

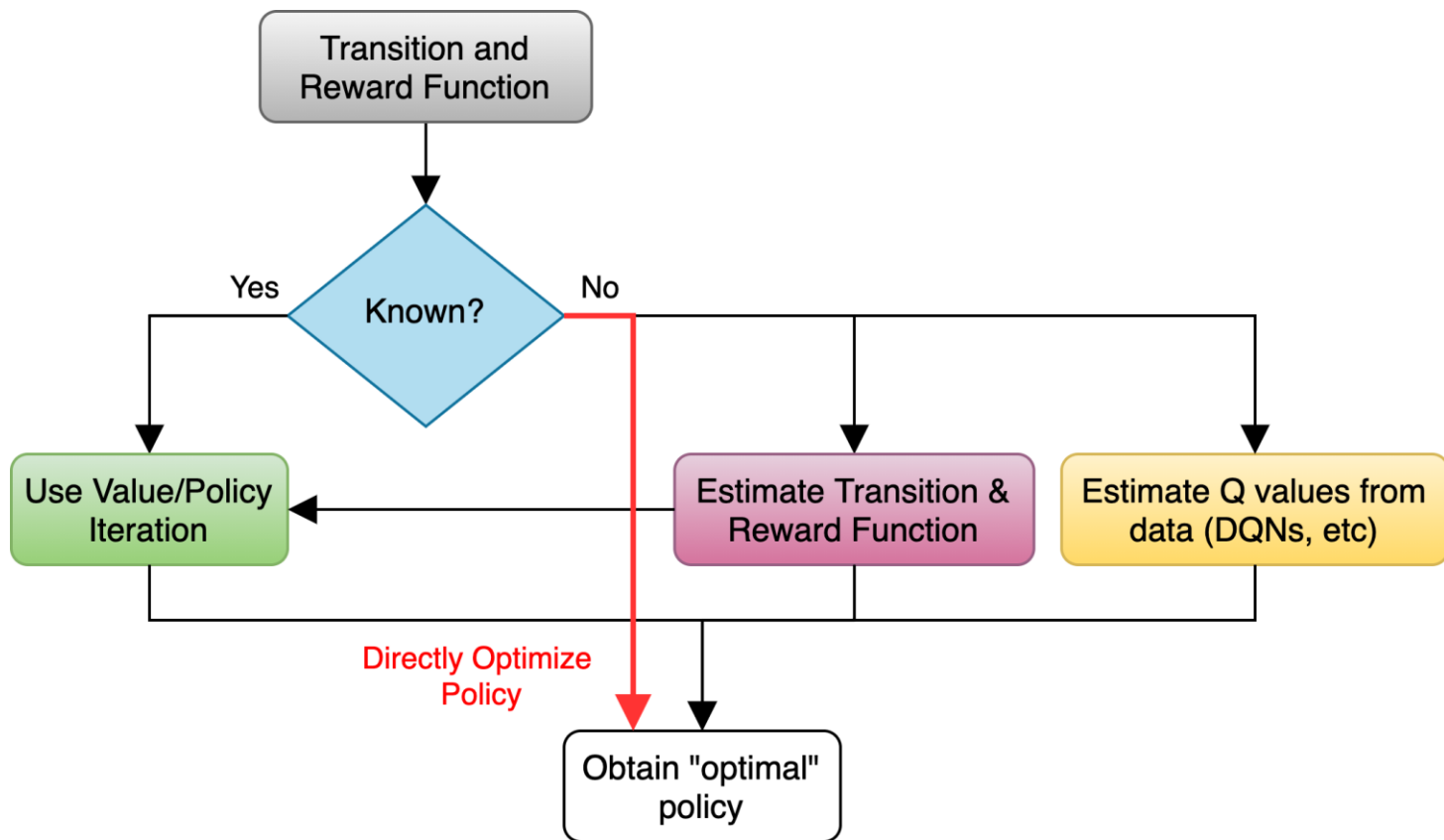
Topics:

- Reinforcement Learning Part 3
  - **Policy gradient**
  - **Actor-Critic**
  - **Proximal Policy Optimization (PPO)**
  - **(Probably Wed.) RL from Human Feedback for ChatGPT**

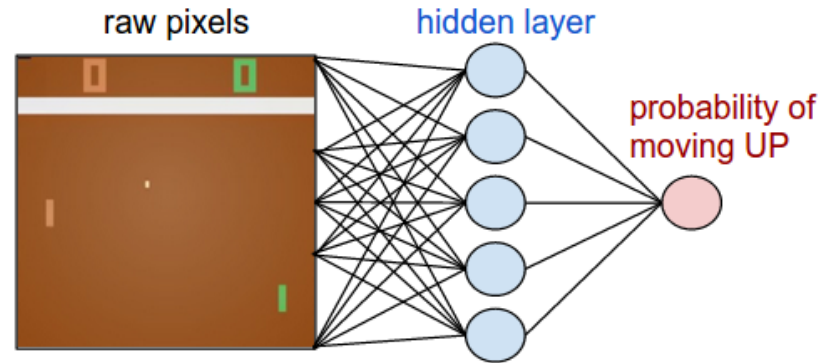
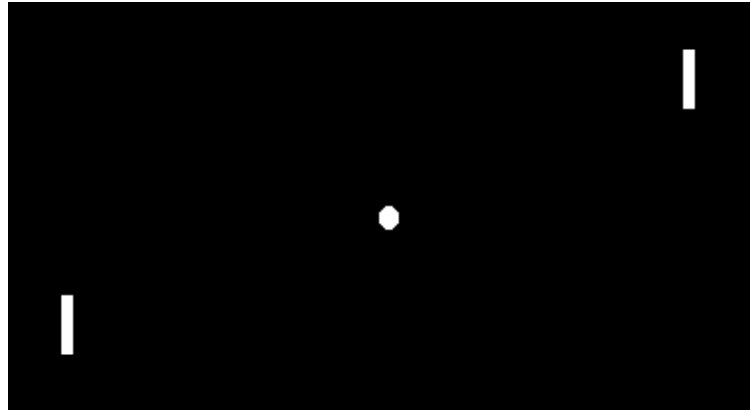
**CS 4644-DL / 7643-A**  
**ZSOLT KIRA**

## Admin

- Projects!
- CLOS is coming!



## Overview



## Pong from Pixels

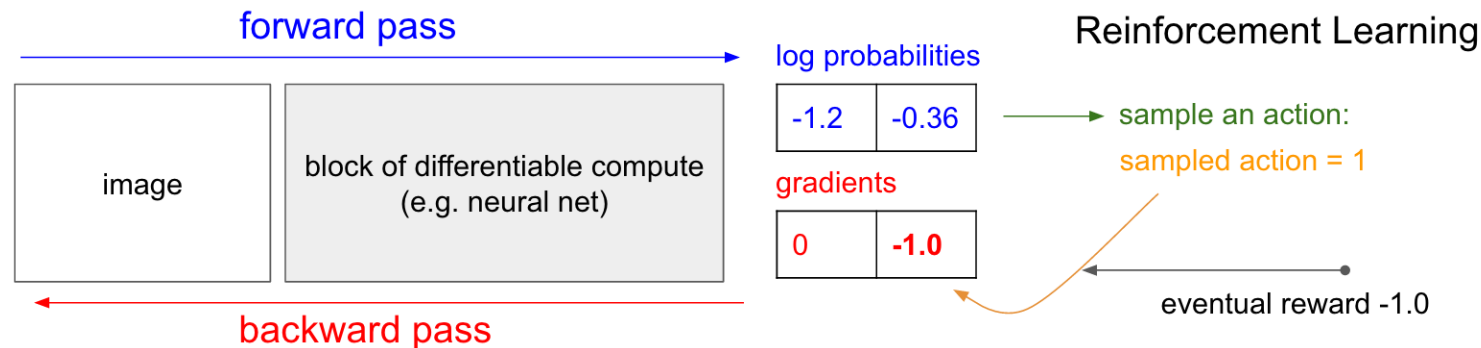
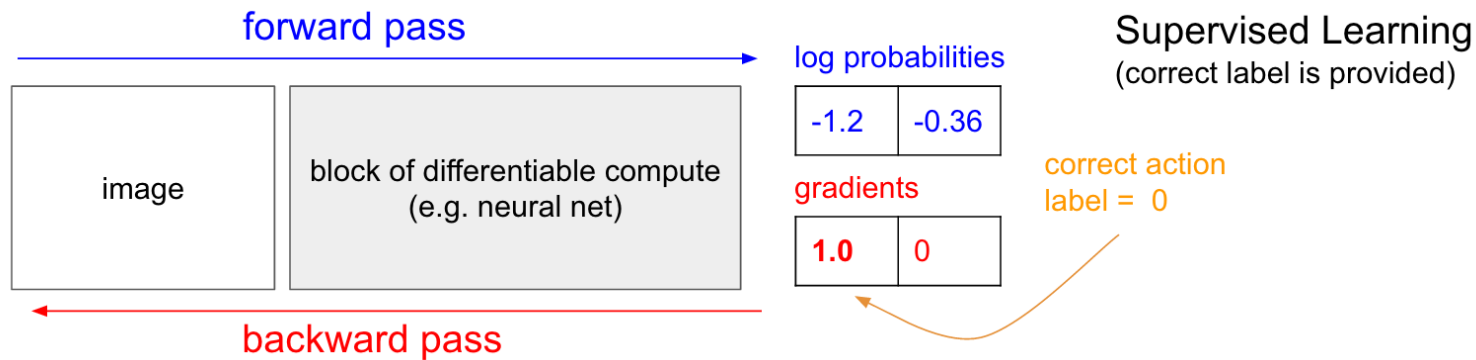


Image Source: <http://karpathy.github.io/2016/05/31/rl/>

## Policy Gradient: Loss Function

- Slightly re-writing the notation

Let  $\tau = (s_0, a_0, \dots s_T, a_T)$  denote a trajectory

$$\begin{aligned}\pi_{\theta}(\tau) &= p_{\theta}(\tau) = p_{\theta}(s_0, a_0, \dots s_T, a_T) \\ &= p(s_0) \prod_{t=0}^T p_{\theta}(a_t \mid s_t) \cdot p(s_{t+1} \mid s_t, a_t)\end{aligned}$$

$$\arg \max_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)]$$

$$\begin{aligned}
 J(\theta) &= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)] \\
 &= \mathbb{E}_{a_t \sim \pi(\cdot | s_t), s_{t+1} \sim p(\cdot | s_t, a_t)} \left[ \sum_{t=0}^T \mathcal{R}(s_t, a_t) \right]
 \end{aligned}$$

● How to gather data?

● We already have a policy:  $\pi_{\theta}$

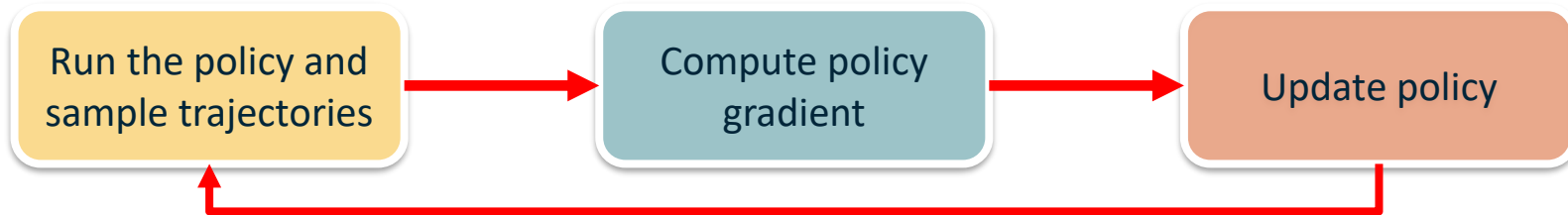
● Sample  $N$  trajectories  $\{\tau_i\}_{i=1}^N$  by acting according to  $\pi_{\theta}$

$$\approx \frac{1}{N} \sum_{i=1}^N \sum_{t=1}^T r(s_t^i, a_t^i)$$

- Sample trajectories  $\tau_i = \{s_1, a_1, \dots, s_T, a_T\}_i$  by acting according to  $\pi_\theta$
- Compute policy gradient as

$$\nabla_\theta J(\theta) \approx ?$$

- Update policy parameters:  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$



Slide credit: Sergey Levine



$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\mathcal{R}(\tau)]$$

$$= \nabla_{\theta} \int \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau$$

Expectation as integral

$$= \int \nabla_{\theta} \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau$$

Exchange integral and gradient

$$= \int \nabla_{\theta} \pi_{\theta}(\tau) \cdot \frac{\pi_{\theta}(\tau)}{\pi_{\theta}(\tau)} \cdot \mathcal{R}(\tau) d\tau$$

$$= \int \pi_{\theta}(\tau) \nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau) d\tau$$

$$\nabla_{\theta} \log \pi(\tau) = \frac{\nabla_{\theta} \pi(\tau)}{\pi(\tau)}$$

$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\nabla_{\theta} \log \pi_{\theta}(\tau) \mathcal{R}(\tau)]$$

$$\pi_{\theta}(\tau) = p(s_0) \prod_{t=0}^T p_{\theta}(a_t | s_t) \cdot p(s_{t+1} | s_t, a_t)$$

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} [\underbrace{\nabla_{\theta} \log \pi_{\theta}(\tau)}_{\text{Doesn't depend on Transition probabilities!}} \mathcal{R}(\tau)]$$

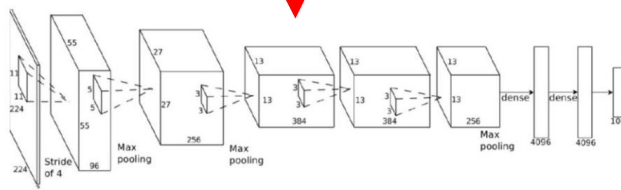
$$\nabla_{\theta} \left[ \cancel{\log p(s_0)} + \sum_{t=1}^T \log \pi_{\theta}(a_t | s_t) + \sum_{t=1}^T \cancel{\log p(s_{t+1} | s_t, a_t)} \right]$$

Doesn't depend on  
Transition probabilities!

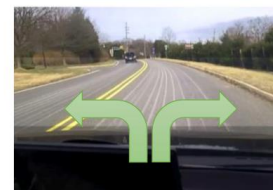
$$= \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \sum_{t=1}^T \mathcal{R}(s_t, a_t) \right]$$



$s_t$



$\pi_{\theta}(\mathbf{a}_t | s_t)$



$\mathbf{a}_t$

Continuous Action Space?

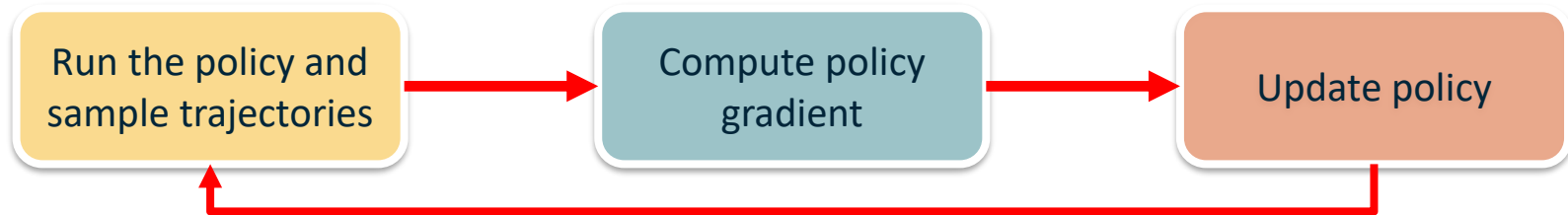
## Deriving The Policy Gradient

Sample trajectories  $\tau_i = \{s_1, a_1, \dots, s_T, a_T\}_i$  by acting according to  $\pi_\theta$

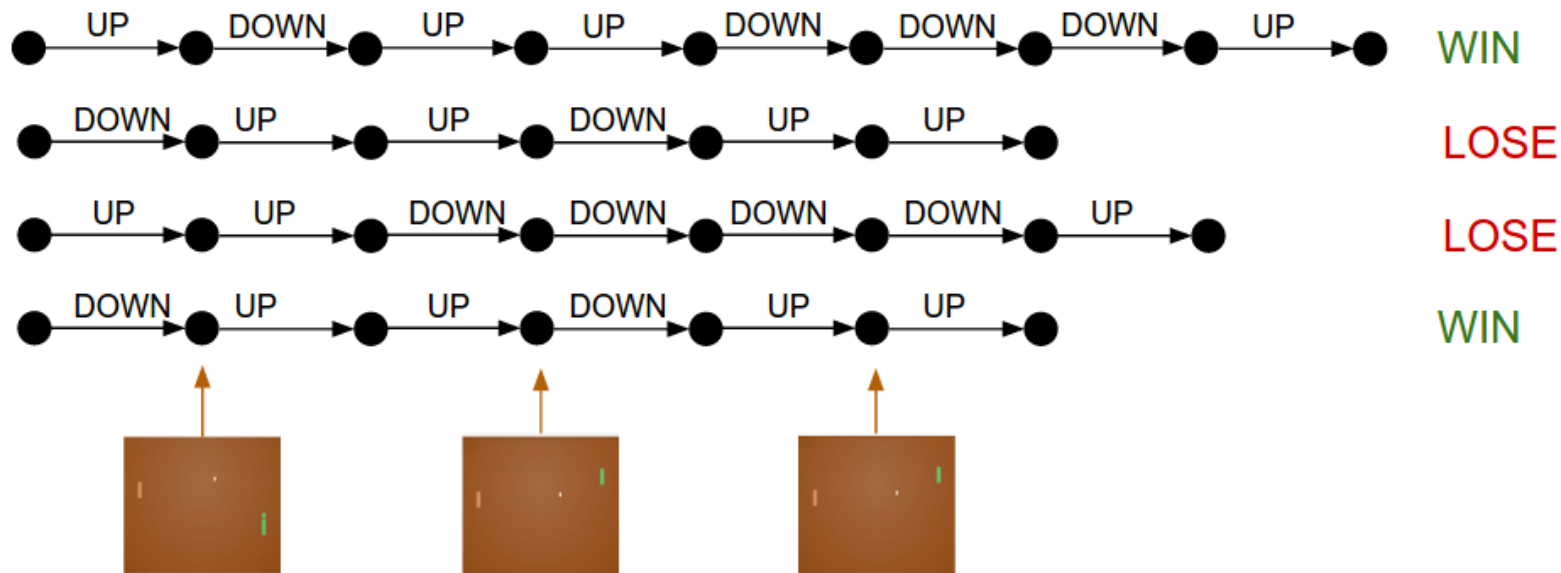
Compute policy gradient as

$$\nabla_\theta J(\theta) \approx \frac{1}{N} \sum_i \left[ \sum_{t=1}^T \nabla_\theta \log \pi_\theta (a_t^i | s_t^i) \cdot \sum_{t=1}^T \mathcal{R}(s_t^i | a_t^i) \right]$$

Update policy parameters:  $\theta \leftarrow \theta + \alpha \nabla_\theta J(\theta)$



Slide credit: Sergey Levine



Slide credit: Dhruv Batra

## Drawbacks of Policy Gradients

# Issues with Policy Gradients

- Credit assignment is hard!
  - Which specific action led to increase in reward
  - Suffers from high variance → leading to unstable training
- **Next time:** How to fix these issues

# Variance reduction

Gradient estimator:  $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

**First idea:** Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

# Variance reduction

Gradient estimator:  $\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} r(\tau) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$

**First idea:** Push up probabilities of an action seen, only by the cumulative future reward from that state

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

**Second idea:** Use discount factor  $\gamma$  to ignore delayed effects

$$\nabla_{\theta} J(\theta) \approx \sum_{t \geq 0} \left( \sum_{t' \geq t} \gamma^{t'-t} r_{t'} \right) \nabla_{\theta} \log \pi_{\theta}(a_t | s_t)$$

- Credit assignment is hard!
  - Which specific action led to increase in reward
  - Suffers from **high variance**, leading to unstable training
- How to reduce the variance?
  - Subtract an action independent baseline from the reward

$$\nabla_{\theta} J(\theta) = \mathbb{E}_{\tau \sim p_{\theta}(\tau)} \left[ \sum_{t=1}^T \nabla_{\theta} \log \pi_{\theta}(a_t | s_t) \cdot \sum_{t=1}^T (\mathcal{R}(s_t, a_t) - b(s_t)) \right]$$

- Why does it work? **Normalization constant (expected value doesn't change)**
- What is the best choice of b?



# How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

# How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

# Actor-Critic

- Learn both policy and Q function
  - Use the “actor” to sample trajectories
  - Use the Q function to “evaluate” or “critic” the policy

# Actor-Critic

- Learn both policy and Q function
  - Use the “actor” to sample trajectories
  - Use the Q function to “evaluate” or “critic” the policy
- REINFORCE:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \mathcal{R}(s, a)]$
- Actor-critic:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a)]$

# Actor-Critic

- Learn both policy and Q function
  - Use the “actor” to sample trajectories
  - Use the Q function to “evaluate” or “critic” the policy
- REINFORCE:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \mathcal{R}(s, a)]$
- Actor-critic:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a)]$
- Q function is unknown too! Update using  $\mathcal{R}(s, a)$

# Actor-Critic

- Initialize  $s$ ,  $\theta$  (policy network) and  $\beta$  (Q network)

# Actor-Critic

- Initialize  $s$ ,  $\theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot|s)$

# Actor-Critic

- Initialize  $s, \theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot|s)$
- For each step:
  - Sample reward  $\mathcal{R}(s, a)$  and next state  $s' \sim p(s'|s, a)$



# Actor-Critic

- Initialize  $s, \theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot|s)$
- For each step:
  - Sample reward  $\mathcal{R}(s, a)$  and next state  $s' \sim p(s'|s, a)$
  - evaluate “actor” using “critic”  $Q_{\beta}(s, a)$

# Actor-Critic

- Initialize  $s, \theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot | s)$
- For each step:
  - Sample reward  $\mathcal{R}(s, a)$  and next state  $s' \sim p(s' | s, a)$
  - evaluate “actor” using “critic”  $Q_{\beta}(s, a)$  **and update policy:**

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\beta}(s, a)$$

- Initialize  $s$ ,  $\theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot | s)$
- For each step:
  - Sample reward  $\mathcal{R}(s, a)$  and next state  $s' \sim p(s' | s, a)$
  - evaluate “actor” using “critic”  $Q_{\beta}(s, a)$  and update policy:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\beta}(s, a)$$

- Update “critic”: MSE Loss  $:= \left( Q_{new}(s, a) - (r + \max_a Q_{old}(s', a)) \right)^2$ 
  - Recall Q-learning

- Initialize  $s$ ,  $\theta$  (policy network) and  $\beta$  (Q network)
- sample action  $a \sim \pi_{\theta}(\cdot|s)$
- For each step:
  - Sample reward  $\mathcal{R}(s, a)$  and next state  $s' \sim p(s'|s, a)$
  - evaluate “actor” using “critic”  $Q_{\beta}(s, a)$  and update policy:

$$\theta \leftarrow \theta + \alpha \nabla_{\theta} \log \pi_{\theta}(a | s) Q_{\beta}(s, a)$$

- Update “critic”:

- Recall Q-learning  $\text{MSE Loss} := \left( Q_{\text{new}}(s, a) - (r + \max_a Q_{\text{old}}(s', a)) \right)^2$

$$a \leftarrow a', s \leftarrow s'$$

- Update  $\beta$  Accordingly

# How to choose the baseline?

A better baseline: Want to push up the probability of an action from a state, if this action was better than the **expected value of what we should get from that state**.

Q: What does this remind you of?

A: Q-function and value function!

Intuitively, we are happy with an action  $a_t$  in a state  $s_t$  if  $Q^\pi(s_t, a_t) - V^\pi(s_t)$  is large. On the contrary, we are unhappy with an action if it's small.

Using this, we get the estimator:  $\nabla_\theta J(\theta) \approx \sum_{t \geq 0} (Q^{\pi_\theta}(s_t, a_t) - V^{\pi_\theta}(s_t)) \nabla_\theta \log \pi_\theta(a_t | s_t)$

# Actor-critic

- In general, replacing the policy evaluation or the “critic” leads to different flavors of the actor-critic
  - REINFORCE:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) \mathcal{R}(s, a)]$
  - Q – Actor Critic  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) Q^{\pi_{\theta}}(s, a)]$
  - Advantage Actor Critic:  $\nabla_{\theta} J(\pi_{\theta}) = \mathbb{E}_{a \sim \pi_{\theta}} [\nabla_{\theta} \log \pi_{\theta}(a|s) A^{\pi_{\theta}}(s, a)]$   
 $= Q^{\pi_{\theta}}(s, a) - V^{\pi_{\theta}}(s)$

# Advanced policy gradient methods

- **Trust Region Policy Gradient** (TRPO, Schulman 2017)
- Issue with vanilla actor critic: policy may receive huge update!
  - Big parameter update -> drastic change in behavior -> may stuck in low-reward region!
- Idea: Anchor policy updates to past!

$$J(\theta) = \mathbb{E}_{s \sim \rho^{\pi_{\theta_{\text{old}}}}, a \sim \pi_{\theta_{\text{old}}}} \left[ \frac{\pi_{\theta}(a|s)}{\pi_{\theta_{\text{old}}}(a|s)} \hat{A}_{\theta_{\text{old}}}(s, a) \right]$$

- Idea: constrain the update to a *trust region* using off-policy policy gradient
- Subject to:

$$\mathbb{E}_{s \sim \rho^{\pi_{\theta_{\text{old}}}}} [D_{\text{KL}}(\pi_{\theta_{\text{old}}}(\cdot | s) \parallel \pi_{\theta}(\cdot | s))] \leq \delta$$

importance sampling

$$\begin{aligned} E_{x \sim p(x)}[f(x)] &= \int p(x) f(x) dx \\ &= \int \frac{q(x)}{q(x)} p(x) f(x) dx \\ &= \int q(x) \frac{p(x)}{q(x)} f(x) dx \\ &= E_{x \sim q(x)} \left[ \frac{p(x)}{q(x)} f(x) \right] \end{aligned}$$

Optimizing this objective requires calculating Hessian (second-order optimization)!

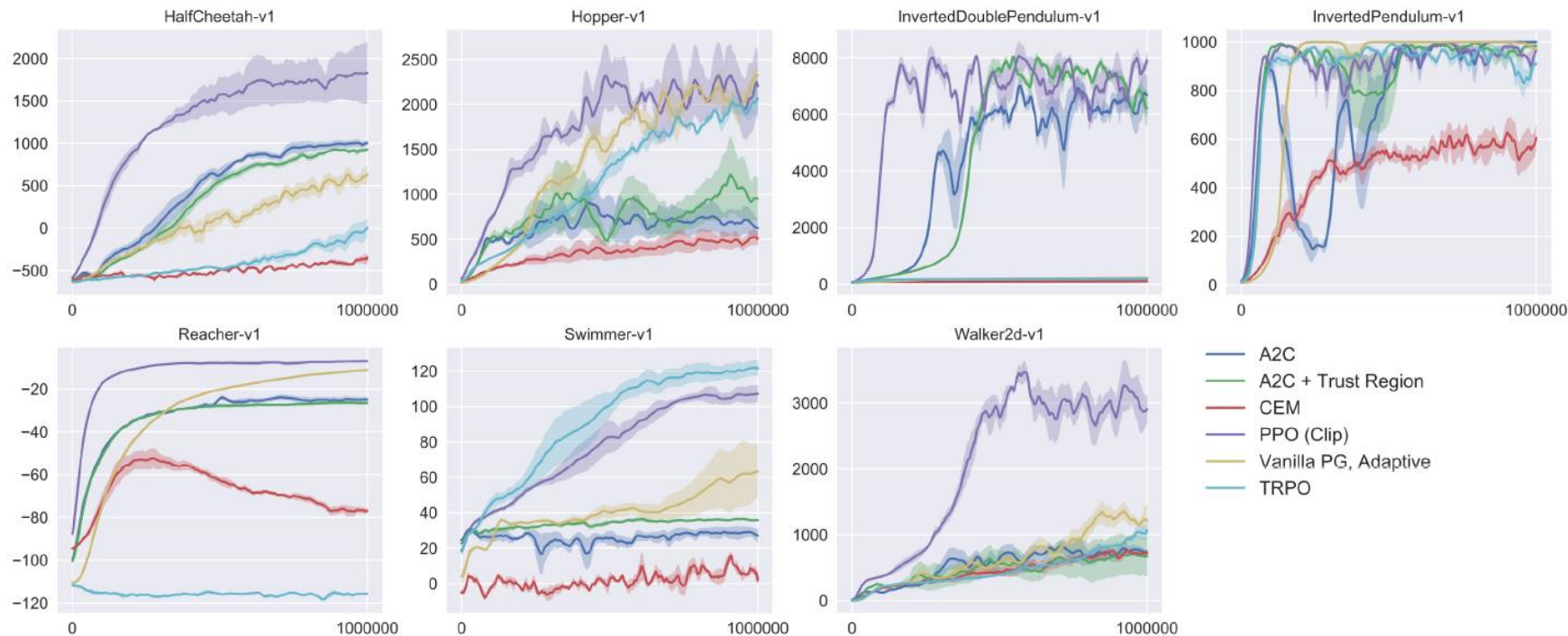
# Advanced policy gradient methods

- **Proximal Policy Optimization** (PPO, Schulman 2017)
- Issue with TRPO: objective too complicated! Requires second-order optimization (calculating Hessian).
- Idea: Approximate trust-region constraint with a penalty term

$$\underset{\theta}{\text{maximize}} \quad \hat{\mathbb{E}}_t \left[ \frac{\pi_{\theta}(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} \hat{A}_t \right] - \beta \hat{\mathbb{E}}_t [\text{KL}[\pi_{\theta_{\text{old}}}(\cdot | s_t), \pi_{\theta}(\cdot | s_t)]]$$



# Advanced policy gradient methods



Schulman 2017

# Applications to Language Modeling

- One of the key benefits of RL is ability to tune to **evaluative feedback**
- Used to turn *unaligned* language models to ones that do what we want: Answer our questions
- Reinforcement Learning from Human Feedback (RLHF)

# From Language Models to Assistants

1. Zero-Shot (ZS) and Few-Shot (FS) In-Context Learning
2. Instruction finetuning
3. Reinforcement Learning from Human Feedback (RLHF)

4. What's next?

# RL to Improve Language modeling $\neq$ 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.

Language models are not *aligned* with user intent [[Ouyang et al., 2022](#)].

# Language modeling $\neq$ assisting users

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

**COMPLETION**

**Human**

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.

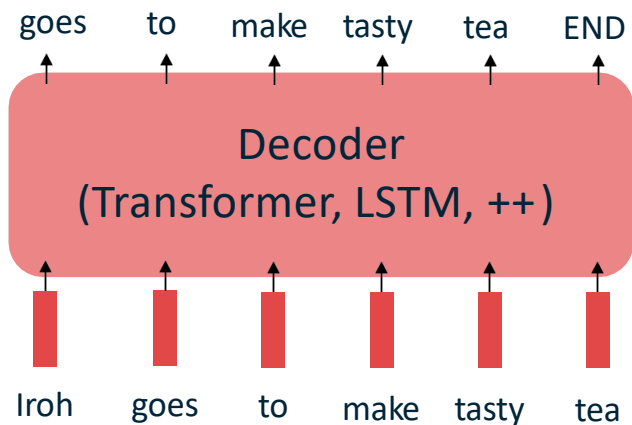
Language models are not *aligned* with user intent [[Ouyang et al., 2022](#)].

Finetuning to the rescue!

Pretraining can improve NLP applications by serving as parameter initialization.

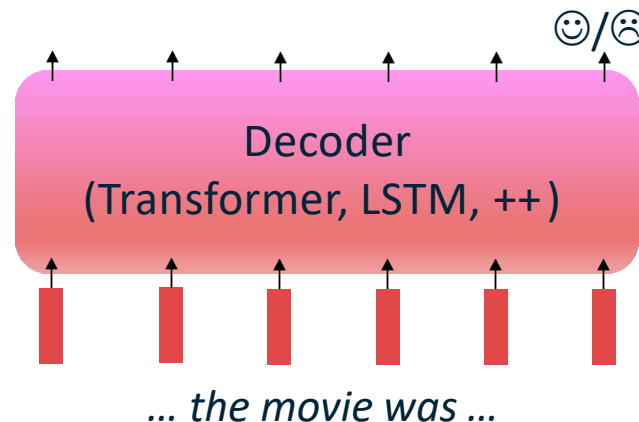
### Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



### Step 2: Finetune (on your task)

Not many labels; adapt to the task!

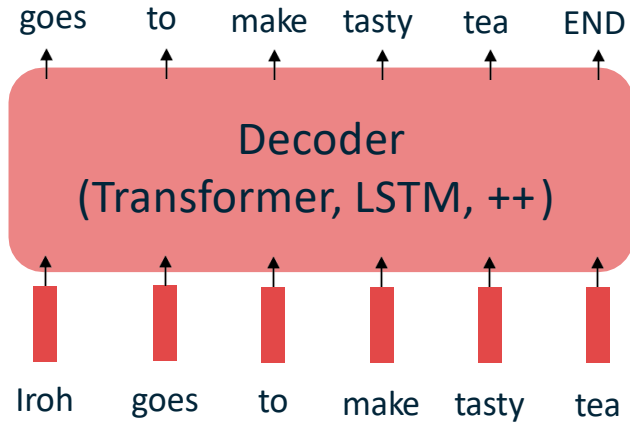


# Scaling up finetuning

Pretraining can improve NLP applications by serving as parameter initialization.

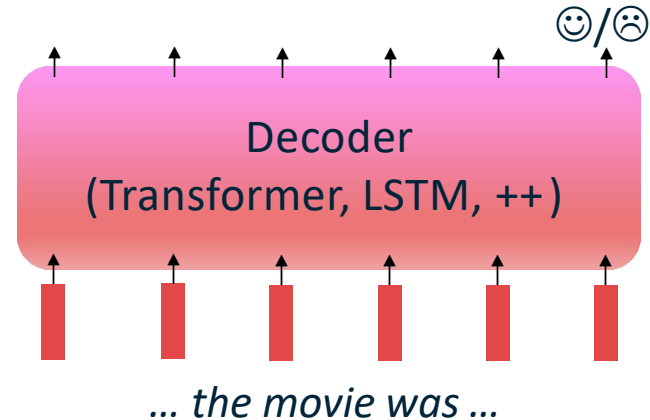
## Step 1: Pretrain (on language modeling)

Lots of text; learn general things!



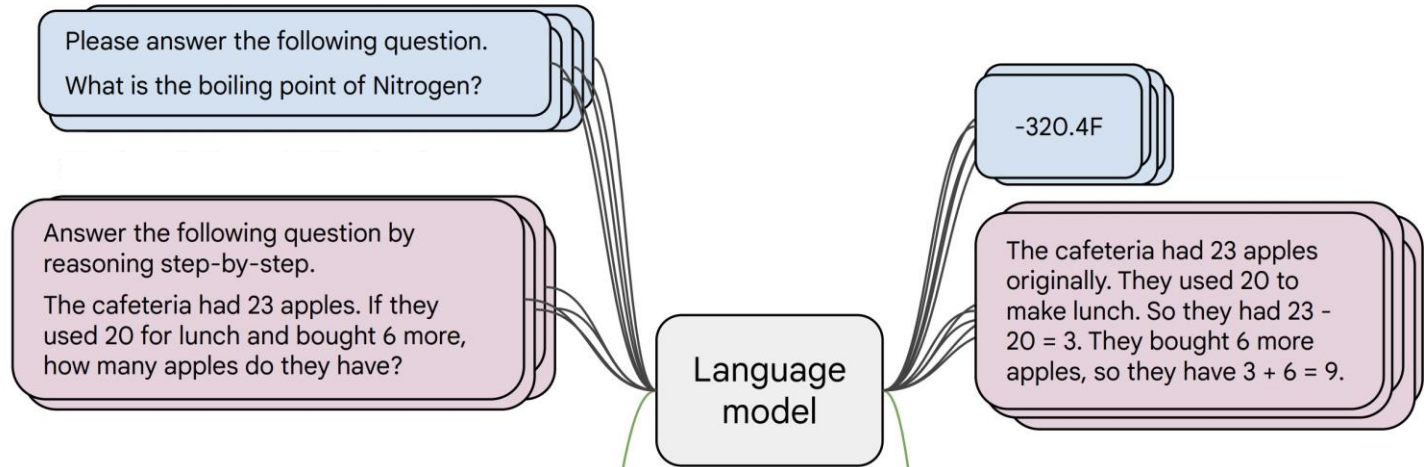
## Step 2: Finetune (on **many tasks**)

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



# Instruction finetuning

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



- Evaluate on **unseen tasks**



[FLAN-T5; [Chung et al., 2022](#)]



# Instruction finetuning

## 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)

Highly recommend trying FLAN-T5 out to get a sense of its capabilities:

<https://huggingface.co/google/flan-t5-xxl>

# Instruction finetuning

## 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

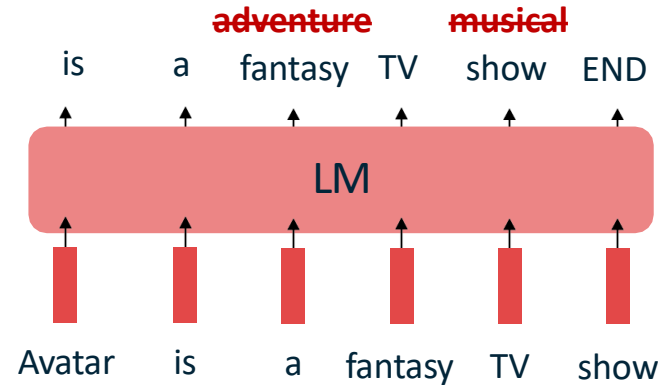
The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✓

Highly recommend trying FLAN-T5 out to get a sense of its capabilities:

<https://huggingface.co/google/flan-t5-xxl>

# Limitations of instruction finetuning?

- One limitation of instruction finetuning is obvious: it's **expensive** to collect ground-truth data for tasks.
- But there are other, subtler limitations too. Can you think of any?
- **Problem 1:** tasks like open-ended creative generation have no right answer.
  - *Write me a story about a dog and her pet grasshopper.*
- **Problem 2:** language modeling penalizes all token-level mistakes equally, but some errors are worse than others.
- Even with instruction finetuning, there is a mismatch between the LM objective and the objective of “satisfy human preferences”!
- Can we **explicitly attempt to satisfy human preferences**?



# From Language Models to Assistants

## 1. Zero-Shot (ZS) and Few-Shot (FS) In-Context Learning

- + No finetuning needed, prompt engineering (e.g. CoT) can improve performance
- Limits to what you can fit in context
- Complex tasks will probably need gradient steps

## 2. Instruction finetuning

- + Simple and straightforward, generalize to unseen tasks
- Collecting demonstrations for so many tasks is expensive
- Mismatch between LM objective and human preferences

## 3. Reinforcement Learning from Human Feedback (RLHF)

## 4. What's next?

# From Language Models to Assistants

## 1. Zero-Shot (ZS) and Few-Shot (FS) In-Context Learning

- + No finetuning needed, prompt engineering (e.g. CoT) can improve performance
- Limits to what you can fit in context
- Complex tasks will probably need gradient steps

## 2. Instruction finetuning

- + Simple and straightforward, generalize to unseen tasks
- Collecting demonstrations for so many tasks is expensive
- Mismatch between LM objective and human preferences

## 3. Reinforcement Learning from Human Feedback (RLHF)

## 4. What's next?

# Optimizing for human preferences

- Let's say we were training a language model on some task (e.g. summarization).
- For each LM sample  $s$ , imagine we had a way to obtain a *human reward* of that summary:  $R(s) \in \mathbb{R}$ , higher is better.

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.

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

- Now we want to maximize the expected reward of samples from our LM:

[   ]

Note: for mathematical simplicity  
we're assuming only one "prompt"

# How do we model human preferences?

- Awesome: now for any **arbitrary, non-differentiable reward function**  $R(s)$ , we can train our language model to maximize expected reward.
- Not so fast! (Why not?)
- **Problem 1:** human-in-the-loop is expensive!
- **Solution:** instead of directly asking humans for preferences, **model their preferences** as a separate (NLP) problem! [Knox and Stone, 2009]

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

$$R(s_1) = 8.0$$



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

$$R(s_2) = 1.2$$



Train an LM  $RM_\phi(s)$  to  
predict human  
preferences from an  
annotated dataset, then  
optimize for  $RM_\phi$  instead.

# How do we model human preferences?

- **Problem 2:** human judgments are noisy and miscalibrated!
- **Solution:** instead of asking for direct ratings, ask for **pairwise comparisons**, which can be more reliable [[Phelps et al., 2015](#); [Clark et al., 2018](#)]

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

$$R(s_3) = \begin{matrix} s_3 \\ 4.1? & 6.6? & 3.2? \\ ? \end{matrix}$$



# How do we model human preferences?

- **Problem 2:** human judgments are noisy and miscalibrated!
- **Solution:** instead of asking for direct ratings, ask for **pairwise comparisons**, which can be more reliable [[Phelps et al., 2015](#); [Clark et al., 2018](#)]

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

>

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

>

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

$S_1$

1.2

$S_3$

$S_2$

Bradley-Terry [1952] paired comparison model

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

“winning”

sample

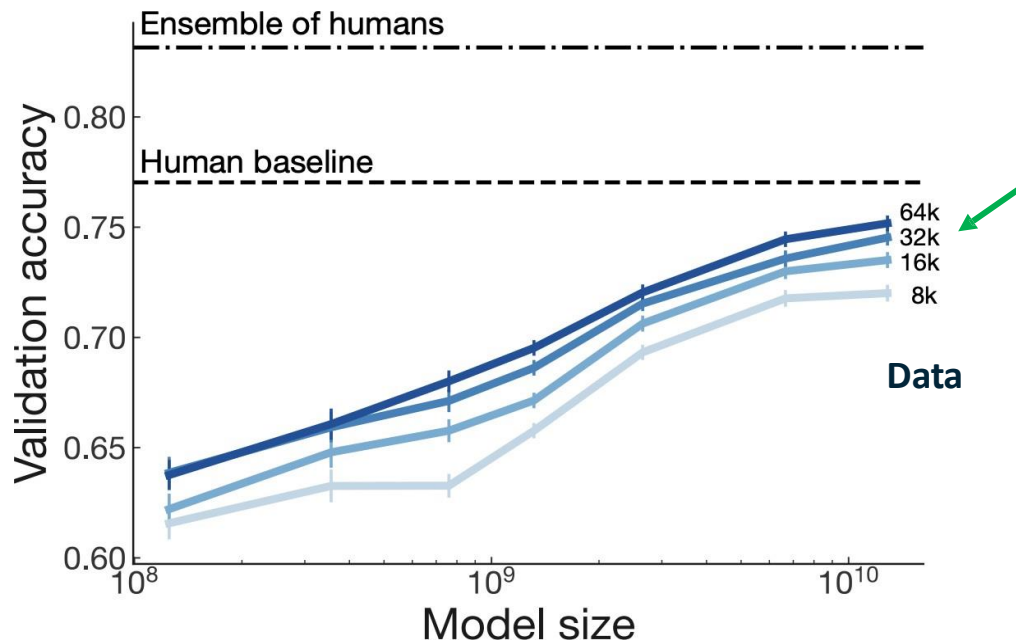
“losing”

sample

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

# Make sure your reward model works first!

Evaluate RM on predicting outcome of held-out human judgments

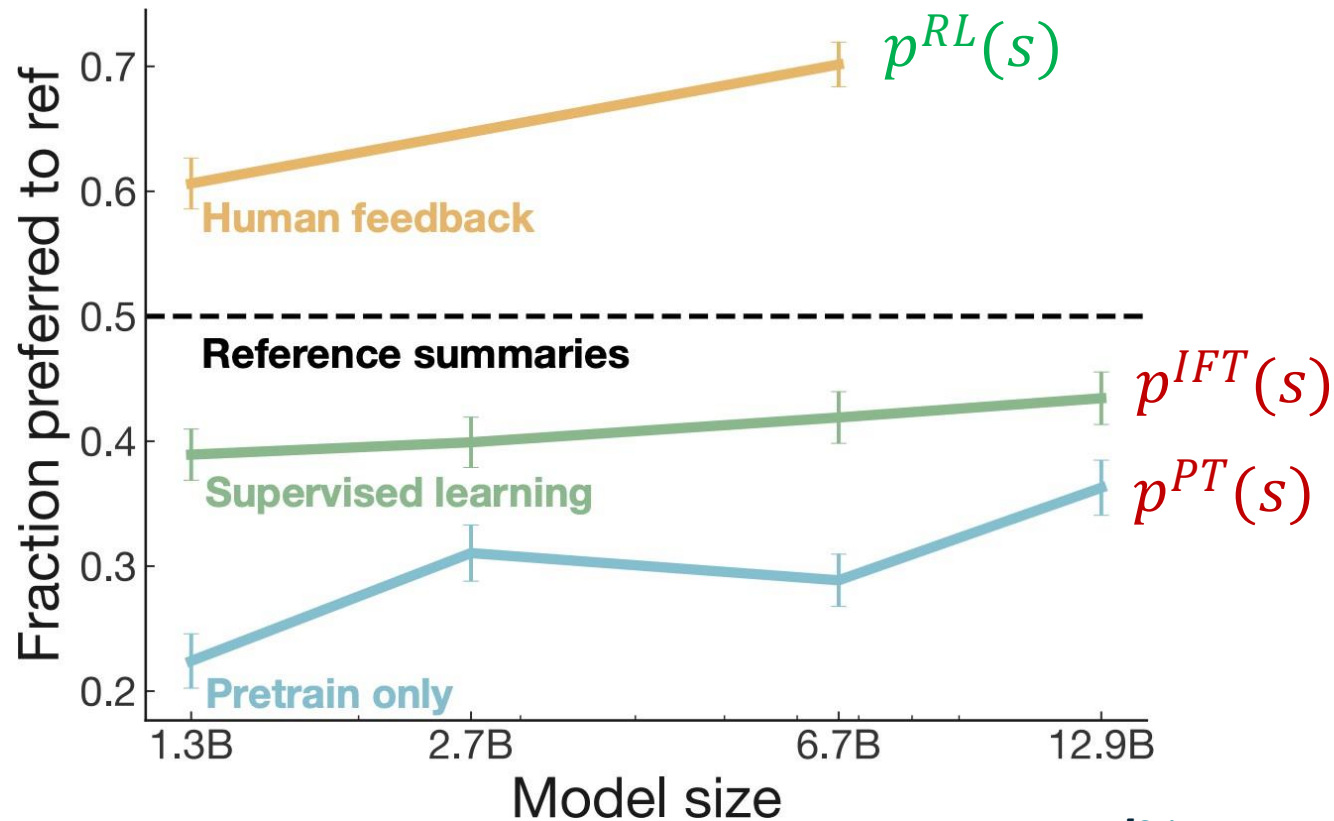


Large enough RM  
trained on enough  
data approaching  
single human perf

# RLHF: Putting it all together [[Christiano et al., 2017](#); [Stiennon et al., 2020](#)]

- 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 reward with RL

# RLHF provides gains over pretraining + finetuning



[Stiennon et al., 2020]

# InstructGPT: scaling up RLHF to tens of thousands of tasks

30k  
tasks!

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

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.

Once upon a time...

The reward model  
calculates a  
reward for  
the output.

RM

The reward is  
used to update  
the policy  
using PPO.

[Ouyang et al., 2022]

# InstructGPT: scaling up RLHF to tens of thousands of tasks

## Tasks collected from labelers:

- **Plain:** We simply ask the labelers to come up with an arbitrary task, while ensuring the tasks had sufficient diversity.
- **Few-shot:** We ask the labelers to come up with an instruction, and multiple query/response pairs for that instruction.
- **User-based:** We had a number of use-cases stated in waitlist applications to the OpenAI API. We asked labelers to come up with prompts corresponding to these use cases.

| Use-case      | Prompt  |
|---------------|---|
| Brainstorming | List five ideas for how to regain enthusiasm for my career  |
| Generation    | Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home. |

54

# InstructGPT

**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.

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.

# InstructGPT

**PROMPT** *Write a short poem about a wise frog.*

**COMPLETION** GPT-3

Write a short story in which a character has two different names.

Write a short story in which you try to get something back that you have lost.

Write a short story in which a character has a bad dream.

InstructGPT

The frog is a symbol of wisdom  
He knows all the secrets of the world  
He is a master of disguise  
And a great teacher of life  
He is a symbol of transformation  
And the bringer of change  
He is the frog who has seen it all  
And knows the meaning of it all



# ChatGPT: Instruction Finetuning + 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.

**(Instruction finetuning!)**

# ChatGPT: Instruction Finetuning + 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

To create a reward model for reinforcement learning, we needed to collect comparison data, which consisted of two or more model responses ranked by quality. To collect this data, we took conversations that AI trainers had with the chatbot. We randomly selected a model-written message, sampled several alternative completions, and had AI trainers rank them. Using these reward models, we can fine-tune the model using Proximal Policy Optimization. We performed several iterations of this process.

**(RLHF!)**

# Summary

- Reinforcement Learning has wide range of uses!
  - Can optimize any non-differentiable “objective”
  - Can be combined with supervised/unsupervised learning
  - Can be combined with other algorithms such as tree search