

RLHF: Reinforcement Learning based on Human Feedback

Lecture 10
Subir Varma

Contents

- ▶ RLHF: Reinforcement Learning based on Human Feedback
 - ▶ PPO: Proximal Policy Optimization
- 

RLHF

RL based on Human Feedback

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 Askell¹ Peter Welinder Paul Christiano^{*1}

Jan Leike* Ryan Lowe*

OpenAI

Abstract

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not *aligned* with their users. In this paper, we show an avenue for aligning language models with user intent on a wide range of tasks by fine-tuning with human feedback. Starting with a set of labeler-written prompts and prompts submitted through the OpenAI API, we collect a dataset of labeler demonstrations of the desired model behavior, which we use to fine-tune GPT-3 using supervised learning. We then collect a dataset of rankings of model outputs, which we use to further fine-tune this supervised model using reinforcement learning from human feedback. We call the resulting models *InstructGPT*. In human evaluations on our prompt distribution, outputs from the 1.3B parameter *InstructGPT* model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters. Moreover, *InstructGPT* models show improvements in truthfulness and reductions in toxic output generation while having minimal performance regressions on public NLP datasets. Even though *InstructGPT* still makes simple mistakes, our results show that fine-tuning with human feedback is a promising direction for aligning language models with human intent.

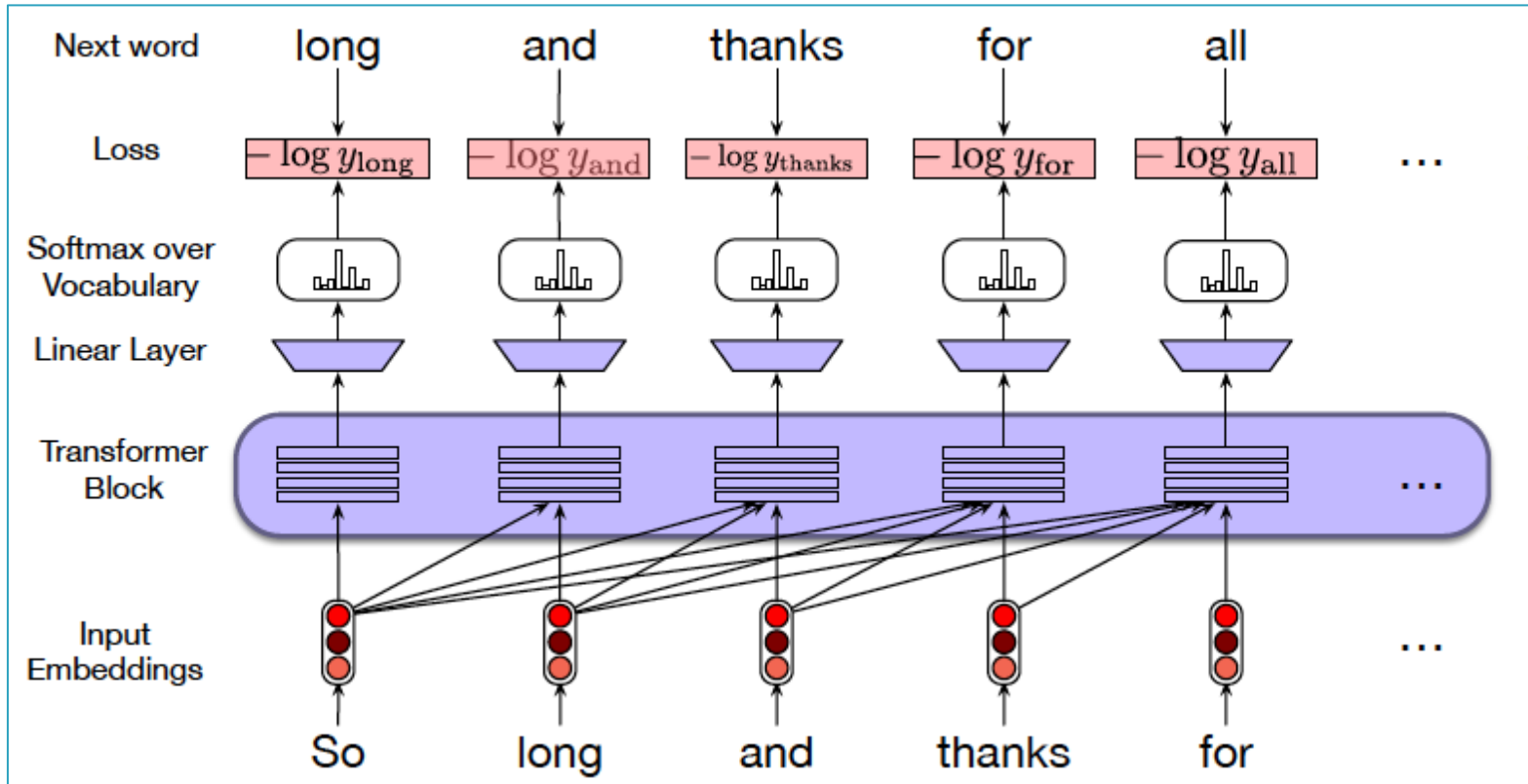
1 Introduction

Large language models (LMs) can be “prompted” to perform a range of natural language processing (NLP) tasks, given some examples of the task as input. However, these models often express unintended behaviors such as making up facts, generating biased or toxic text, or simply not following user instructions (Bender et al., 2021; Bommasani et al., 2021; Kenton et al., 2021; Weidinger et al., 2021; Tamkin et al., 2021; Gehman et al., 2020). This is because the language modeling objective

*Primary authors. This was a joint project of the OpenAI Alignment team. RL and JL are the team leads. Corresponding author: love@openai.com.

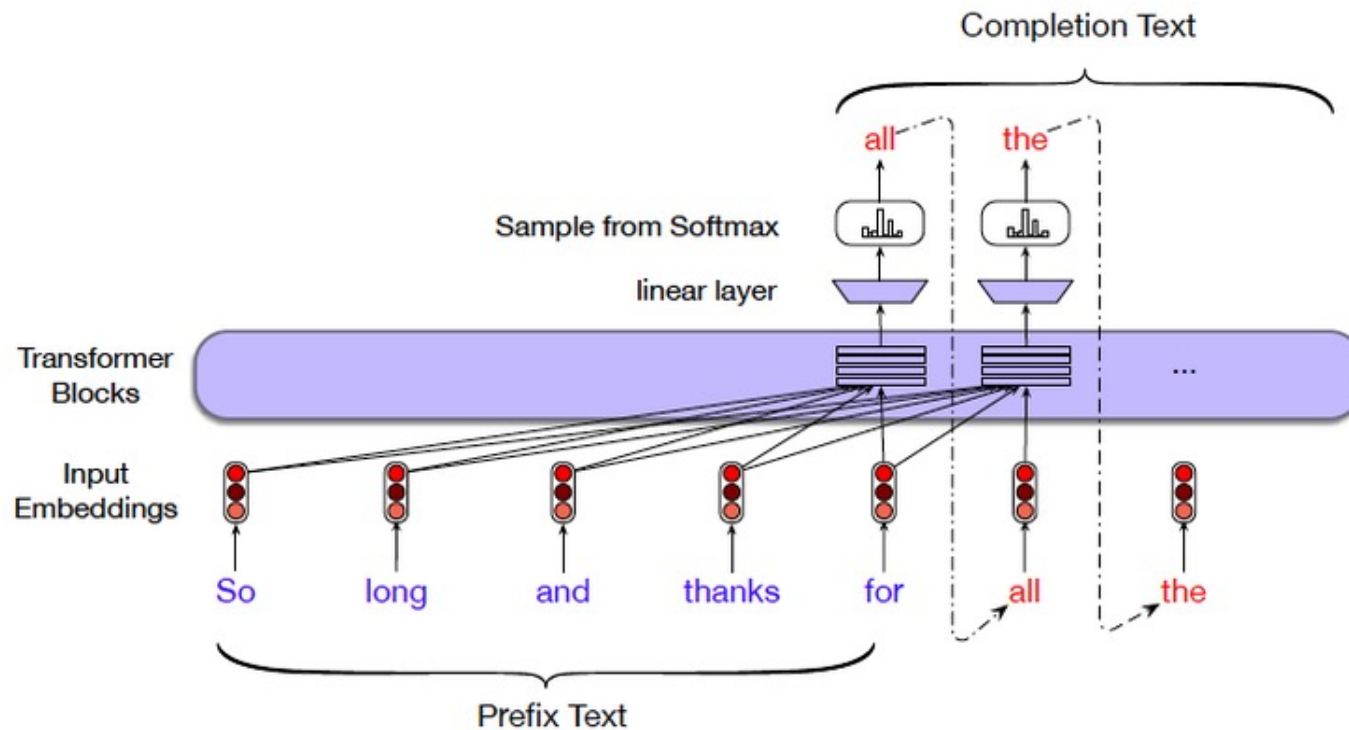
¹Work done while at OpenAI. Current affiliations: AA: Anthropic; PC: Alignment Research Center.

Large Language Models (LLMs) using Transformers: Training



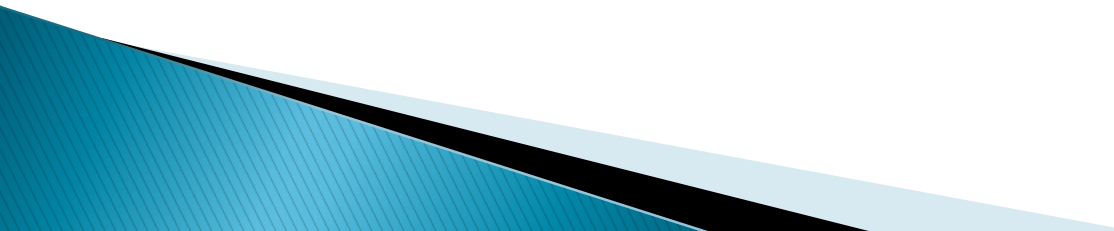
- Trained to predict the next word in a sentence using the Cross entropy Loss function
- Massive amount of text used for Training
- Some LLMs have hundreds of billions of parameters

Inference: Text Completion



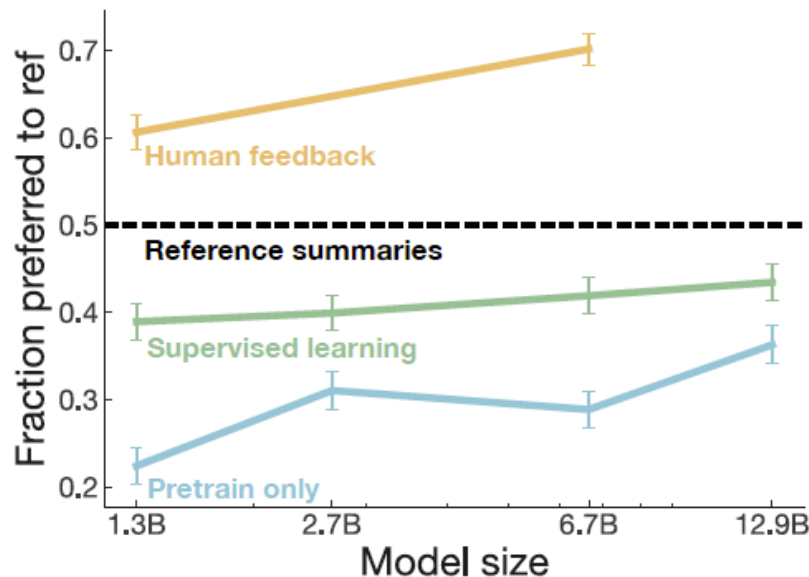
- Sentences are generated one word at a time by sampling from the output probability distribution
- It's thought that LLMs create some sort of model of the world from the training data, which they use during prediction

Problems with LLMs (Pre Chat GPT)

- ▶ Mis-alignment between LLM generated text and generating high-quality text as determined by humans
 - ▶ Reason for mis-alignment: LLM objective which is to predict the next word is different from the objective “follow the users instructions helpfully and safely”
 - ▶ Another reason: The Cross Entropy Loss function does not differentiate between important errors (like making up facts) and un-important errors (selecting a precise word from a set of synonyms)
- 

RLHF Objective

- ▶ Solution using RLHF: Fine tune models from Human Feedback using Reward Learning
- ▶ Make the LLM act according to the user's intention, both explicit instructions such as following instructions and implicit instructions such as staying truthful, and not being biased, toxic or otherwise harmful.



Make the LLM

- Helpful: Help the user
- Honest: No hallucinations
- Harmless: Don't cause harm

Step 0: Choose an LM

As a starting point RLHF use a language model that has already been pretrained with the classical pretraining objectives (see this [blog post](#) for more details).

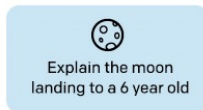
- ▶ OpenAI used a smaller version of GPT-3 for its first popular RLHF model, [InstructGPT](#).
- ▶ Anthropic used transformer models from 10 million to 52 billion parameters trained for this task.
- ▶ DeepMind used their 280 billion parameter model [Gopher](#).

3-Step Procedure

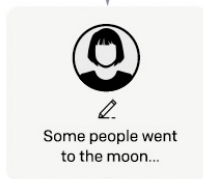
Step 1

Collect demonstration data, and train a supervised policy.

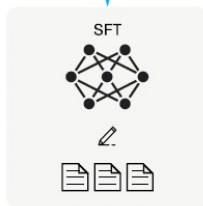
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



This data is used to fine-tune GPT-3 with supervised learning.

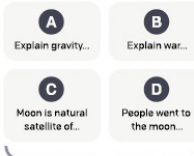
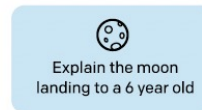


SFT: Supervised Fine Tuning

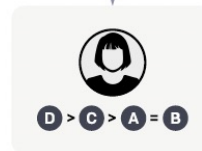
Step 2

Collect comparison data, and train a reward model.

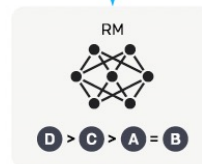
A prompt and several model outputs are sampled.



A labeler ranks the outputs from best to worst.



This data is used to train our reward model.



Step 3

Optimize a policy against the reward model using reinforcement learning.

A new prompt is sampled from the dataset.



The policy generates an output.

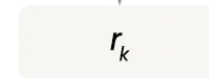


Once upon a time...

The reward model calculates a reward for the output.

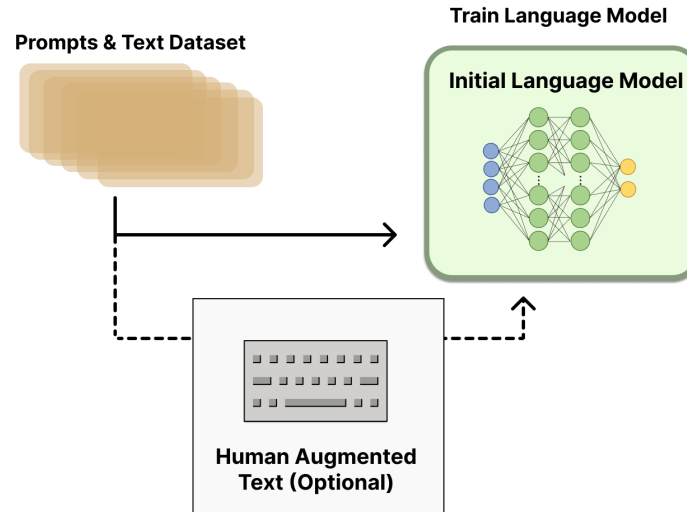


The reward is used to update the policy using PPO.



RLHF

Step 1: Fine Tuning



- ▶ This initial model can be fine-tuned on additional human generated text.
- ▶ This is done using Supervised Learning (with Cross Entropy Loss Function)

RLHF

- ▶ RLHF uses methods from Reinforcement Learning to directly optimize a language model with human feedback.
- ▶ RLHF has enabled language models to begin to align a model trained on a general corpus of text data to that of complex human values.



Can you explain to me how RL with human feedback works for language models like I am five?



Sure! So, imagine you have a robot named Rufus who wants to learn how to talk like a human. Rufus has a language model that helps him understand words and sentences.

First, Rufus will say something using his language model. For example, he might say "I am a robot."

Then, a human will listen to what Rufus said and give him feedback on whether it sounded like a natural sentence a human would say. The human might say, "That's not quite right, Rufus. Humans don't usually say 'I am a robot.' They might say 'I'm a robot' or 'I am a machine.'"

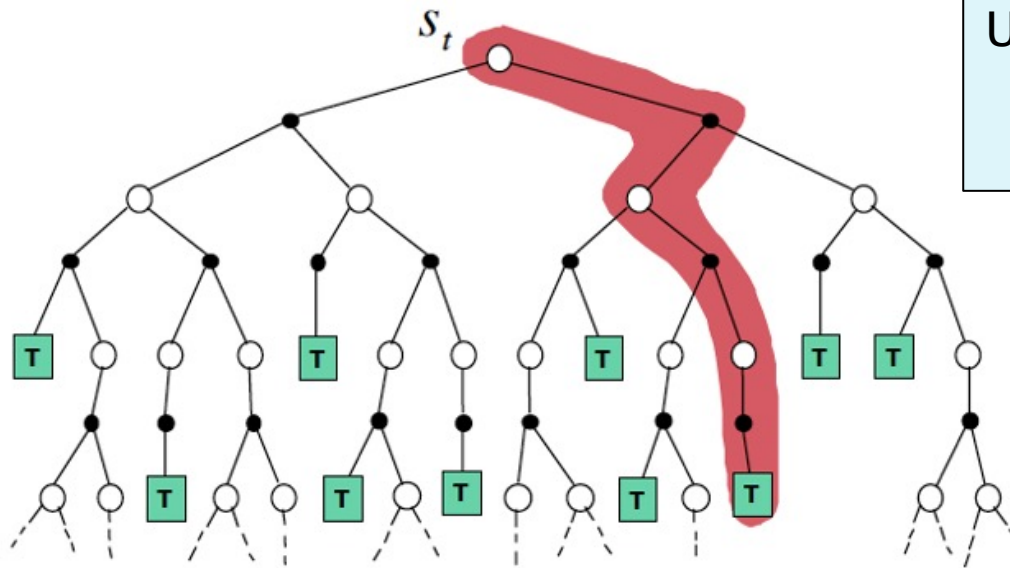
Rufus will take this feedback and use it to update his language model. He will try to say the sentence again, using the new information he received from the human. This time, he might say "I'm a robot."

The human will listen again and give Rufus more feedback. This process will continue until Rufus can say sentences that sound natural to a human.

Over time, Rufus will learn how to talk like a human thanks to the feedback he receives from humans. This is how language models can be improved using RL with human feedback.

ChatGPT's explanation
For RLHF!

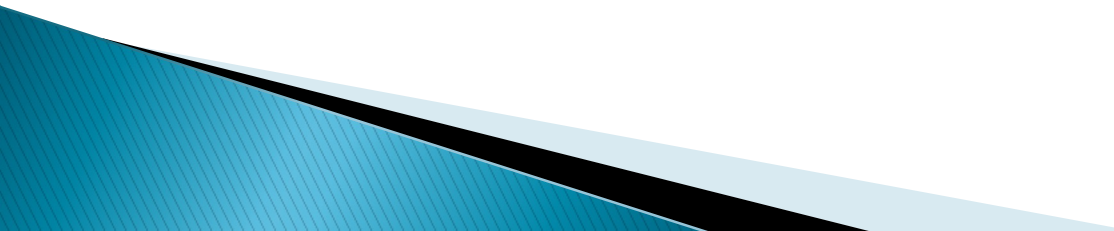
Basic Idea behind RLHF Algorithm: Use Policy Gradients on Generated Text



- Treat the LLM as a Policy Function $\pi(S)$ with words as Actions
- An episode corresponds to the text being generated by the LLM
- State S corresponds to the (incomplete) text that has been generated so far
- Action A correspond to the choice of the next word to be added
- Reward R is for the entire text at the end of the episode
- After each episode, modify the Policy (i.e. the LLM) so that in the next iteration it generates text with a higher reward

Where did the reward come from?

RLHF Steps

- ▶ **Step 1: Collect samples from existing policies and send comparisons to humans:** For each Reddit post, sample summaries from several sources including the current policy, initial policy, original reference summaries and various baselines. Send a batch of pairs of summaries to human evaluators, who are tasked with selecting the best summary of a given Reddit post.
 - ▶ **Step 2: Learn a reward model from human comparisons:** Given a post and a candidate summary, train a reward model to predict the log odds that this summary is the better one, as judged by human labelers.
 - ▶ **Step 3: Optimize a policy against the reward model:** Treat the logit output of the reward model as a reward that is used to optimize using Reinforcement Learning, specifically with the PPO algorithm.
- 

RLHF Steps

1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



"j is better than k"

2 Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward r for each summary.



r_j

r_k

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$\text{loss} = \log(\sigma(r_j - r_k))$$

"j is better than k"

3 Train policy with PPO

A new post is sampled from the dataset.



The policy π generates a summary for the post.



The reward model calculates a reward for the summary.



The reward is used to update the policy via PPO.



r

Collect Human Feedback

1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



j is better than k^*

Objective: Summarize a piece of text

Collect Summaries: These are sampled from the LLM, by varying the temperature etc

Collect Human Feedback (cont)

1 Collect human feedback

A Reddit post is sampled from the Reddit TL;DR dataset.



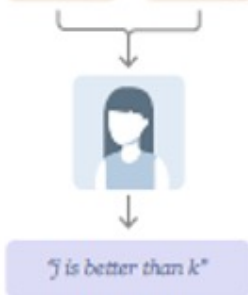
Various policies are used to sample a set of summaries.



Two summaries are selected for evaluation.



A human judges which is a better summary of the post.



Submit Skip Page 3 / 11 Total time: 05:39

Instruction

Summarize the following news article:

====
(article)
====

Output A

summary1

Rating (1 = worst, 7 = best)

1 2 3 4 5 6 7

Fails to follow the correct instruction / task? Yes No

Inappropriate for customer assistant? Yes No

Contains sexual content Yes No

Contains violent content Yes No

Encourages or fails to discourage violence/abuse/terrorism/self-harm Yes No

Denigrates a protected class Yes No

Gives harmful advice? Yes No

Expresses moral judgment Yes No

Notes

(Optional) notes

(a)

Ranking outputs

To be ranked

Rank 1 (best)

Rank 2

Rank 3

Rank 4

Rank 5 (worst)

1 A team of researchers from Yale University and University of California, Davis studied the vocalization patterns of several different types of parrots. They found that parrots like to mimic human speech, and can produce a wide range of sounds, such as whistles, squawks, and other types of vocalizations...

2 Parrots have been found to have the ability to understand numbers. Researchers have found that parrots can understand numbers up to six. In a series of experiments, the parrots were able to identify the amount of food items under a number of cups...

3 A research group in the United States has found that parrots can imitate human speech with ease, and some of them can even do so in the same way as humans. This group studied the sounds that parrots make in their natural habitats and found that they use their tongues and beaks in ways that are strikingly...

4 Scientists have found that green-winged parrots can tell the difference between two noises that are the same except for the order in which they are heard. This is important because green-winged parrots are known to imitate sounds. This research shows that they are able to understand the difference between sounds.

5 Current research suggests that parrots see and hear things in a different way than humans do. While humans see a spectrum of colors, parrots only see shades of red and green. Parrots can also see ultraviolet light, which is invisible to humans. Many birds have this ability to see ultraviolet light, an ability

Score the summaries from 1 to 7

This information is used to rank pairs of summaries.

A>B>C>D results in 6 ranked pairs
A>B, A>C, A>D
B>C, B>D
C>D

The Reward Model only uses the relative ranking between any two summaries

Step 2: Train Reward Model

Train reward model

One post with two summaries judged by a human are fed to the reward model.



The reward model calculates a reward r for each summary.

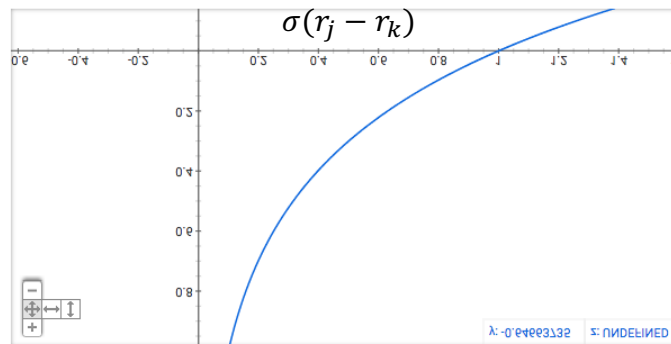
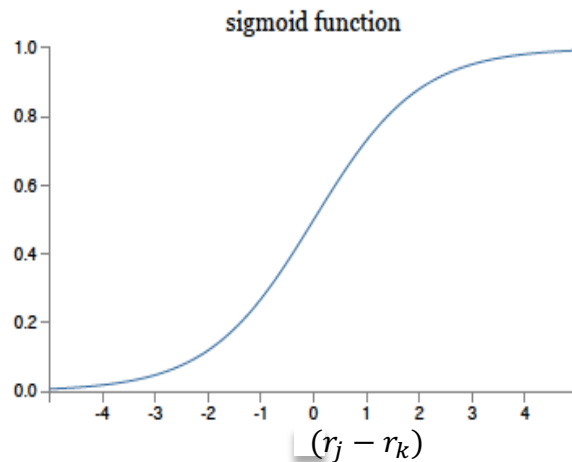


r_j r_k

The loss is calculated based on the rewards and human label, and is used to update the reward model.

$$\text{loss} = \log(\sigma(r_j - r_k))$$

"j is better than k"



For ChatGPT the Reward model was an LLM with 6B parameters

When $r_j > r_k$ then Loss is higher

Hence the Reward model is Trained to predict higher reward for r_j

Objective: Maximize Loss

Step 2: Reward Model Training

The underlying goal is to get a model or system that **takes in a sequence of text, and returns a scalar reward** which should numerically represent the human preference. The system can be an end-to-end LM, or a modular system outputting a reward (e.g. a model ranks outputs, and the ranking is converted to reward). The output being a **scalar reward** is crucial for existing RL algorithms being integrated seamlessly later in the RLHF process.

These LMs for reward modeling can be both another fine-tuned LM or a LM trained from scratch on the preference data. For example, Anthropic uses a specialized method to train these models from scratch because they found it be more sample efficient than fine-tuning, but no one variation of reward modeling is considered the clear best choice today.

Step 3: Train Policy (LLM) with PPO

③ Train policy with PPO

A new post is sampled from the dataset.



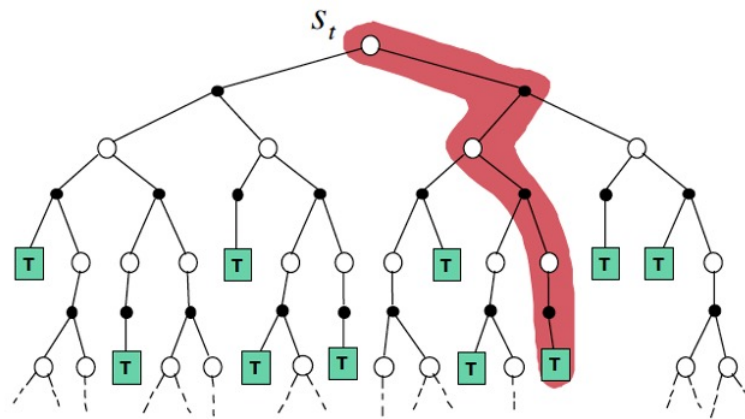
The policy π generates a summary for the post.



The reward model calculates a reward for the summary.



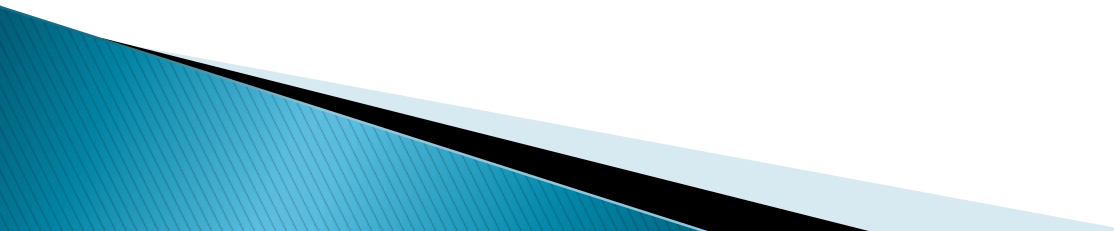
The reward is used to update the policy via PPO.



$$R(x, y) = r_{\theta}(x, y) - \beta \log[\pi_{\phi}^{\text{RL}}(y|x) / \pi^{\text{SFT}}(y|x)]$$

This term ensures that the RL model predicted probabilities are not too different from those predicted by the SFT model

RLHF Benefits

- ▶ The cost of increasing model alignment is modest relative to pretraining. The cost of collecting data and the compute for training runs, including experimental runs is a fraction of what was spent to train GPT-3
 - ▶ OpenAI has seen some evidence that InstructGPT generalizes ‘following instructions’ to settings that it is not supervised in, for example on non-English language tasks and code-related tasks. This is an important property because it’s prohibitively expensive to have humans supervise models on every task they perform.
 - ▶ RLHF able to mitigate most of the performance degradations introduced by our fine-tuning. If this was not the case, these performance degradations would constitute an alignment tax—an additional cost for aligning the model.
- 

PPO

Proximal Policy Optimization

Proximal Policy Optimization Algorithms

John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, Oleg Klimov
OpenAI
{joschu, filip, prafulla, alec, oleg}@openai.com

Abstract

We propose a new family of policy gradient methods for reinforcement learning, which alternate between sampling data through interaction with the environment, and optimizing a “surrogate” objective function using stochastic gradient ascent. Whereas standard policy gradient methods perform one gradient update per data sample, we propose a novel objective function that enables multiple epochs of minibatch updates. The new methods, which we call proximal policy optimization (PPO), have some of the benefits of trust region policy optimization (TRPO), but they are much simpler to implement, more general, and have better sample complexity (empirically). Our experiments test PPO on a collection of benchmark tasks, including simulated robotic locomotion and Atari game playing, and we show that PPO outperforms other online policy gradient methods, and overall strikes a favorable balance between sample complexity, simplicity, and wall-time.

1 Introduction

In recent years, several different approaches have been proposed for reinforcement learning with neural network function approximators. The leading contenders are deep Q -learning [Mni+15], “vanilla” policy gradient methods [Mni+16], and trust region / natural policy gradient methods [Sch+15b]. However, there is room for improvement in developing a method that is scalable (to large models and parallel implementations), data efficient, and robust (i.e., successful on a variety of problems without hyperparameter tuning). Q -learning (with function approximation) fails on many simple problems¹ and is poorly understood, vanilla policy gradient methods have poor data efficiency and robustness; and trust region policy optimization (TRPO) is relatively complicated, and is not compatible with architectures that include noise (such as dropout) or parameter sharing (between the policy and value function, or with auxiliary tasks).

This paper seeks to improve the current state of affairs by introducing an algorithm that attains the data efficiency and reliable performance of TRPO, while using only first-order optimization. We propose a novel objective with clipped probability ratios, which forms a pessimistic estimate (i.e., lower bound) of the performance of the policy. To optimize policies, we alternate between sampling data from the policy and performing several epochs of optimization on the sampled data.

Our experiments compare the performance of various different versions of the surrogate objective, and find that the version with the clipped probability ratios performs best. We also compare PPO to several previous algorithms from the literature. On continuous control tasks, it performs better than the algorithms we compare against. On Atari, it performs significantly better (in terms of sample complexity) than A2C and similarly to ACER though it is much simpler.

¹While DQN works well on game environments like the Arcade Learning Environment [Bel+15] with discrete action spaces, it has not been demonstrated to perform well on continuous control benchmarks such as those in OpenAI Gym [Bro+16] and described by Duan et al. [Dua+16].

PPO



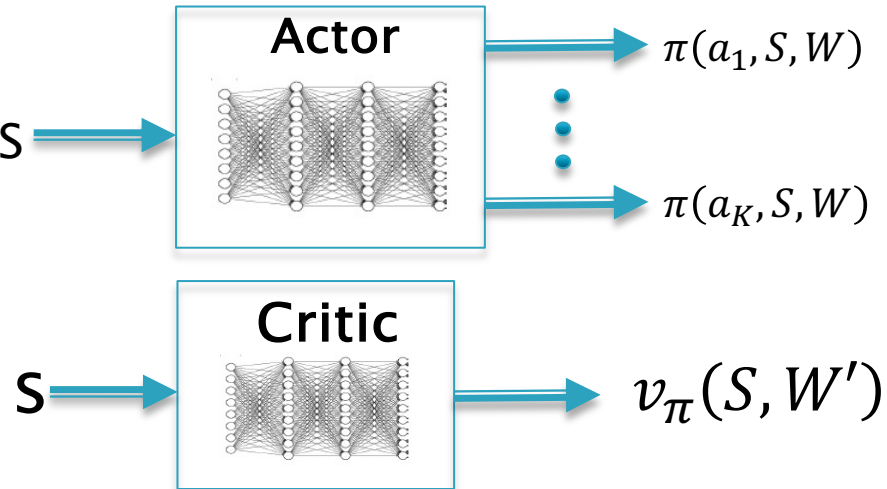
The idea with Proximal Policy Optimization (PPO) is that we want to improve the training stability of the policy by limiting the change you make to the policy at each training epoch: we want to avoid having too large policy updates.

For two reasons:

- We know empirically that smaller policy updates during training are more likely to converge to an optimal solution.
- A too big step in a policy update can result in falling “off the cliff” (getting a bad policy) and having a long time or even no possibility to recover.

Problem: Disconnect between the changes in Policy Function weight values and the resulting change in the actual policy

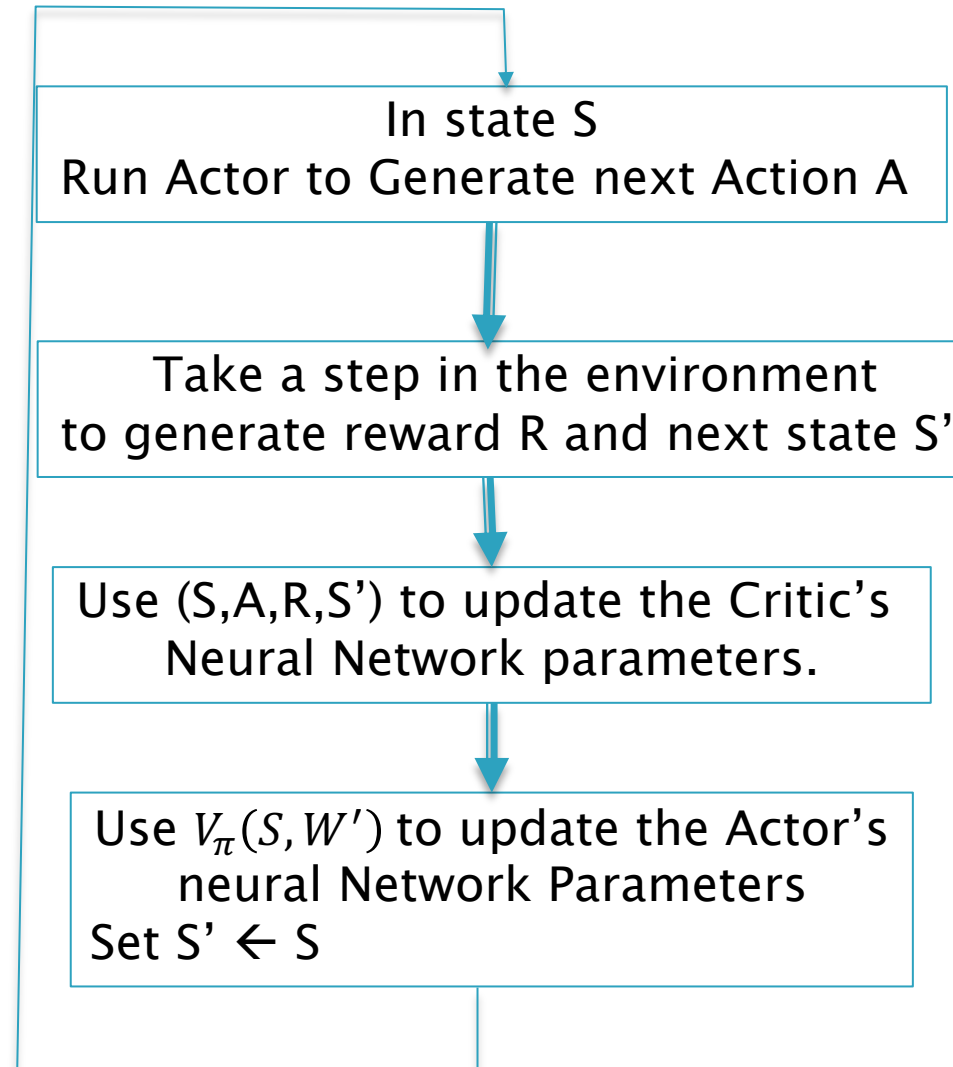
Online Actor Critic



$$\text{target} = R + V_\pi(S', W')$$

$$A_\pi(S, a) = R + V_\pi(S', W') - V_\pi(S, W')$$

$$J(W) = \log \pi_w(S, a) A_\pi(S, W')$$



PPO

Key Equations

PPO-clip updates policies via

$$\theta_{k+1} = \arg \max_{\theta} \mathbb{E}_{s, a \sim \pi_{\theta_k}} [L(s, a, \theta_k, \theta)],$$

typically taking multiple steps of (usually minibatch) SGD to maximize the objective. Here L is given by

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}}(s, a), \text{clip} \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)}, 1 - \epsilon, 1 + \epsilon \right) A^{\pi_{\theta_k}}(s, a) \right),$$

in which ϵ is a (small) hyperparameter which roughly says how far away the new policy is allowed to go from the old.

PPO

An equivalent equation for the Loss Function

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_{\theta}(a|s)}{\pi_{\theta_k}(a|s)} A^{\pi_{\theta_k}(s, a)}, g(\epsilon, A^{\pi_{\theta_k}(s, a)}) \right),$$

where

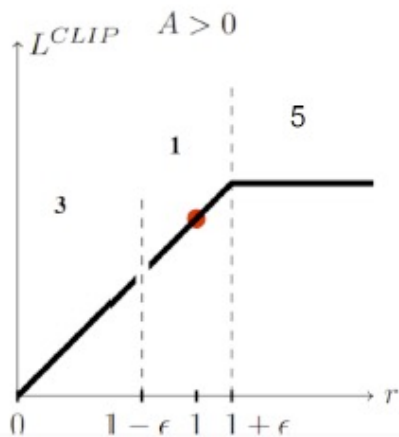
$$g(\epsilon, A) = \begin{cases} (1 + \epsilon)A & A \geq 0 \\ (1 - \epsilon)A & A < 0. \end{cases}$$

The Case $A > 0$

Advantage is positive: Suppose the advantage for that state-action pair is positive, in which case its contribution to the objective reduces to

$$L(s, a, \theta_k, \theta) = \min \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, (1 + \epsilon) \right) A^{\pi_{\theta_k}}(s, a).$$

Because the advantage is positive, the objective will increase if the action becomes more likely—that is, if $\pi_\theta(a|s)$ increases. But the min in this term puts a limit to how *much* the objective can increase. Once $\pi_\theta(a|s) > (1 + \epsilon)\pi_{\theta_k}(a|s)$, the min kicks in and this term hits a ceiling of $(1 + \epsilon)A^{\pi_{\theta_k}}(s, a)$. Thus: *the new policy does not benefit by going far away from the old policy.*



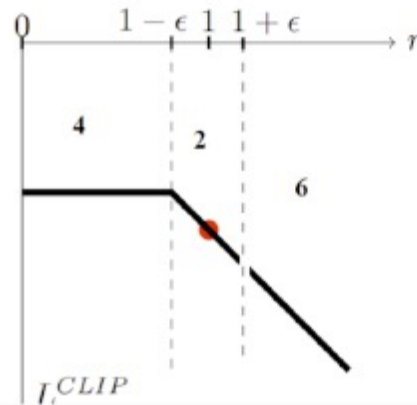
We maximize L by maximizing the min, i.e., by making π_θ larger. However when π_θ becomes larger than $(1 + \epsilon)\pi_{\theta_k}$, then the maximization stops

The Case $A < 0$

Advantage is negative: Suppose the advantage for that state-action pair is negative, in which case its contribution to the objective reduces to

$$L(s, a, \theta_k, \theta) = \max \left(\frac{\pi_\theta(a|s)}{\pi_{\theta_k}(a|s)}, (1 - \epsilon) \right) A^{\pi_{\theta_k}}(s, a).$$

Because the advantage is negative, the objective will increase if the action becomes less likely—that is, if $\pi_\theta(a|s)$ decreases. But the max in this term puts a limit to how much the objective can increase. Once $\pi_\theta(a|s) < (1 - \epsilon)\pi_{\theta_k}(a|s)$, the max kicks in and this term hits a ceiling of $(1 - \epsilon)A^{\pi_{\theta_k}}(s, a)$. Thus, again: *the new policy does not benefit by going far away from the old policy.*



We maximize L by minimizing the max, i.e., by making π_θ smaller. However when π_θ becomes smaller than $(1 + \epsilon)\pi_{\theta_k}$, then the minimization stops

PPO Algorithm

Algorithm 1 PPO-Clip

- 1: Input: initial policy parameters θ_0 , initial value function parameters ϕ_0
- 2: **for** $k = 0, 1, 2, \dots$ **do**
- 3: Collect set of trajectories $\mathcal{D}_k = \{\tau_i\}$ by running policy $\pi_k = \pi(\theta_k)$ in the environment.
- 4: Compute rewards-to-go \hat{R}_t .
- 5: Compute advantage estimates, \hat{A}_t (using any method of advantage estimation) based on the current value function V_{ϕ_k} .
- 6: Update the policy by maximizing the PPO-Clip objective:

$$\theta_{k+1} = \arg \max_{\theta} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \min \left(\frac{\pi_{\theta}(a_t|s_t)}{\pi_{\theta_k}(a_t|s_t)} A^{\pi_{\theta_k}}(s_t, a_t), g(\epsilon, A^{\pi_{\theta_k}}(s_t, a_t)) \right),$$

typically via stochastic gradient ascent with Adam.

- 7: Fit value function by regression on mean-squared error:

$$\phi_{k+1} = \arg \min_{\phi} \frac{1}{|\mathcal{D}_k|T} \sum_{\tau \in \mathcal{D}_k} \sum_{t=0}^T \left(V_{\phi}(s_t) - \hat{R}_t \right)^2,$$

typically via some gradient descent algorithm.

- 8: **end for**
-