

Instruct-GPT: Reinforcement Learning with Human Feedback and Proximal Policy Optimization

Erhard Hinrichs

Seminar für Sprachwissenschaft
Eberhard-Karls Universität Tübingen

References

The subsequent slides are taken from the following materials:

- ▶ L. Ouyang et al. (2022). Training language models to follow instructions with human feedback. arXiv:2203.02155v1
- ▶ Rafael Rafailov et al. (2023). Direct Preference Optimization: Your Language Model is Secretly a Reward Model.
<https://arxiv.org/abs/2305.18290>
- ▶ Luis Serrano, Proximal Policy Optimization (PPO).
https://www.youtube.com/watch?v=TjHH_--7l8g
- ▶ Luis Serrano, Re-inforcement Learning with Human Feedback.
https://www.youtube.com/watch?v=Z_JUqJBpVOk
- ▶ Luis Serrano, Direct Preference Optimization (DPO) - How to fine-tune LLMs directly without reinforcement learning
<https://www.youtube.com/watch?v=k2pD3k1485A>

RLHF and PPO

- ▶ Reinforcement Learning with Human Feedback (RLHF) is a method used for training Large Language Models (LLMs), including chat-GPT and its predecessor Instruct-GPT.
- ▶ RLHF in turn is based on gradient-based reinforcement learning method called Proximal Policy Optimization (PPO).
- ▶ In PPO two networks are trained simultaneously
 - ▶ a value network
 - ▶ a policy network

Problems with Prompt-Driven Large Language Models (LLMs)

Misalignment of user intention and model response

- ▶ hallucinations: generated content that is nonsensical or unfaithful to the provided source content; "making up facts"
- ▶ generating biased text
- ▶ toxic language: communication that is harmful to others
- ▶ due to different objectives of LLMs ("predict the next word or sentence") and of machine-human interaction ("follow the user's instructions helpfully and safely")

Instruct-GPT: Fine-Tuning, Reward Modelling, and Reinforcement Learning

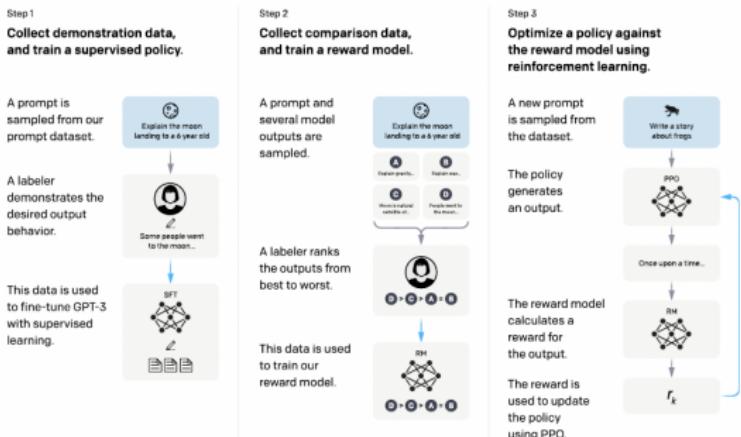


Figure 2: A diagram illustrating the three steps of our method: (1) supervised fine-tuning (SFT), (2) reward model (RM) training, and (3) reinforcement learning via proximal policy optimization (PPO) on this reward model. Blue arrows indicate that this data is used to train one of our models. In Step 2, boxes A-D are samples from our models that get ranked by labelers. See Section 3 for more details on our method.

Figure taken from: L. Ouyang et al. (2022).

Results of Human Evaluation of Prompt Distributions

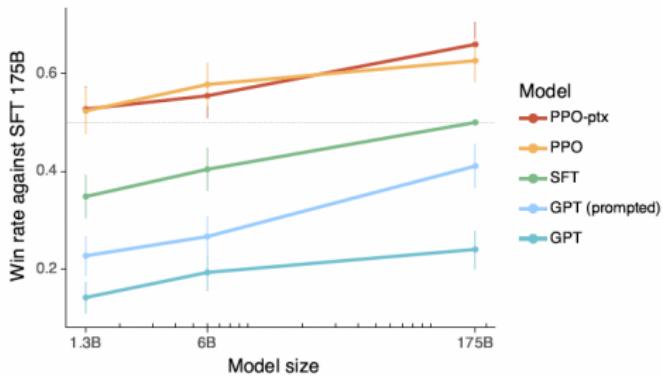


Figure 1: Human evaluations of various models on our API prompt distribution, evaluated by how often outputs from each model were preferred to those from the 175B SFT model. Our InstructGPT models (PPO-ptx) as well as its variant trained without pretraining mix (PPO) significantly outperform the GPT-3 baselines (GPT, GPT prompted); outputs from our 1.3B PPO-ptx model are preferred to those from the 175B GPT-3. Error bars throughout the paper are 95% confidence intervals.

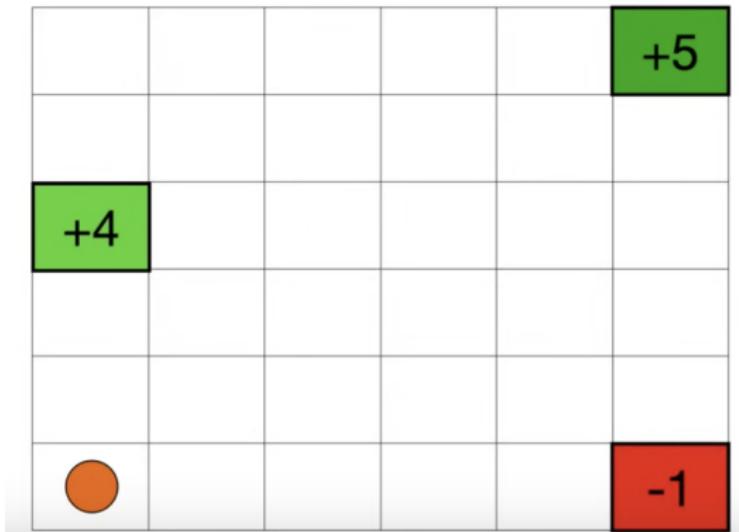
Figure taken from: L. Ouyang et al. (2022). Training language models to follow instructions with human feedback.

arXiv:2203.02155v1

GridWorld Game

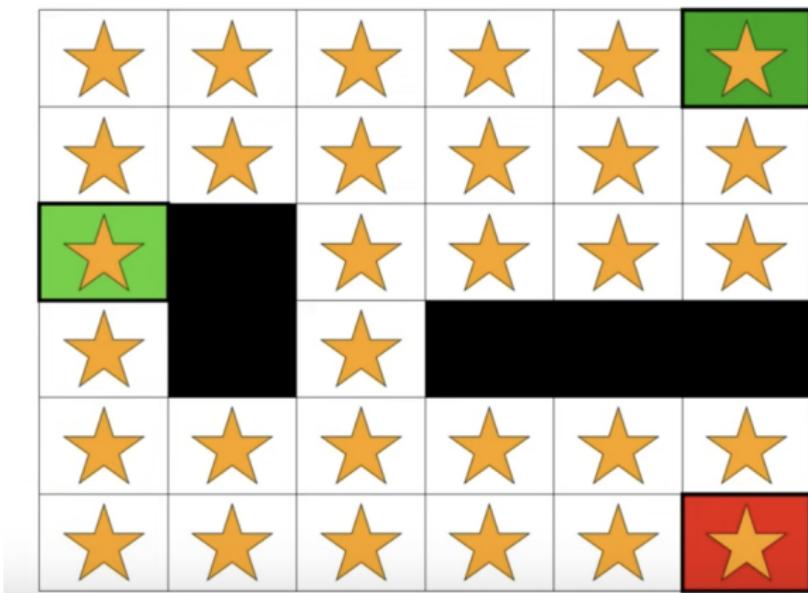
Objective: To get a player to learn how to play the game in an optimal way

Gridworld



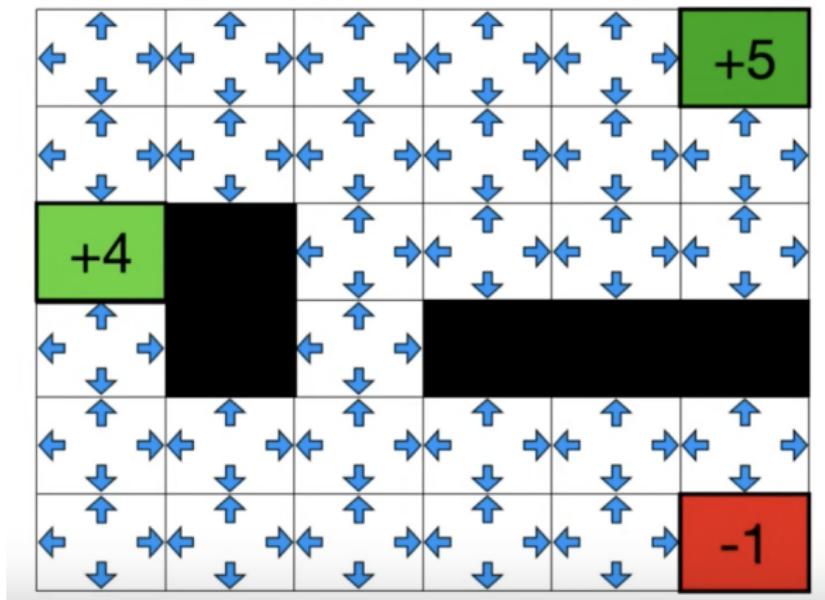
GridWorld Game: States

States



GridWorld Game: Actions

Actions



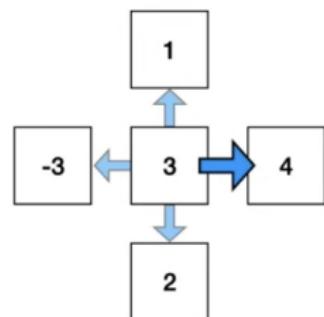
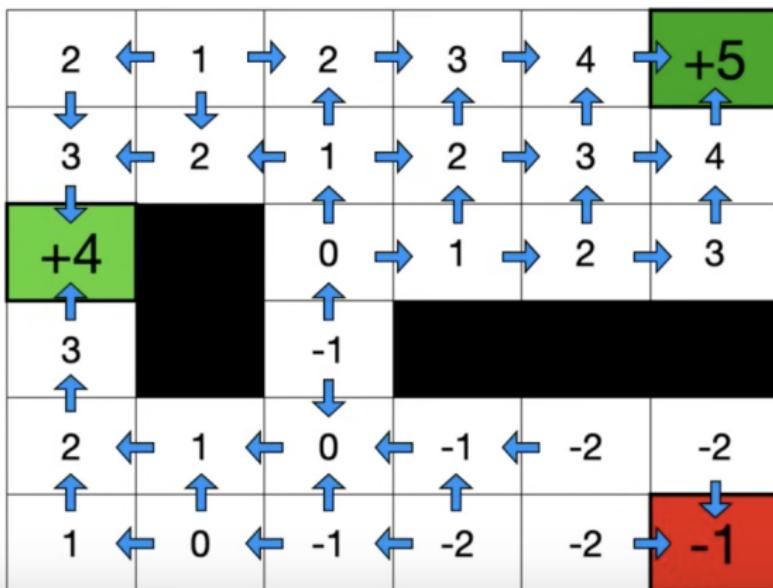
GridWorld Game: Values

How to calculate all the values?

2	1	2	3	4	+5
3	2	1	2	3	4
+4		0	1	2	3
3		-1			
2	1	0	-1	-2	-2
1	0	-1	-2	-2	-1

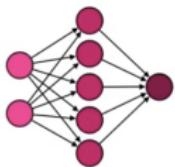
GridWorld Game: Policy

Where is the best place to move?



Two neural networks

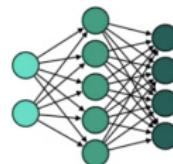
Value
neural
network



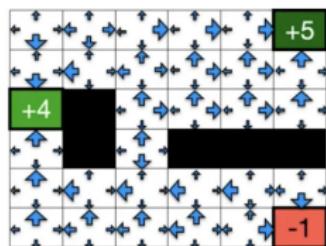
Approximates
values

2	1	2	3	4	+5
3	2	1	2	3	4
+4		0	1	2	3
3		-1			
2	1	0	-1	-2	-2
1	0	-1	-2	-2	-1

Policy
neural
network



Approximates
policy



Training the value neural network

$$\left(\text{Value neural network} - \text{Actual values} \right)^2$$

The diagram illustrates the calculation of the squared error loss function for training a value neural network. It shows three components:

- Value neural network:** A 3x3 matrix of predicted values:

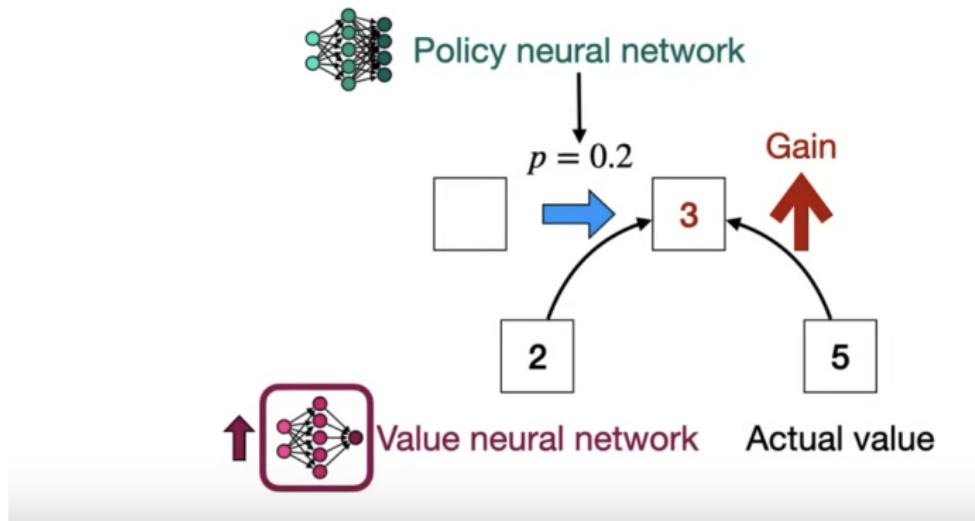
0.4	0.5	1.2
-0.3	0.6	
0.8		
- Actual values:** A 3x3 matrix of target values:

2	3	4
0	1	
-1		
- Squared Error:** The result of the subtraction followed by squaring:

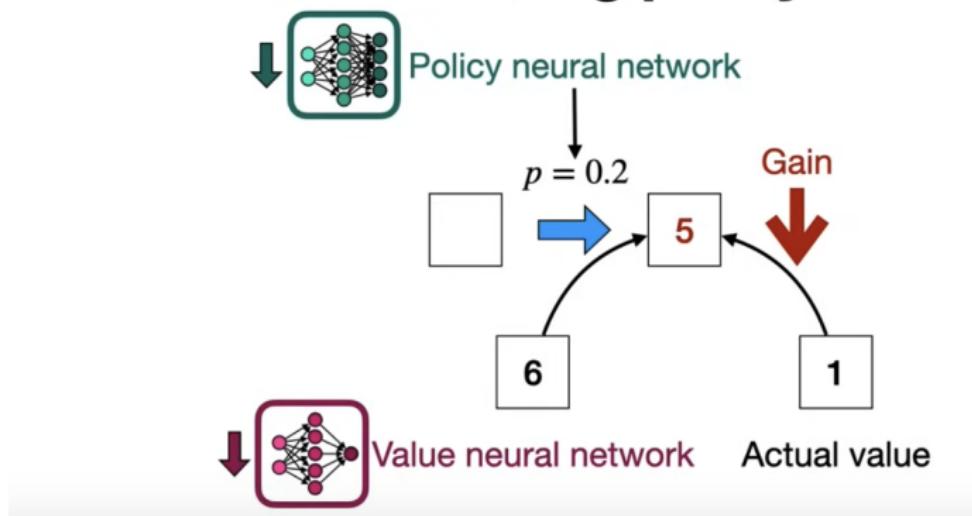
2.56	6.25	7.84
0.09	0.16	
3.24		

The final result is the sum of all these squared differences.

Main idea for training policy

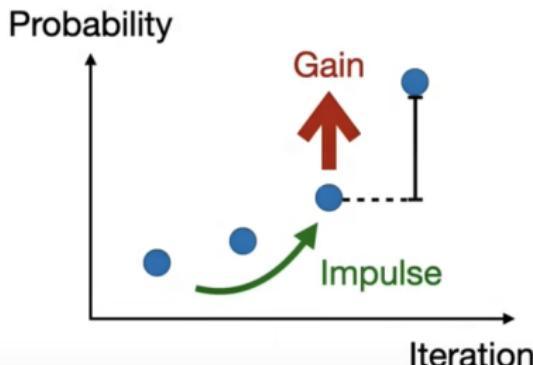


Main idea for training policy

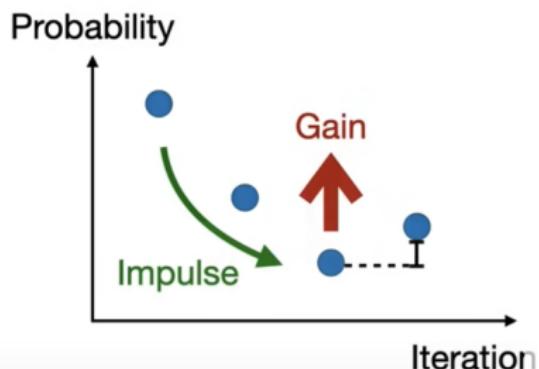


Impulse (Momentum)

Scenario 1 →

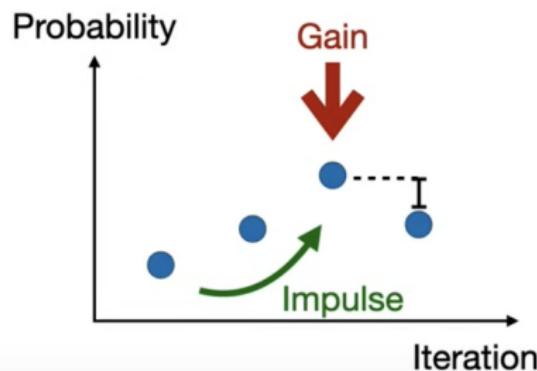


Scenario 2 →

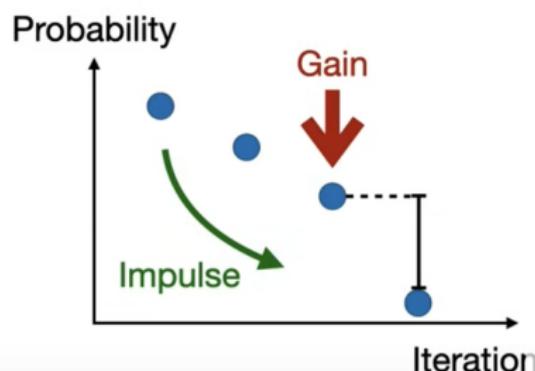


Impulse (Momentum)

Scenario 3 



Scenario 4 



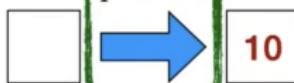
Impulse (Momentum)

Scenario 1 →

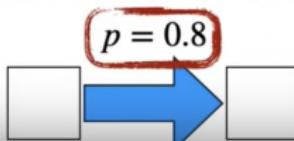
Iteration 99



Iteration 100



Iteration 101

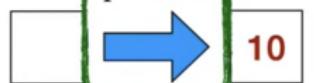


Scenario 2 →

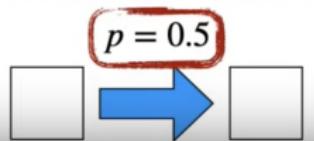
Iteration 99



Iteration 100



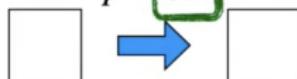
Iteration 101



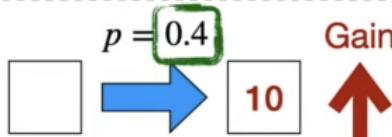
Dividing by the previous probability

Scenario 1 →

Iteration 99



Iteration 100



Surrogate
objective
function

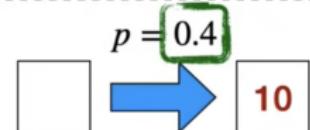
$$\frac{0.4}{0.2} \textcolor{red}{10} = \textcolor{red}{20}$$

Scenario 2 →

Iteration 99



Iteration 100

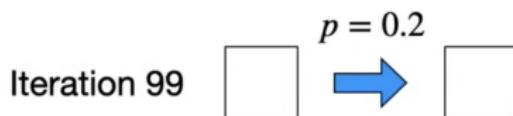


$$\frac{0.4}{0.8} \textcolor{red}{10} = \textcolor{red}{5}$$

Surrogate Objective Function

Train neural network to make
this as small as possible

$$L_{\text{policy}}(\theta) = \mathbb{E} \frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} A_t$$



$$\frac{0.4}{0.2} \quad \textcolor{brown}{10} = \boxed{20}$$

Iteration 100

$p = 0.4$

Solution: Clipping

Current probability

$$\pi_{\theta}(a_t | s_t)$$

Previous probability

$$\pi_{\theta_{\text{old}}}(a_t | s_t)$$

If between 0.7 and 1.3, keep

If too small: Make it 0.7

If too large: Make it 1.3



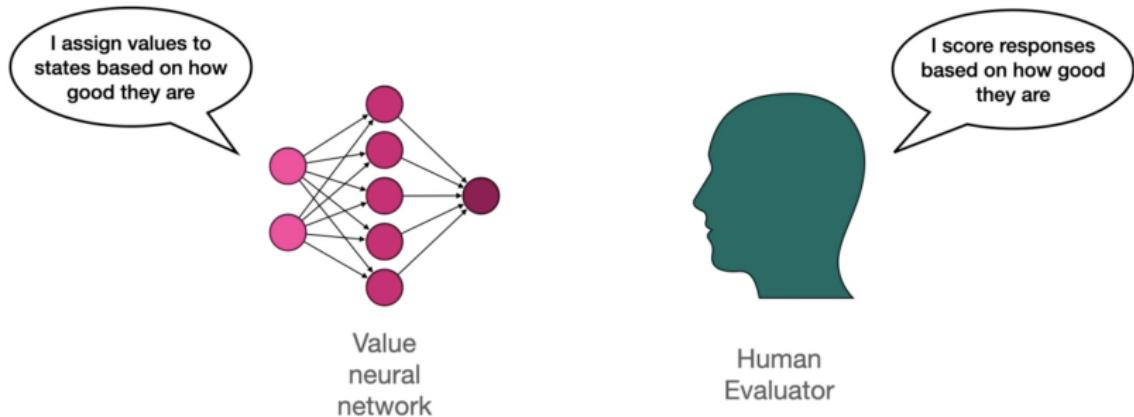
Solution: Clipping

$$L_{\text{policy}}^{\text{CLIP}}(\theta) = \mathbb{E} \left[\min \left(\frac{\pi_\theta(a_t | s_t)}{\pi_{\theta_{\text{old}}}(a_t | s_t)} A_t, \text{clip}(\rho, 1 - \epsilon, 1 + \epsilon) A_t \right) \right]$$

 $\left\{ \begin{array}{l} \text{If between } 1 - \epsilon \text{ and } 1 + \epsilon, \text{ keep} \\ \text{If smaller than } 1 - \epsilon, \text{ keep } 1 - \epsilon \\ \text{If larger than } 1 + \epsilon: \text{ Make it } 1 + \epsilon \end{array} \right\}$

Re-Information Learning with Human Feedback

Value network learns to mimic the human evaluator



Training the value network

			What color is the sky? Red		
		What color is the sky? Blue	What color is the sky?	What color is the sky? Banana	
			What color is the sky		
	What color	What color is	What color is the		
<start>	What				



What color is the sky?
Blue



What color is the sky?
Red



What color is the sky?
Banana



Training the value network

			+2		
	+3	What color is the sky?		+1	
		What color is the sky			
<start>	What		What color is the		

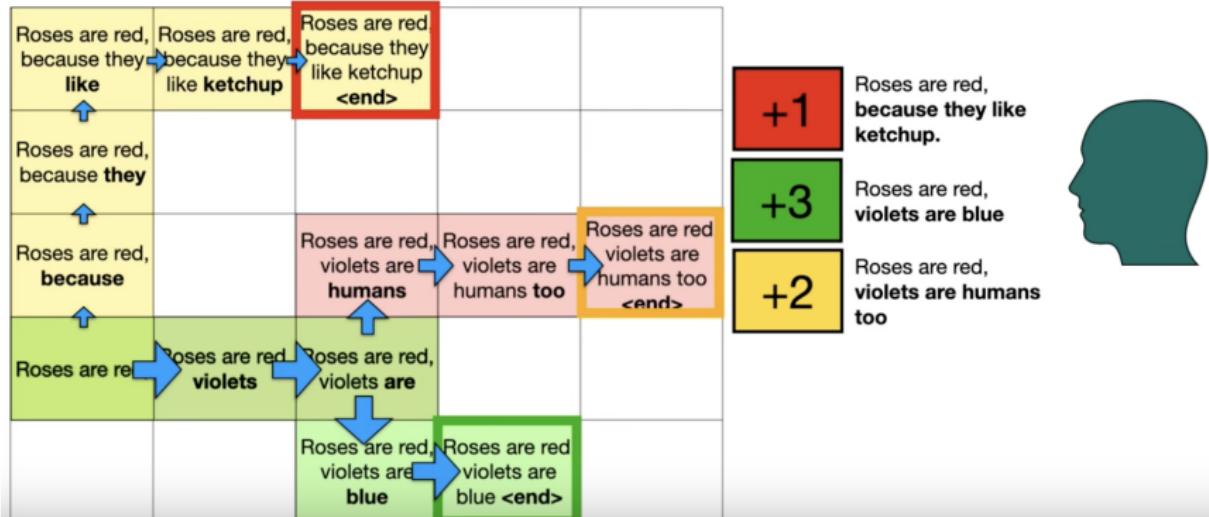
- +3 What color is the sky?
Blue
- +2 What color is the sky?
Red
- +1 What color is the sky?
Banana



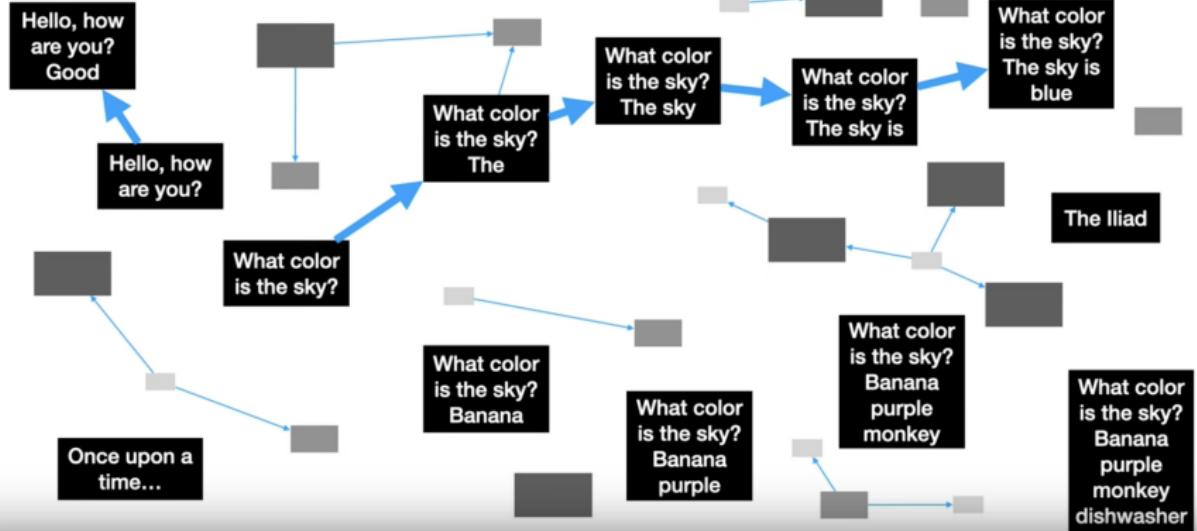
Value
neural
network



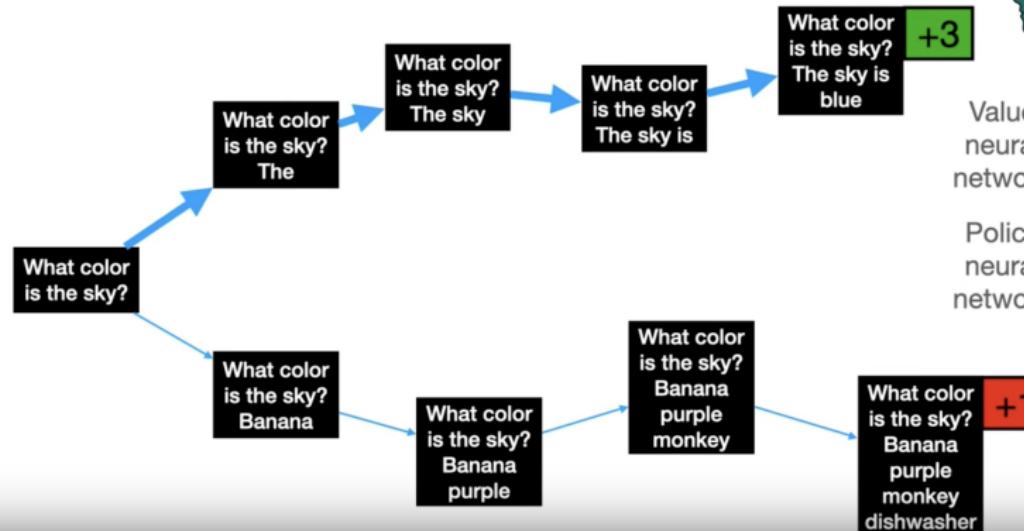
Longer paths



Obviously it's not a grid



Obviously it's not a grid



Value
neural
network



Policy
neural
network



Illustrative user prompts from InstructGPT/GPT-3 distributions

Use Case	Example
brainstorming	indie movie ideas <ul style="list-style-type: none">- A guy travels to South America to become a shaman.- A documentary about the world of juggling.
brainstorming	Tell me a list of topics related to: <ul style="list-style-type: none">- interior design- sustainable ecosystem- fake plants
classification	The following is a list of products and the kind of product they are. Product: product. Type: type Product: product. Type: type Product: product. Type:

Illustrative user prompts from InstructGPT/GPT-3 distributions

Use Case	Example
classification	This is a tweet sentiment classifier. tweet Sentiment: negative ==== tweet Sentiment: neutral ==== tweet Sentiment:

Categories for user prompts from InstructGPT/GPT-3

Brainstorming

Chat

Classification

Extract

Generation

QA, closed

QA, open

Rewrite

Summarization

Other

Illustrative user prompts from InstructGPT/GPT-3 distributions

Use Case	Example
generation	<p>Write an outline for an essay about John von Neumann and his contributions to computing:</p> <p>I. Introduction, his life and background</p> <p>A: His early life</p> <p>B:</p>
rewrite	<p>Rephrase this for me: "I can't seem to find out how to work this darn thing."</p> <p>Alternate phrasing: "</p>
rewrite	<p>Covert my resume into a profile overview.</p> <p>{ resume }</p> <p>Profile overview:</p>

Illustrative user prompts from InstructGPT/GPT-3 distributions

Use Case	Example
chat	<p>The following is a conversation with an AI assistant.</p> <p>The assistant is helpful, creative, clever, and very friendly.</p> <p>Human: Hello, who are you?</p> <p>AI: I am an AI created by OpenAI. How can I help you today?</p> <p>Human: I'm feeling kind of down today.</p> <p>AI:</p>
closed qa	<p>When you drop a heavy stone from a tree, what happens?</p> <p>A. The stone falls to the ground.</p> <p>B: The stone stays in the tree.</p> <p>C: The stone floats.</p> <p>D: Nothing happens.</p> <p>Answer:</p>

Illustrative user prompts from InstructGPT/GPT-3 distributions

Use Case	Example
closed qa	<p>Text: { article describing what yoga mats to buy }</p> <p>Question: What are the things I should consider when buying a yoga mat?</p> <p>Answer:</p>
open qa	<p>Q: Who is Leonardo da Vinci?</p> <p>A:</p>

InstructGPT vs. GPT-3 (Ouyang et al. 2022)

- ▶ Labelers significantly prefer InstructGPT outputs over outputs from GPT-3.
- ▶ InstructGPT models show improvements in truthfulness over GPT-3.
- ▶ InstructGPT shows small improvements in toxicity over GPT-3, but not bias.
- ▶ InstructGPT models show promising generalization to instructions outside of the RLHF finetuning distribution.
- ▶ InstructGPT models generalize to the preferences of “held-out” labelers that did not produce any training data.
- ▶ InstructGPT still makes simple mistakes.

Direct Preference Optimization

Quote from the abstract of Rafailov et al. 2023:

In this paper we introduce a new parameterization of the reward model in RLHF that enables extraction of the corresponding optimal policy in closed form, allowing us to solve the standard RLHF problem with only a simple classification loss. The resulting algorithm, which we call Direct Preference Optimization (DPO), is stable, performant, and computationally lightweight, eliminating the need for sampling from the LM during fine-tuning or performing significant hyperparameter tuning. Our experiments show that DPO can fine-tune LMs to align with human preferences as well as or better than existing methods. Notably, fine-tuning with DPO exceeds PPO-based RLHF in ability to control sentiment of generations, and matches or improves response quality in summarization and single-turn dialogue while being substantially simpler to implement and train.