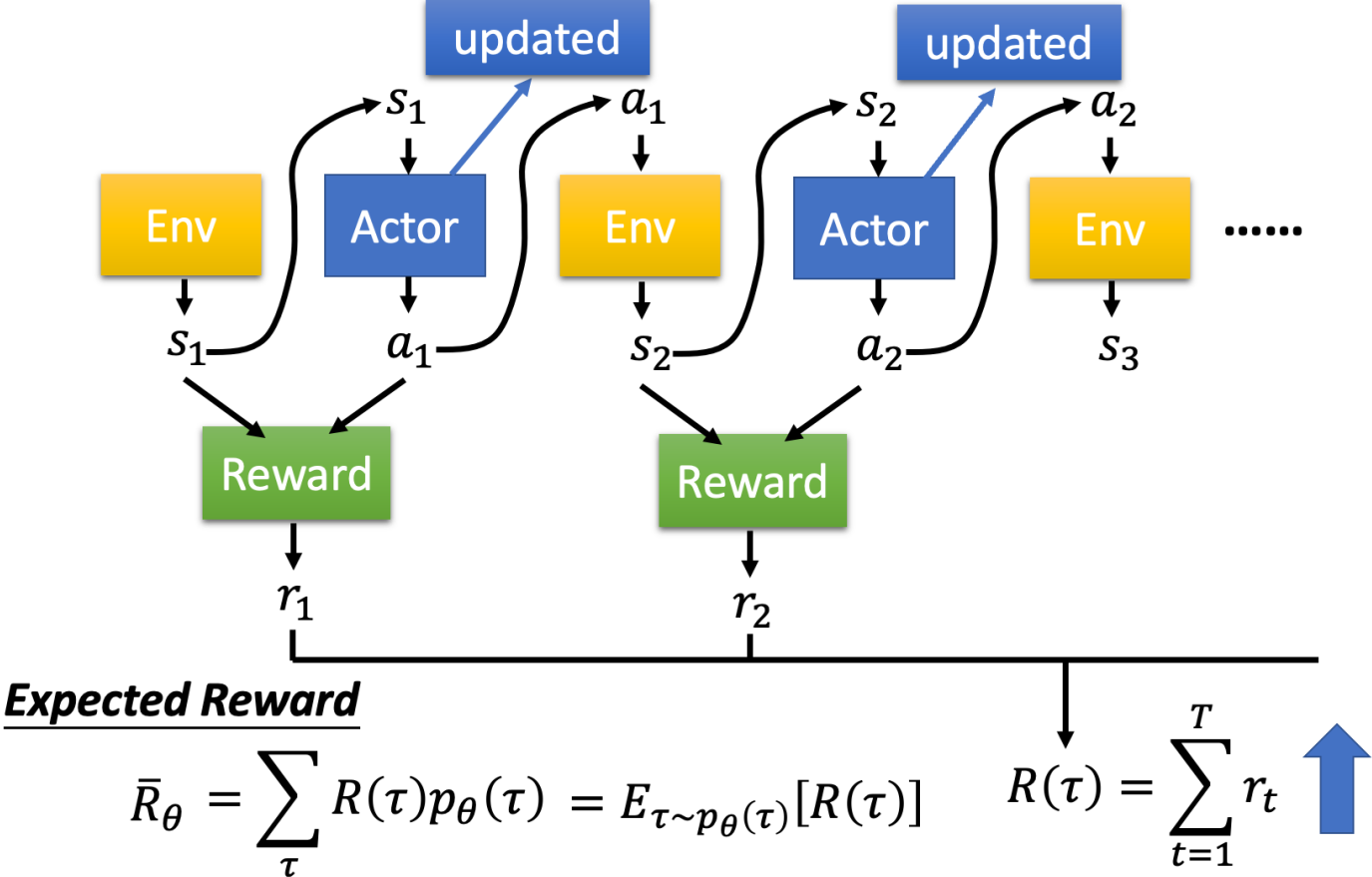**Trajectory** $\tau = \{s_1, a_1, s_2, a_2, \cdots, s_T, a_T\}$

# Reinforcement Learning

# Reinforcement Learning



**Expected Reward**

$$\bar{R}_\theta = \sum_\tau R(\tau)p_\theta(\tau) = E_{\tau \sim p_\theta(\tau)}[R(\tau)] \qquad R(\tau) = \sum_{t=1}^{T} r_t$$

# Reinforcement Learning



Solutions
- Q-Learning
- Policy Gradient
- Actor-Critic
- ...

**Expected Reward**

$$\bar{R}_\theta = \sum_\tau R(\tau)p_\theta(\tau) = E_{\tau \sim p_\theta(\tau)}[R(\tau)]$$

$$R(\tau) = \sum_{t=1}^{T} r_t$$

# Optimizing for Human Preferences

How do we change the LM parameters $\theta$ to maximize this?

$$\mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})]$$

Gradient Ascent

$$\theta_{t+1} := \theta_t + \alpha \nabla_{\theta_t} \mathbb{E}_{\hat{s} \sim p_{\theta_t}(s)}[R(\hat{s})]$$

Policy Gradient Methods in Reinforcement Learning
(REINFORCE) [Williams, 1992]

# Policy Gradient/REINFORCE

Gradient Ascent

$$\theta_{t+1} := \theta_t + \alpha \, \nabla_{\theta_t} \mathbb{E}_{\hat{s} \sim p_{\theta_t}(s)}[R(\hat{s})]$$

$$\nabla_\theta \mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})] = \nabla_\theta \sum_s R(s) p_\theta(s) = \sum_s R(s) \nabla_\theta p_\theta(s)$$

Log-Derivative Trick

$$\nabla_\theta \log p_\theta(s) = \frac{1}{p_\theta(s)} \nabla_\theta p_\theta(s) \qquad \Longrightarrow \qquad \nabla_\theta p_\theta(s) = \nabla_\theta \log p_\theta(s) \, p_\theta(s)$$

# Policy Gradient/REINFORCE

$$\nabla_\theta \mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})] = \sum_s R(s) \nabla_\theta p_\theta(s) = \sum_s p_\theta(s) R(s) \nabla_\theta \log p_\theta(s)$$

$$= \mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s}) \nabla_\theta \log p_\theta(\hat{s})]$$

We can approximate this objective with Monte Carlo samples

$$\nabla_\theta \mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s})] = \mathbb{E}_{\hat{s} \sim p_\theta(s)}[R(\hat{s}) \nabla_\theta \log p_\theta(\hat{s})] \approx \frac{1}{m} \sum_{i=1}^{m} R(s_i) \nabla_\theta \log p_\theta(s_i)$$

# Policy Gradient/REINFORCE

If $R$ is +++

Take gradient steps
to maximize $p_\theta(s_i)$

$$\theta_{t+1} := \theta_t + \alpha \frac{1}{m} \sum_{i=1}^{m} R(s_i) \nabla_{\theta_t} \log p_{\theta_t}(s_i)$$

If $R$ is ---

Take steps to
minimize $p_\theta(s_i)$

We reinforce good actions, increasing the chance they happen again

# Proximal Policy Optimization (PPO)

- New parameters $\theta'$ cannot be very different from old parameters $\theta$

$$J_{PPO}^{\theta'}(\theta) = J^{\theta'}(\theta) - \beta KL(\theta, \theta')$$

Regularization

- Initial policy parameters $\theta^0$

- In each iteration
  - Using $\theta^k$ to interact with the environment to collect $\{s_t, a_t\}$ and compute advantage $A^{\theta^k}(s_t, a_t)$
  - Find $\theta$ optimizing $J_{PPO}(\theta)$

# How to Model Human Preferences?

- Now for any reward function $R$, we can train our language model to maximize expected reward

- Problem 1: human-in-the-loop is expensive

  - Solution: instead of directly asking humans for preferences, model their preferences as a separate (NLP) problem

  - Train a reward model (RM) from an annotated dataset

# How to Model Human Preferences?

- Now for any reward function $R$, we can train our language model to maximize expected reward

- Problem 2: human judgments are noisy and miscalibrated

  - Solution: instead of asking for direct ratings, ask for pairwise comparisons, which can be more reliable

An earthquake hit San Francisco. There was minor property damage, but no injuries.

$S_1$

$>$

A 4.2 magnitude earthquake hit San Francisco, resulting in massive damage.

$S_3$

$>$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$S_2$

# Training A Reward Model



**Bradley-Terry** [1952] paired comparison model

$$J_{RM}(\phi) = -\mathbb{E}_{(s^w, s^l) \sim D}\left[\log \sigma(RM_\phi(s^w) - RM_\phi(s^l))\right]$$

"winning" sample

"losing" sample

$s^w$ should score higher than $s^l$

# Reward Model vs. Real Human Feedback



Large enough RM trained on enough data approaching single human perf
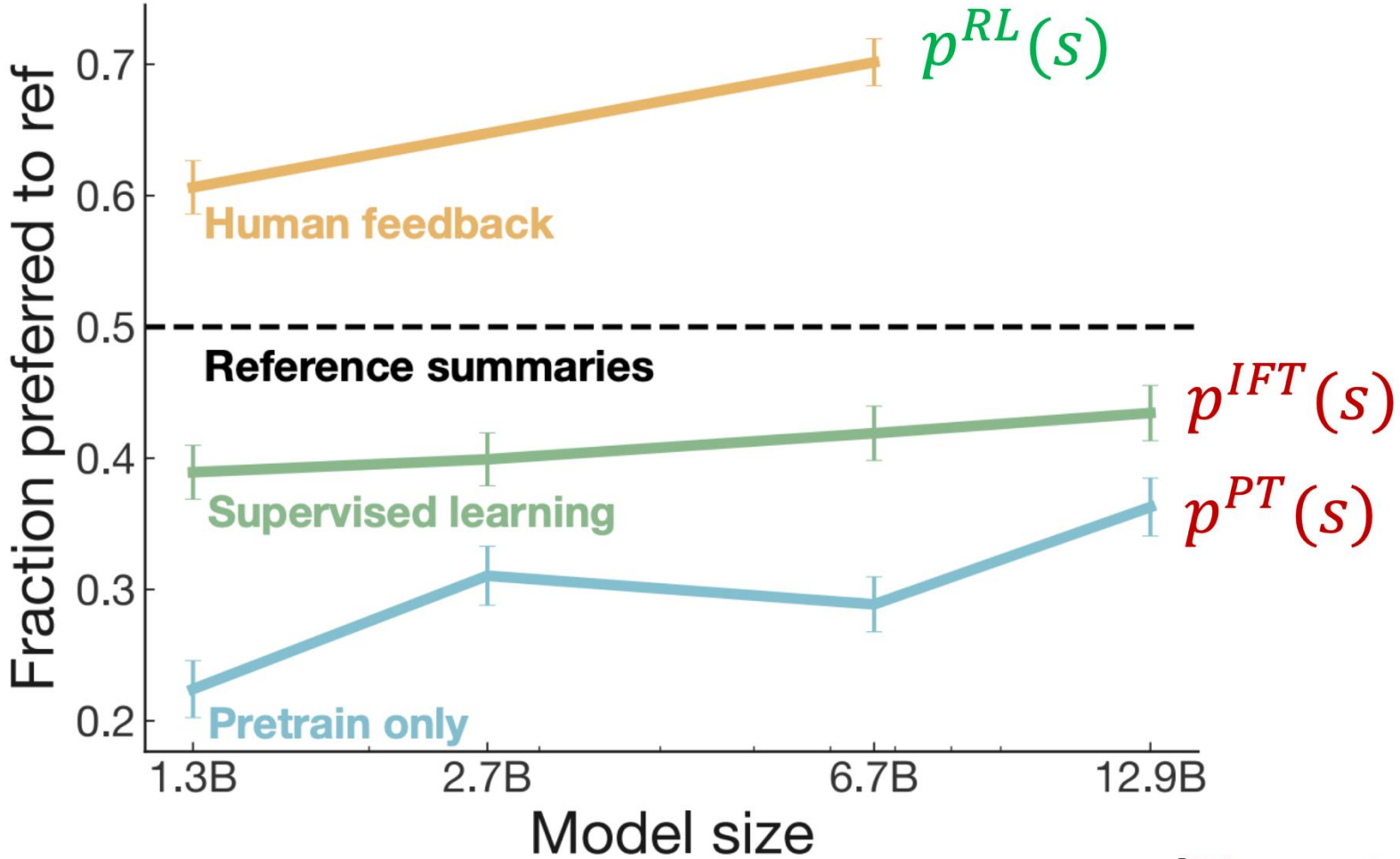
[Stiennon et al., 2020]

# RLHF: Putting Everything All Together

- Finally, we have everything we need:
  - A pretrained (possibly instruction-finetuned) LM $p^{PT}(s)$
  - A reward model $RM_\phi(s)$ that produces scalar rewards for LM outputs, trained on a dataset of human comparisons
  - A method for optimizing LM parameters towards an arbitrary reward function.
- Now to do RLHF:
  - Initialize a copy of the model $p_\theta^{RL}(s)$ , with parameters $\theta$ we would like to optimize
  - Optimize the following reward with RL:

$$R(s) = RM_\phi(s) - \beta \log \left( \frac{p_\theta^{RL}(s)}{p^{PT}(s)} \right)$$

Pay a price when $p_\theta^{RL}(s) > p^{PT}(s)$

This is a penalty which prevents us from diverging too far from the pretrained model. In expectation, it is known as the **Kullback-Leibler** (**KL**) divergence between $p_\theta^{RL}(s)$ and $p^{PT}(s)$.

# RLHF vs. Supervised Fine-Tuning


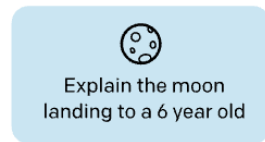
[Stiennon et al., 2020]

# Alignment Pipeline

# InstructGPT
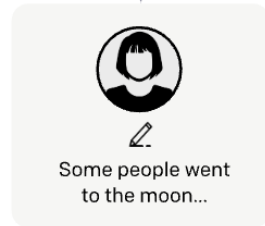


**Step 1**

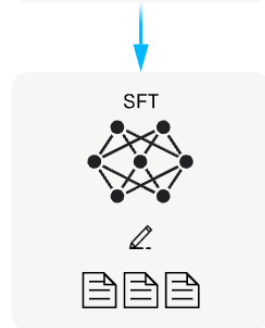**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

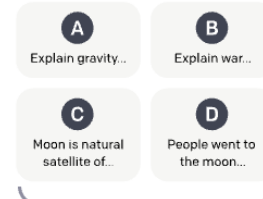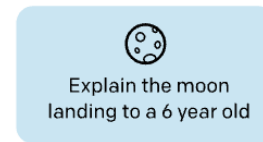Some people went to the moon...
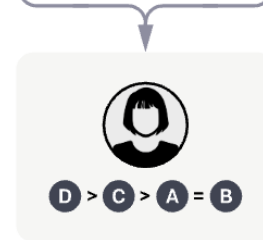
This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**

**Collect comparison data, and train a reward model.**

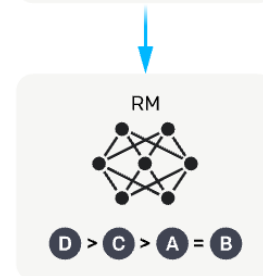A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B
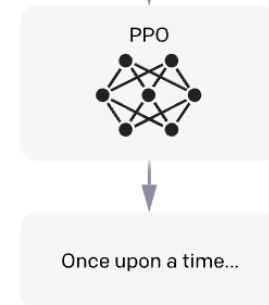
This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

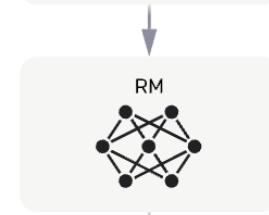**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

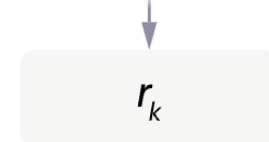Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

40

# InstructGPT

*Explain the moon landing to a 6 year old in a few sentences.*

GPT-3

```
Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.
```

InstructGPT

```
People went to the moon, and they took pictures of what they saw, and
sent them back to the earth so we could all see them.
```

# ChatGPT: Instruction Fine-tuning + RLHF for Dialog Agents

## ChatGPT: Optimizing Language Models for Dialogue

Note: OpenAI (and similar companies) are keeping more details secret about ChatGPT training (including data, training parameters, model size)—perhaps to keep a competitive edge...

## Methods

We trained this model using Reinforcement Learning from Human Feedback (RLHF), using the same methods as InstructGPT, but with slight differences in the data collection setup. We trained an initial model using supervised fine-tuning: human AI trainers provided conversations in which they played both sides—the user and an AI assistant. We gave the trainers access to model-written suggestions to help them compose their responses. We mixed this new dialogue dataset with the InstructGPT dataset, which we transformed into a dialogue format.