

10-423/623: Generative AI Lecture 11 – Reinforcement Learning with Human Feedback

Henry Chai & Matt Gormley

10/2/24

Front Matter

- Announcements:
 - HW2 released ~~9/23~~ 9/24, due 10/7 at 11:59 PM

Recall: Few-shot Learning with LLMs

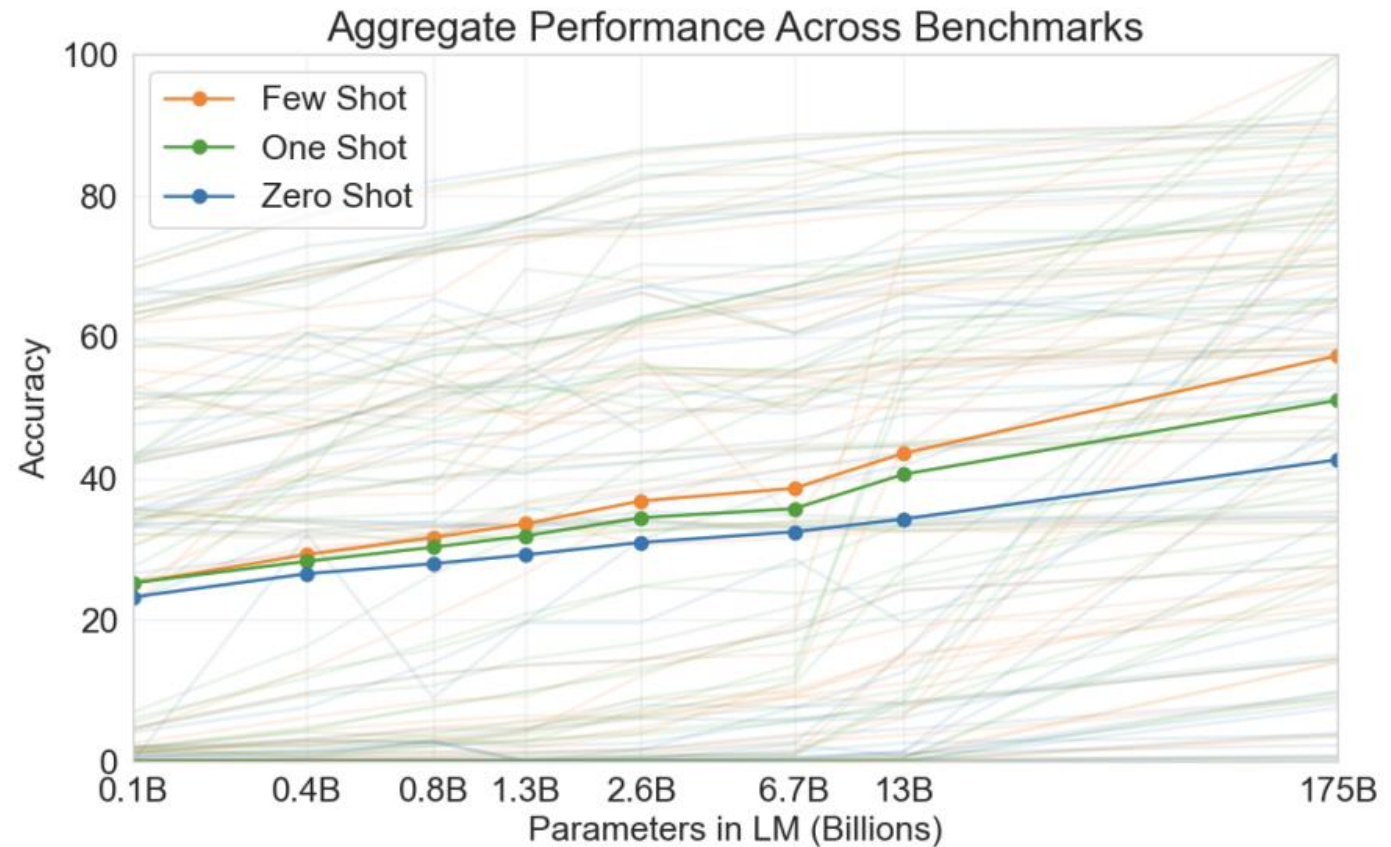
- Suppose you have...
 1. a small labelled dataset (i.e., a few-shot setting), \mathcal{D}
 2. a very large pre-trained language model
- There are two ways to “learn”:
 - A. Supervised fine-tuning i.e., updating the LLM’s parameters using
 1. a standard supervised objective
 2. backpropagation to compute gradients
 3. your favorite optimizer (e.g., Adam)
- **Con:** backpropagation requires $\sim 3x$ the memory and computation time as the forward computation
- **Con:** you might not have access to the model parameters

Recall: Few-shot Learning with LLMs

- Suppose you have...
 1. a small labelled dataset (i.e., a few-shot setting), \mathcal{D}
 2. a very large pre-trained language model
- There are two ways to “learn”:
 - B. In-context learning i.e., feeding the training dataset to the LLM as a prompt and taking the output as a prediction
 - the LLM (hopefully) infers patterns in the training dataset during inference (i.e., decoding)
 - **Pro:** no backpropagation required and only one pass through the training dataset *per test example*
 - **Pro:** does not require access to the model parameters, only API access to the model itself
 - **Con:** the prompt may be very long and Transformer LMs require $O(N^2)$ time/space where N = length of context

Recall: Few-shot Learning via In-context Learning with LLMs

- Standard setup: a set of input/output pairs from a training dataset are presented in sequence to an LLM, typically along with a plain-text task description



Recall: Few-shot In-context Learning with LLMs

- In-context learning is surprisingly *sensitive* to...
 1. the order the training examples are presented in
 2. label imbalance (e.g. # of positive vs. # of negative)
 3. the number of unique labels in the training dataset
- In-context learning is surprisingly *insensitive* to...
 1. the correctness of the labels!
 2. the amount of training data used in the prompt!

So why does *this* work? Why is it better than zero-shot learning?

- In-context learning is surprisingly *sensitive* to...
 1. the order the training examples are presented in
 2. label imbalance (e.g. # of positive vs. # of negative)
 3. the number of unique labels in the training dataset
- In-context learning is surprisingly *insensitive* to...
 1. the correctness of the labels!
 2. the amount of training data used in the prompt!

Few-shot In-context Learning with LLMs

- Min et al. (2022) identified four meaningful factors:

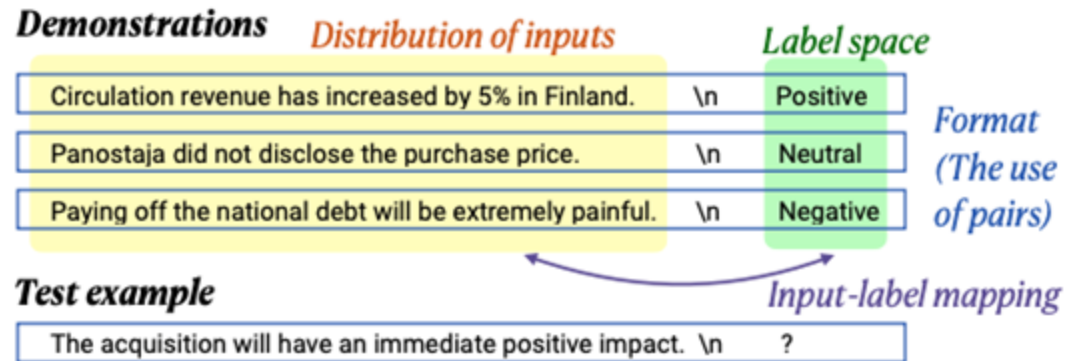


Figure 7: Four different aspects in the demonstrations: the input-label mapping, the distribution of the input text, the label space, and the use of input-label pairing as the format of the demonstrations.

- Another potentially meaningful aspect of in-context learning: what *exactly* are we asking the LLM?

Prompt Engineering

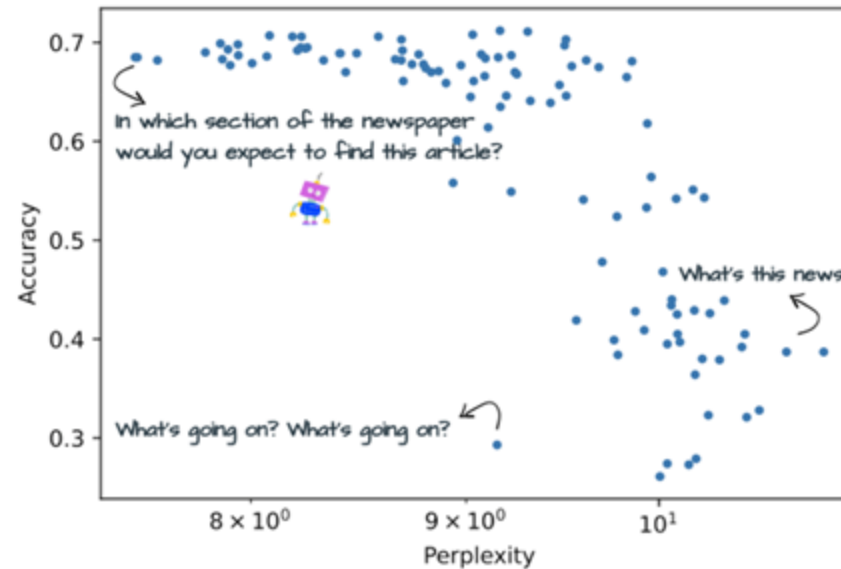
- Not all prompts are equally good!
- Example: zero-shot news article classification using OPT-175B on the AG News dataset

Prompt	Accuracy
What is this piece of news regarding?	40.9
What is this article about?	52.4
What is the best way to describe this article?	68.2
What is the most accurate label for this news article?	71.2

- What affects the accuracy associated with using a prompt?
- One potential answer: how likely the prompt is under the learned model's implied distribution over sequences

Prompt Engineering

- Not all prompts are equally good!
- Example: zero-shot news article classification using OPT-175B on the AG News dataset



- Perplexity is the exponentiated average negative log-likelihood of a sequence
- Lower perplexity = higher likelihood

Prompt Engineering

- Not all prompts are equally good!
- Example: zero-shot news article classification using OPT-175B on the AG News dataset

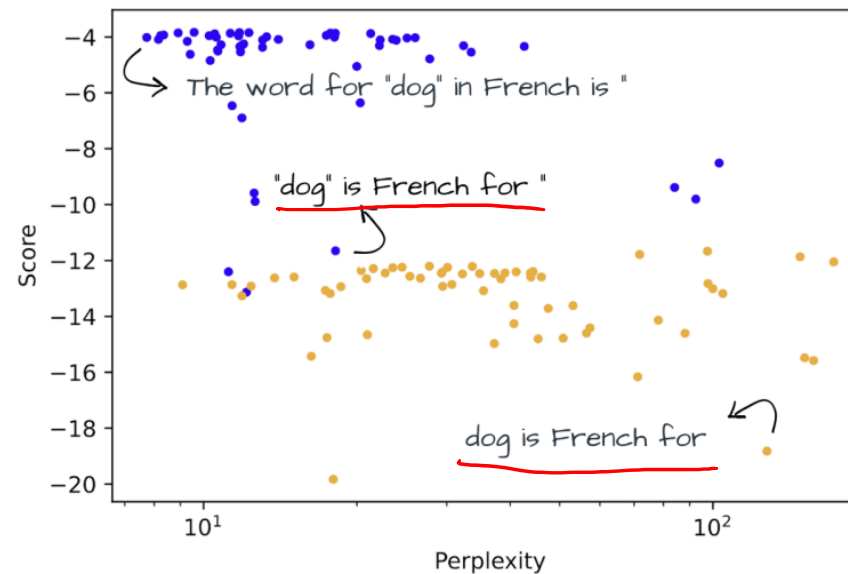


Figure 2: Score of correct label vs. perplexity for the word-level translation task in French with OPT 175B. The x axis is in log scale. The blue points stand for prompts with quotation marks for the words, while the yellow points are of prompts without quotation marks.

- Perplexity is the exponentiated average negative log-likelihood of a sequence
- Lower perplexity = higher likelihood

Learning to Prompt

- Some ways of *learning* better prompts for your task:
 1. Prompt paraphrasing – programmatically generate and test many different prompts from a paraphrase model, then pick the one that works best
 2. Gradient-based search – use optimization to search for the discrete representation of the prompt that makes the desired output most likely
 3. Prompt tuning – directly optimize the embeddings that are input into the LLM, without bothering to construct a discrete representation of the prompt

Chain-of-Thought Prompting

- Insight: asking an LLM to reason about its answer can improve its in-context performance
- Chain-of-thought prompting provides examples of reasoning in the in-context training examples

Standard Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

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

Model Output

A: The answer is 27. ❌

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

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

Model Output

A: 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 $3 + 6 = 9$. The answer is 9. ✅

Chain-of-Thought Prompting

- Insight: asking an LLM to reason about its answer can improve its in-context performance
- Chain-of-thought prompting provides examples of reasoning in the in-context training examples

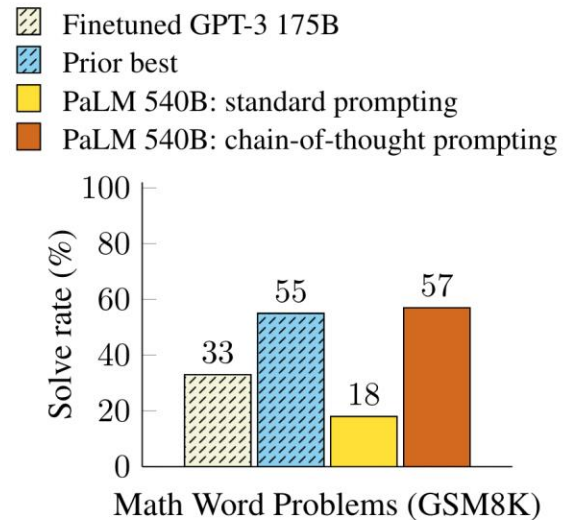


Figure 2: PaLM 540B uses chain-of-thought prompting to achieve new state-of-the-art performance on the GSM8K benchmark of math word problems. Finetuned GPT-3 and prior best are from Cobbe et al. (2021).

Chain-of-Thought Prompting

Model Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

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

Model Output

A: 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 $3 + 6 = 9$. The answer is 9. ✓

Chain-of-Thought Prompting

- Insight: asking an LLM to reason about its answer can improve its in-context performance
- Chain-of-thought prompting provides examples of reasoning in the in-context training examples

(a) Few-shot

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) *The answer is 8. ❌*

(b) Few-shot-CoT

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. $5 + 6 = 11$. The answer is 11.

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A:

(Output) *The juggler can juggle 16 balls. Half of the balls are golf balls. So there are $16 / 2 = 8$ golf balls. Half of the golf balls are blue. So there are $8 / 2 = 4$ blue golf balls. The answer is 4. ✔️*

(c) Zero-shot

Q: A juggler can juggle 16 balls. Half of the balls are golf balls, and half of the golf balls are blue. How many blue golf balls are there?

A: The answer (arabic numerals) is

(Output) *8 ❌*

Chain-of-Thought Prompting

- Insight: asking an LLM to reason about its answer can improve its in-context performance
- Chain-of-thought prompting provides examples of reasoning in the in-context training examples

	MultiArith	GSM8K
Zero-Shot	17.7	10.4
Few-Shot (2 samples)	33.7	15.6
Few-Shot (8 samples)	33.8	15.6
Zero-Shot-CoT	78.7	40.7
Few-Shot-CoT (2 samples)	84.8	41.3

Chain-of-Thought Prompting

- Insight: asking an LLM to reason about its answer can improve its in-context performance
- Chain-of-thought prompting provides examples of reasoning in the in-context training examples

September 12, 2024

Learning to Reason with LLMs

We are introducing OpenAI o1, a new large language model trained with reinforcement learning to perform complex reasoning. o1 thinks before it answers —it can produce a long internal chain of thought before responding to the user.

Recall: Zero-shot Learning with LLMs

- Summarization

Prompt:

Story: Jason listened to the weather and heard it was going to be sunny. He thought the kids might like to go swimming. He gathered up the swimsuits, towels and sunscreen. Jason and the kids got into the truck and drove to the beach. They spent the next 2 hours playing and splashing in the surf.

One-sentence Summary:

Model output (Llama-2-70B):

Jason took the kids swimming at the beach after hearing the weather forecast, gathering necessary items and driving there.



Model output (Llama-2-7B Chat):

Sure! Here is a one-sentence summary of the story:

Jason and the kids went on a fun-filled day at the beach, playing and splashing in the surf after gathering all the necessary items, including swimsuits, towels, and sunscreen.

Instruction Fine-Tuning

- LLMs are trained to reduce the perplexity of its training corpus i.e., it's good at completing _____
- Issue: an interactive agent (e.g., a chatbot) can't just predict what comes next, it needs to *behave* in a certain way (e.g., conversationally)
 - We need to *align* the LLM with the expectations of a human user for a given task
- Idea: build a “chatbot” training dataset and use it to fine-tune the LLM
 - This approach is called instruction fine-tuning, chat fine-tuning, alignment, behavioral fine-tuning, etc...

So how could we build such a dataset? Where can we get...

- Prompts?

- from real people
- ask an LLM
- in structured settings

- Response?

- from real people
- ask an LLM?
- in structured settings

- LLMs are trained to reduce the perplexity of its training corpus i.e., it's good at completing _____
- Issue: an interactive agent (e.g., a chatbot) can't just predict what comes next, it needs to *behave* in a certain way (e.g., conversationally)
 - We need to *align* the LLM with the expectations of a human user for a given task
- Idea: build a **“chatbot” training dataset** and use it to fine-tune the LLM
 - This approach is called instruction fine-tuning, chat fine-tuning, alignment, behavioral fine-tuning, etc...

Instruction Fine-Tuning Dataset: InstructGPT

- Consists of 13k prompt/response pairs
- Human labelers wrote both an instruction prompt and a demonstration response
- Some prompts were taken from early users of the OpenAI API, and labelers wrote the demo response
- All data is *closed source*

Table 1: Distribution of use case categories from our API prompt dataset.

Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

Instruction Fine-Tuning Dataset: Dolly

- Consists of 15k prompt/response pairs
- Employees of DataBricks wrote all the prompt/response pairs
- All data is *open source!*

We set up a contest, where the top 20 labelers would get a big award. We also outlined 7 very specific tasks:

- Open Q&A: For instance, "Why do people like comedy movies?" or "What is the capital of France?" In some cases, there's not a correct answer, and in others, it requires drawing on knowledge of the world at large.
- Closed Q&A: These are questions that can be answered using only the information contained in a passage of reference text. For instance, given a paragraph from Wikipedia on the atom, one might ask, "What is the ratio between protons and neutrons in the nucleus?"
- Extract information from Wikipedia: Here an annotator would copy a paragraph from Wikipedia and extract entities or other factual information such as weights or measurements from the passage.
- Summarize information from Wikipedia: For this, annotators provided a passage from Wikipedia and were asked to distill it to a short summary.
- Brainstorming: This task asked for open-ended ideation and an associated list of possible options. For instance, "What are some fun activities I can do with my friends this weekend?"
- Classification: For this task, annotators were asked to make judgments about class membership (e.g. are the items in a list animals, minerals or vegetables) or to judge the properties of a short passage of text, such as the sentiment of a movie review.
- Creative writing: This task would include things like writing a poem or a love letter.

Instruction Fine-Tuning Dataset: Dolly

Category	Instruction	Response
General Q&A	What is best Playstation or Xbox?	Playstation and Xbox are two competing game consoles, where which is "best" will depend on the preferences of the gamer. From a graphics and performance perspective they are comparable, with Sony's Playstation technically leading, though most gamers would be unlikely to notice a significant difference. From a game library perspective, both consoles have their respective first-party titles that appeal to different demographics. For example, Microsoft has seen great popularity due to popular title Halo being exclusively available on the console. On the other hand, Playstation has popular exclusive titles such as God of War and Gran Turismo. Ultimately the decision of which console is best will likely come down to which games titles the purchaser prefers.

Instruction Fine-Tuning Dataset: Dolly

Category	Instruction	Response
Open Q&A	What individual has won the most Olympic gold medals in the history of the games?	Michael Phelps has won the most gold medals of all time with 23 golds.

Instruction Fine-Tuning Dataset: Dolly

Category	Instruction	Context	Response
Closed Q&A	When was the first Reading railway station opened?	Reading railway station is a major transport hub in Reading, Berkshire, England. It is on the northern edge of the town centre, near the main retail and commercial areas and the River Thames, 36 miles (58 km) from London Paddington. The first Reading station was opened on 30 March 1840 as the temporary western terminus of the original line of the Great Western Railway (GWR). Reading is the ninth-busiest station in the UK outside London and the second busiest interchange station outside London with over 3.8 million passengers changing trains at the station annually.	The first Reading railway station was opened on the 30th of March, 1840.

Instruction Fine-Tuning Dataset: Dolly

Category	Instruction	Context	Response
Information Extraction	Extract all of the dates mentioned in this paragraph and list them using bullets in the format {Date} - {Description}	Daniel Leavitt was born November 16, 1813, at Rye, New Hampshire, the son of Benning Leavitt, influential businessman, state senator, county commissioner and later Chicopee Selectman, and his wife Olive (Jenness) Leavitt. Daniel Leavitt married in 1838 at West Springfield, Massachusetts, Ruth Jeannette Ball. They had three children.	<ul style="list-style-type: none">• November 16, 1813 - This is when Daniel Leavitt was born• 1838 - This is when Daniel Leavitt married his wife Ruth Jeannette Ball.

Instruction Fine-Tuning Dataset: Dolly

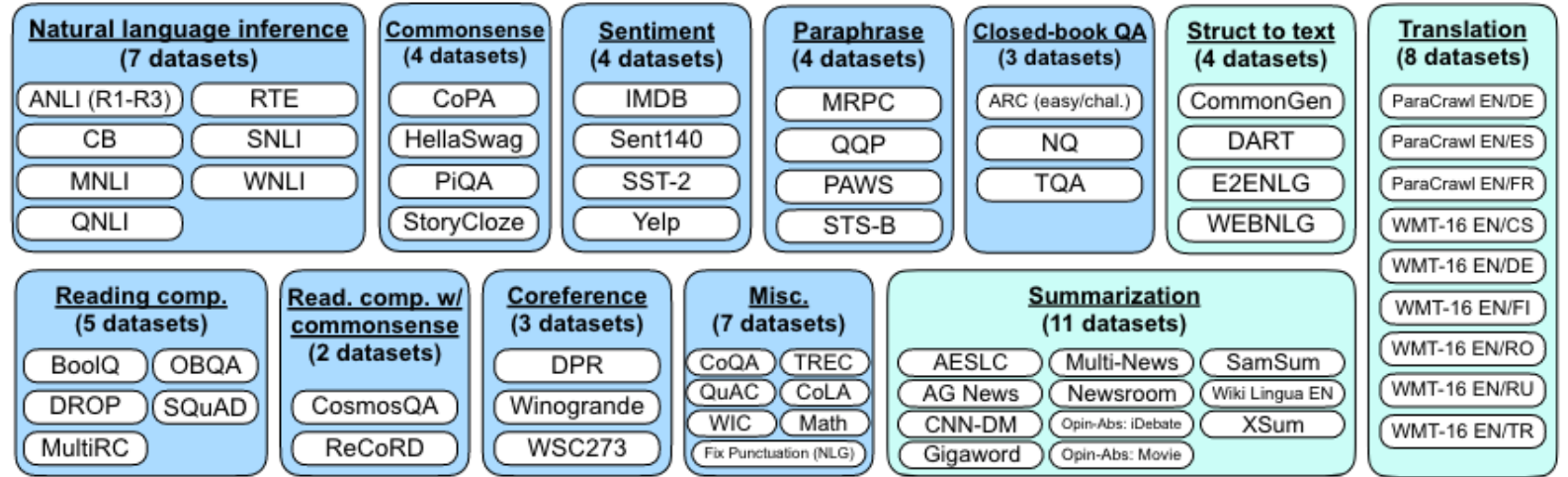
Category	Instruction	Response
Brainstorming	What are some unique curtain tie backs that you can make yourself?	There are many items you can use to make a curtain tie back. Some ideas for this include a chain, a circle and pin, jute rope, wooden beaded rope, a necklack or bracelet, a door knob, a length of a leather belt, macrame rope, or a string of artificial flowers.
Classification	Identify which instrument is string or percussion: Cantaro, Gudok	Gudok is string, Cantaro is percussion.

Instruction Fine-Tuning Dataset: Dolly

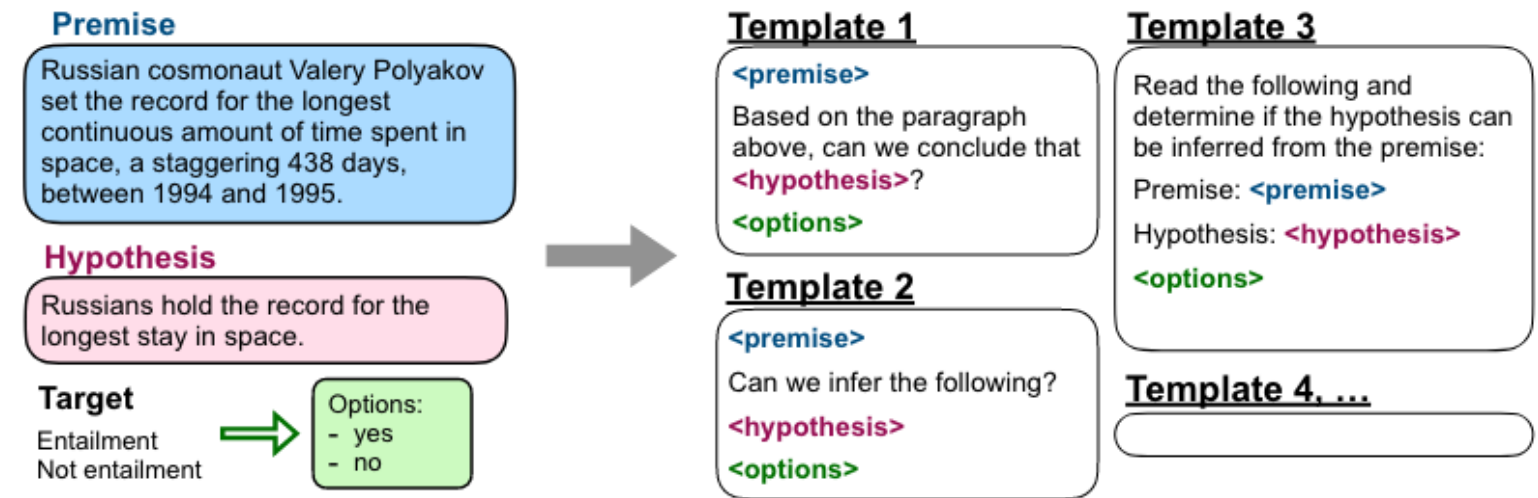
Category	Instruction	Response
Creative Writing	Write a haiku about how much I love pickles.	Caressed by sweet brine, you are crispy and crunchy. Oh, won't you be mine?
Creative Writing	What should I do on a free afternoon in San Francisco?	You can visit and take photos at San Francisco's landmarks such as Pier 39 or Golden Gate Bridge. If you enjoy walking outdoors, there are several large parks in the city. You can visit Golden Gate Park and the various attractions and museums located in the park. The Presidio is also a popular park where you can get some of the best views of Golden Gate Bridge.

Instruction Fine-Tuning Dataset: Flan

- **Key idea:** *programmatically* build instruction fine-tuning examples from existing NLP tasks (12 tasks, 62 datasets)



- For each NLP task/dataset, Flan created 10 templates e.g.,



Instruction Fine-Tuning Datasets

Release	Collection	Model Details				Data Collection & Training Details				
		Model	Base	Size	Public?	Prompt Types	Tasks in Flan	# Exs	Methods	
2020 05	UnifiedQA	UnifiedQA	RoBerta	110-340M	P	ZS	46 / 46	750k		
2021 04	CrossFit	BART-CrossFit	BART	140M	NP	FS	115 / 159	71.M		
2021 04	Natural Inst v1.0	Gen. BART	BART	140M	NP	ZS / FS	61 / 61	620k	+ Detailed k-shot Prompts	
2021 09	Flan 2021	Flan-LaMDA	LaMDA	137B	NP	ZS / FS	62 / 62	4.4M	+ Template Variety	
2021 10	P3	T0, T0+, T0++	T5-LM	3-11B	P	ZS	62 / 62	12M	+ Template Variety + Input Inversion	
2021 10	MetalCL	MetalCL	GPT-2	770M	P	FS	100 / 142	3.5M	+ Input Inversion + Noisy Channel Opt	
2021 11	ExMix	ExT5	T5	220M-11B	NP	ZS	72 / 107	500k	+ With Pretraining	
2022 04	Super-Natural Inst.	Tk-Instruct	T5-LM, mT5	11-13B	P	ZS / FS	1556 / 1613	5M	+ Detailed k-shot Prompts + Multilingual	
2022 10	GLM	GLM-130B	GLM	130B	P	FS	65 / 77	12M	+ With Pretraining + Bilingual (en, zh-cn)	
2022 11	xP3	BLOOMz, mT0	BLOOM, mT5	13-176B	P	ZS	53 / 71	81M	+ Massively Multilingual	
2022 12	Unnatural Inst. [†]	T5-LM-Unnat. Inst.	T5-LM	11B	NP	ZS	~20 / 117	64k	+ Synthetic Data	
2022 12	Self-Instruct [†]	GPT-3 Self Inst.	GPT-3	175B	NP	ZS	Unknown	82k	+ Synthetic Data + Knowledge Distillation	
2022 12	OPT-IML Bench [†]	OPT-IML	OPT	30-175B	P	ZS + FS CoT	~2067 / 2207	18M	+ Template Variety + Input Inversion + Multilingual	
2022 10	Flan 2022 (ours)	Flan-T5, Flan-PaLM	T5-LM, PaLM	10M-540B	P NP	ZS + FS CoT	1836	15M	+ Template Variety + Input Inversion + Multilingual	

Figure 2: A **Timeline of Public Instruction Tuning Collections** specifies the collection release date, detailed information on the finetuned models (the base model, their size, and whether the model itself is Public (P) or Not Public (NP)), what prompt specification they were trained for (zero-shot, few-shot, or Chain-of-Thought), the number of tasks contained in the Flan 2022 Collection (released with this work), and core methodological contributions in each work.

Note that the number of tasks and of examples vary under different assumptions and so are estimates. For instance, the definition of “task” and “task category” vary by work, and are not easily simplified to one ontology. The reported counts for the number of tasks are reported using task definitions from the respective works.

[†] indicates concurrent work.

perplexity (x)

$$= e^{-\frac{1}{T} \sum_{t=1}^T \log p(x_t | x_{t-1})}$$

Instruction Fine-Tuning Models

Instruction fine-tuned LLMs	# Params	Base Model	Fine-tuning Trainset		
			Self-build	Dataset Name	Size
Instruct-GPT (Ouyang et al., 2022)	176B	GPT-3 (Brown et al., 2020b)	Yes	-	-
BLOOMZ (Muennighoff et al., 2022) ¹	176B	BLOOM (Scao et al., 2022)	No	xP3	-
FLAN-T5 (Chung et al., 2022) ²	11B	T5 (Raffel et al., 2019)	No	FLAN 2021	-
Alpaca (Taori et al., 2023) ³	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Vicuna (Chiang et al., 2023) ⁴	13B	LLaMA (Touvron et al., 2023a)	Yes	-	70K
GPT-4-LLM (Peng et al., 2023) ⁵	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Claude (Bai et al., 2022b)	-	-	Yes	-	-
WizardLM (Xu et al., 2023a) ⁶	7B	LLaMA (Touvron et al., 2023a)	Yes	Evol-Instruct	70K
ChatGLM2 (Du et al., 2022) ⁷	6B	GLM (Du et al., 2022)	Yes	-	1.1 Tokens
LIMA (Zhou et al., 2023)	65B	LLaMA (Touvron et al., 2023a)	Yes	-	1K
OPT-IML (Iyer et al., 2022) ⁸	175B	OPT (Zhang et al., 2022a)	No	-	-
Dolly 2.0 (Conover et al., 2023) ⁹	12B	Pythia (Biderman et al., 2023)	No	-	15K
Falcon-Instruct (Almazrouei et al., 2023a) ¹⁰	40B	Falcon (Almazrouei et al., 2023b)	No	-	-
Guanaco (JosephusCheung, 2021) ¹¹	7B	LLaMA (Touvron et al., 2023a)	Yes	-	586K
Minotaur (Collective, 2023) ¹²	15B	StarCoder Plus (Li et al., 2023f)	No	-	-
Nous-Hermes (NousResearch, 2023) ¹³	13B	LLaMA (Touvron et al., 2023a)	No	-	300K+
TÜLU (Wang et al., 2023c) ¹⁴	6.7B	OPT (Zhang et al., 2022a)	No	Mixed	-
YuLan-Chat (YuLan-Chat-Team, 2023) ¹⁵	13B	LLaMA (Touvron et al., 2023a)	Yes	-	250K
MOSS (Tianxiang and Xipeng, 2023) ¹⁶	16B	-	Yes	-	-
Airoboros (Durbin, 2023) ¹⁷	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-
UltraLM (Ding et al., 2023a) ¹⁸	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-

- Key takeaway: instruction fine-tuned models are often very effective at a *much* smaller scale than typical LLMs

Idea: What if we could get an LLM to generate its own fine-tuning dataset? How would we make sure the responses were “good”?

Instruction fine-tuned LLMs	# Params	Base Model	Fine-tuning Trainset		
			Self-build	Dataset Name	Size
Instruct-GPT (Ouyang et al., 2022)	176B	GPT-3 (Brown et al., 2020b)	Yes	-	-
BLOOMZ (Muennighoff et al., 2022) ¹	176B	BLOOM (Scao et al., 2022)	No	xP3	-
FLAN-T5 (Chung et al., 2022) ²	11B	T5 (Raffel et al., 2019)	No	FLAN 2021	-
Alpaca (Taori et al., 2023) ³	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Vicuna (Chiang et al., 2023) ⁴	13B	LLaMA (Touvron et al., 2023a)	Yes	-	70K
GPT-4-LLM (Peng et al., 2023) ⁵	7B	LLaMA (Touvron et al., 2023a)	Yes	-	52K
Claude (Bai et al., 2022b)	-	-	Yes	-	-
WizardLM (Xu et al., 2023a) ⁶	7B	LLaMA (Touvron et al., 2023a)	Yes	Evol-Instruct	70K
ChatGLM2 (Du et al., 2022) ⁷	6B	GLM (Du et al., 2022)	Yes	-	1.1 Tokens
LIMA (Zhou et al., 2023)	65B	LLaMA (Touvron et al., 2023a)	Yes	-	1K
OPT-IML (Iyer et al., 2022) ⁸	175B	OPT (Zhang et al., 2022a)	No	-	-
Dolly 2.0 (Conover et al., 2023) ⁹	12B	Pythia (Biderman et al., 2023)	No	-	15K
Falcon-Instruct (Almazrouei et al., 2023a) ¹⁰	40B	Falcon (Almazrouei et al., 2023b)	No	-	-
Guanaco (JosephusCheung, 2021) ¹¹	7B	LLaMA (Touvron et al., 2023a)	Yes	-	586K
Minotaur (Collective, 2023) ¹²	15B	Starcoder Plus (Li et al., 2023f)	No	-	-
Nous-Hermes (NousResearch, 2023) ¹³	13B	LLaMA (Touvron et al., 2023a)	No	-	300K+
TÜLU (Wang et al., 2023c) ¹⁴	6.7B	OPT (Zhang et al., 2022a)	No	Mixed	-
YuLan-Chat (YuLan-Chat-Team, 2023) ¹⁵	13B	LLaMA (Touvron et al., 2023a)	Yes	-	250K
MOSS (Tianxiang and Xipeng, 2023) ¹⁶	16B	-	Yes	-	-
Airoboros (Durbin, 2023) ¹⁷	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-
UltraLM (Ding et al., 2023a) ¹⁸	13B	LLaMA (Touvron et al., 2023a)	Yes	-	-

- Key takeaway: instruction fine-tuned models are often very effective at a *much* smaller scale than typical LLMs

Reinforcement Learning from Human Feedback

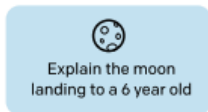
- **InstructGPT** uses Reinforcement Learning with Human Feedback (RLHF) to *further fine-tune an already instruction fine-tuned version of GPT3*
 - “In human evaluations on our prompt distribution, outputs from the 1.3B parameter InstructGPT model are preferred to outputs from the 175B GPT-3, despite having 100x fewer parameters.”

Reinforcement Learning from Human Feedback

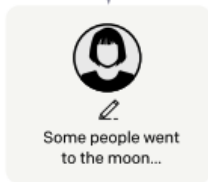
Step 1

Collect demonstration data, and train a supervised policy.

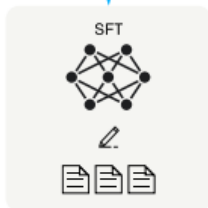
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



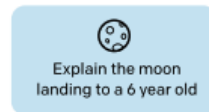
This data is used to fine-tune GPT-3 with supervised learning.



Step 2

Collect comparison data, and train a reward model.

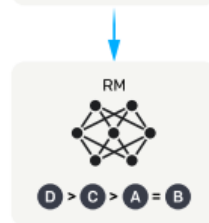
A prompt and several model outputs are sampled.



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.



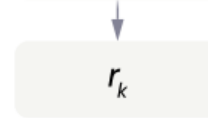
The policy generates an output.



The reward model calculates a reward for the output.



The reward is used to update the policy using PPO.

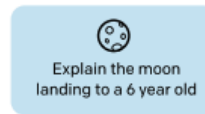


Reinforcement Learning from Human Feedback

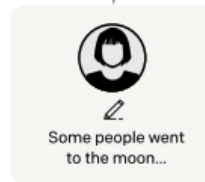
Step 1

Collect demonstration data, and train a supervised policy.

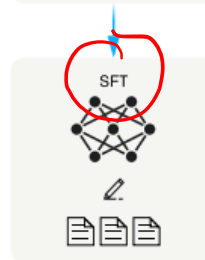
A prompt is sampled from our prompt dataset.



A labeler demonstrates the desired output behavior.



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



- Step 1 performs instruction fine-tuning on 13k training examples
- This aligns the model behavior with what we would expect of a chat agent but the *diversity* of the responses might be limited because of the small (and potentially not representative) training data

Reinforcement Learning from Human Feedback

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

- In Step 2, the model takes 33k prompts and samples $K \in \{4, \dots, 9\}$ responses from the instruction fine-tuned model for each one
- Then, a human labeler ranks all the responses

Reinforcement Learning from Human Feedback: Ranking Interface

Playground task 

Talk to the assistant

Task 1 of 5

Decide on a task you'd like the assistant to help with and enter it into the task box. Interact with the AI assistant. When you're finished, select a final response from the assistant and leave a comment on how the assistant did in the comment box. Click to go to the next task.

Comments

I thought the assistant was ...

Rating

Bad Good

Next Task



Human

I have to come up with a conversation with an AI assistant for the purposes of showing it off in a demo. What sorts of conversation topics would you recommend?



Assistant

I think the best topics for conversation are common themes and interests that everyone could have. Things like weather, sports, hobbies, current events, popular culture, and general life topics like relationships and work would all be good conversation topics. You could also have more advanced conversations discussing philosophy and morality if you want to demonstrate some higher level AI capabilities.



Human

I think I'm more interested in an advanced conversation. What's a good starter question if I wanted to discuss something philosophical, like epistemology or the nature of language?



Assistant

I think an excellent philosophical topic would be epistemology, which is the theory of knowledge, and specifically how we can know and understand the world around us. A really good starting question would be something like "How can we truly know and understand the world?", or "What is the nature of knowledge?". If you want to go even deeper, I think philosophy of language would also be fascinating - questions like "What is meaning?", "How do language and thought relate?", and "How can language mislead and deceive?" would be great discussion topics.



Human

How would you answer a question like: How do language and thought relate?



Choose the most helpful and honest response

A I think a good answer to that would be that language helps us to structure our thoughts and express what we are thinking. Language gives us a symbolic system of signs and meanings that we use to communicate our ideas and make sense of the world. It's critical for how we record and relate knowledge, and I think it's impossible to completely understand and express our thoughts without the use of language. Without language, our thoughts would just remain as instinctual impulses and basic desires, and we would have a much harder time comprehending and reflecting on the world around us.

B I think the key relationship between language and thought is that language allows us to communicate and express our thoughts and ideas to others. Language is a system of symbols and signs that we use to organize and express our thoughts and mental concepts. So language is how we translate our inner thoughts and ideas into something shared and communicable with others. Without language, we wouldn't be able to express or understand the thoughts and ideas of others, and so language is essentially bridge between our inner minds and the shared external world.

A

A

A

A

B

B

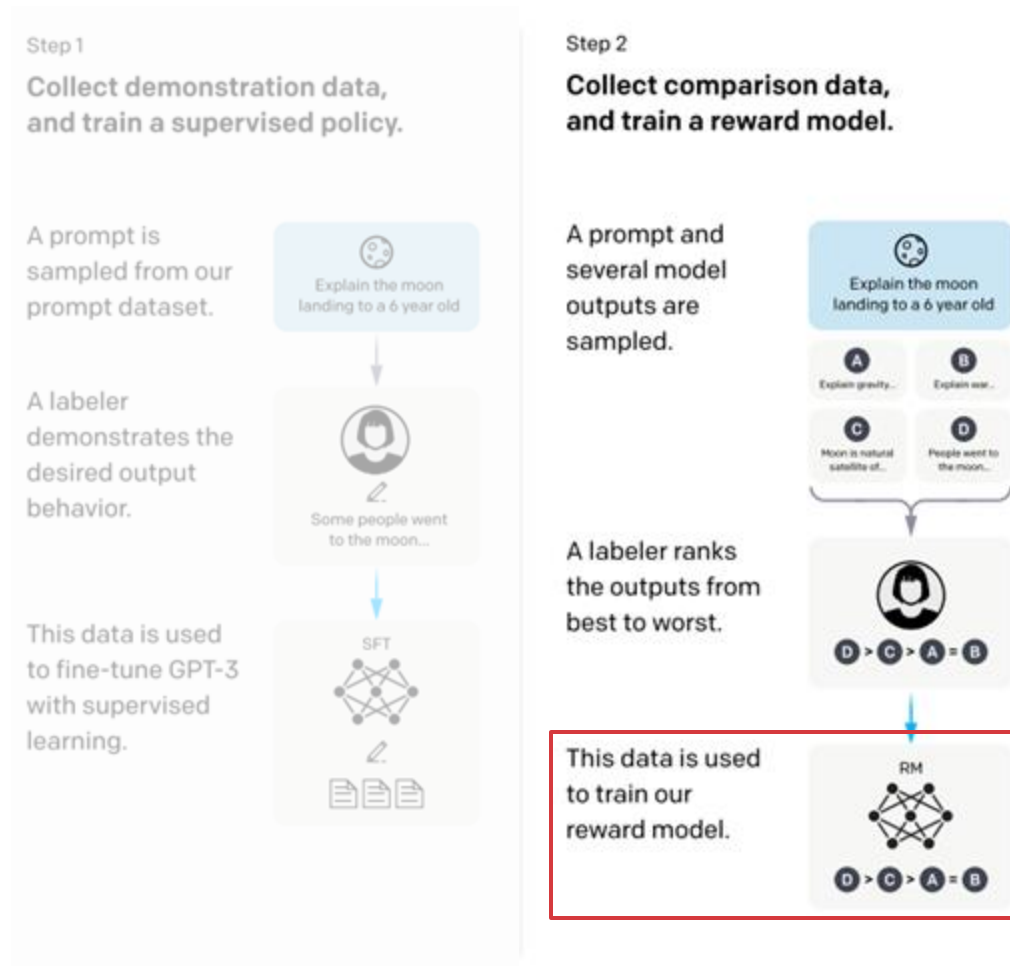
B

B

A is better

B is better

Reinforcement Learning from Human Feedback: Reward Model



- In Step 2, the model takes 33k prompts and samples $K \in \{4, \dots, 9\}$ responses from the instruction fine-tuned model for each one
- Then, a human labeler ranks all the responses

Reinforcement Learning from Human Feedback: Reward Model

- The reward model is a smaller LLM (e.g., 6B vs. 175B parameters), but with the softmax over words replaced with a linear layer so that it outputs a scalar value, i.e. the reward of some (prompt, response) pair
- This *regression* model is trained so that rewards of higher-ranking responses, y_w , are larger than those of the lower-ranking responses, y_l :

$$\ell(\theta) = \frac{1}{\binom{K}{2}} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\theta(x, y_w) - r_\theta(x, y_l))]$$

where x is the prompt, K is the number of responses for prompt x , \mathcal{D} is the set of human labelled preferences, and r_θ is the reward model

Reinforcement Learning from Human Feedback: Reward Model

- The reward model is a smaller LLM (e.g., 6B vs. 175B parameters), but with the softmax over words replaced with a linear layer so that it outputs a scalar value, i.e. the reward of some (prompt, response) pair
- This *regression* model is trained so that rewards of higher-ranking responses, y_w , are larger than those of the lower-ranking responses, y_l :

$$\ell(\theta) = -\frac{1}{\binom{K}{2}} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\theta(x, y_w) - r_\theta(x, y_l))]$$

- During training, all $\binom{K}{2}$ pairwise rankings for a prompt are kept together in a batch for efficiency/stability

Okay so now
what do we do
with this
thing...?

- The reward model is a smaller LLM (e.g., 6B vs. 175B parameters), but with the softmax over words replaced with a linear layer so that it outputs a scalar value, i.e. the reward of some (prompt, response) pair
- This *regression* model is trained so that rewards of higher-ranking responses, y_w , are larger than those of the lower-ranking responses, y_l :

$$\ell(\theta) = -\frac{1}{\binom{K}{2}} \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} [\log \sigma(r_\theta(x, y_w) - r_\theta(x, y_l))]$$

- During training, all $\binom{K}{2}$ pairwise rankings for a prompt are kept together in a batch for efficiency/stability

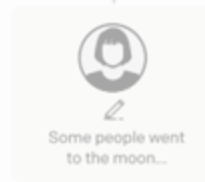
Okay so now
what do we do
with this
thing...?

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.



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.

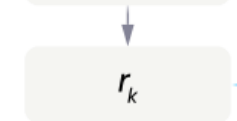


Once upon a time...

The reward model
calculates a
reward for
the output.



The reward is
used to update
the policy
using PPO.



Reinforcement Learning: Problem Formulation

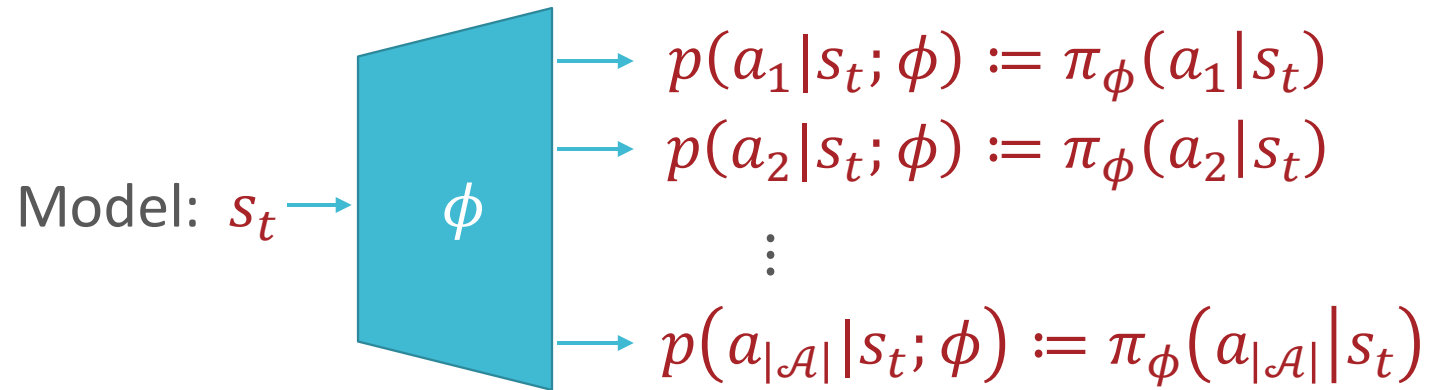
- State space, \mathcal{S}
- Action space, \mathcal{A}
- Reward function
 - Stochastic, $p(r | s, a)$
 - Deterministic, $R: \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$
- Transition function
 - Stochastic, $p(s' | s, a)$
 - Deterministic, $\delta: \mathcal{S} \times \mathcal{A} \rightarrow \mathcal{S}$

Reinforcement Learning: Problem Formulation for Fine-tuning LLMs

- State space, $\mathcal{S} = \{\text{all possible sequences of tokens}\}$
- Action space, $\mathcal{A} = \{\text{vocabulary of next tokens}\}$
- Reward function
 - Stochastic, $p(r \mid s, a)$
 - Deterministic reward based on reward model trained on human feedback, R_θ
 - R_θ is a bit of weird reward function from an RL perspective: it returns $0 \forall a \neq \text{EOS}$ and $r_\theta(x, [s, a] - x)$ otherwise
- Transition function
 - Stochastic, $p(s' \mid s, a)$
 - Deterministic, $\delta(s, a) = [s, a]$

Reinforcement Learning: Object of Interest

- A stochastic, parametrized policy, $\pi_\phi(a | s)$
 - Specifies a distribution over actions in *every* state

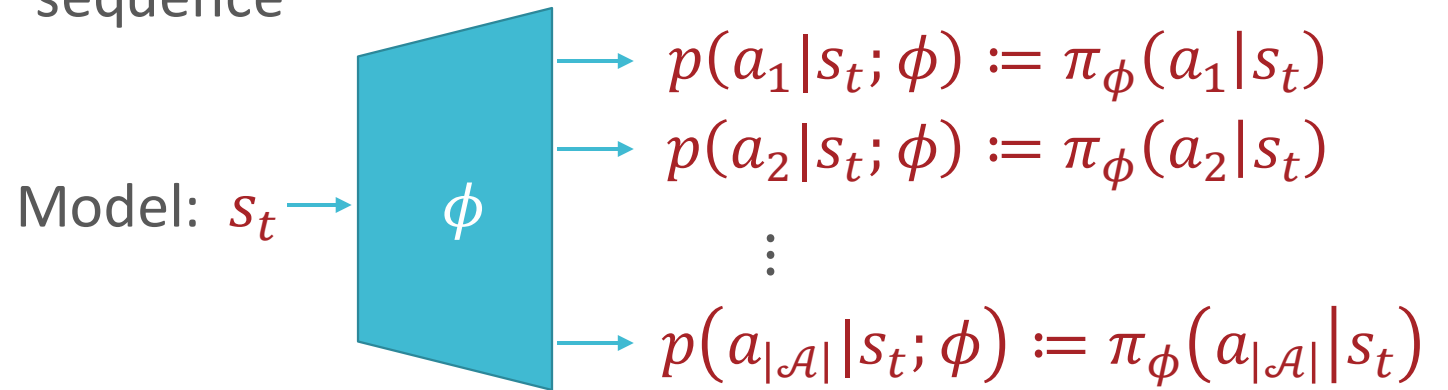


- An *episode* $\mathbf{T} = \{s_0, a_0, s_1, a_1, \dots, s_T\}$ is one run of an agent through the system ending in a terminal state, s_T
- The stochastic policy induces a distribution over episodes

$$\begin{aligned} p_\phi(\mathbf{T}) &= p(\{s_0, a_0, s_1, a_1, \dots, s_T\}) \\ &= p(s_0) \prod_{t=0}^{T-1} p(s_{t+1} | s_t, a_t) \pi_\phi(a_t | s_t) \end{aligned}$$

Reinforcement Learning: Object of Interest for Fine-tuning LLMs

- The LLM to be fine-tuned, $\pi_\phi(a | s)$
 - Specifies a distribution over next tokens given any input sequence



- An *episode* $\mathbf{T} = \{x, a_0, s_1, a_1, \dots, s_T\}$ is one completion of the prompt x , ending in an EOS token
- The LLM induces a distribution over possible completions

$$\begin{aligned} p_\phi(\mathbf{T}) &= p(\{a_0, s_1, a_1, \dots, s_T\} | x := s_0) \\ &= \prod_{t=0}^{T-1} \pi_\phi(a_t | s_t) \end{aligned}$$

Objective function: $\ell(\phi) = -\mathbb{E}_{p_\phi(\mathbf{T})}[R_\theta(\mathbf{T})]$, the negative expected reward of a response

$$\begin{aligned}\nabla_\phi \ell(\phi) &= \nabla_\phi \left(-\mathbb{E}_{p_\phi(\mathbf{T})}[R_\theta(\mathbf{T})] \right) = \nabla_\phi \left(-\int R_\theta(\mathbf{T}) p_\phi(\mathbf{T}) d\mathbf{T} \right) \\ &= -\int R_\theta(\mathbf{T}) \nabla_\phi \left(\prod_{t=0}^{T-1} \pi_\phi(a_t | s_t) \right) d\mathbf{T}\end{aligned}$$

- Issue: $\nabla_\phi p_\phi(\mathbf{T})$ involves taking the gradient of a (hideous) product

Likelihood
Ratio
Method
a.k.a.
REINFORCE
(Williams,
1992)

Objective function: $\ell(\phi) = -\mathbb{E}_{p_\phi(\mathcal{T})}[R_\theta(\mathcal{T})]$, the negative expected reward of a response

$$\begin{aligned}\nabla_\phi \ell(\phi) &= \nabla_\phi \left(-\mathbb{E}_{p_\phi(\mathcal{T})}[R_\theta(\mathcal{T})] \right) = \nabla_\phi \left(-\int R_\theta(\mathcal{T}) p_\phi(\mathcal{T}) d\mathcal{T} \right) \\ &= -\int R_\theta(\mathcal{T}) \nabla_\phi \left(\prod_{t=0}^{T-1} \pi_\phi(a_t | s_t) \right) d\mathcal{T}\end{aligned}$$

• Insight:

$$\nabla_\phi P_\phi(\mathcal{T}) = \frac{P_\phi(\mathcal{T})}{P_\phi(\mathcal{T})} \nabla_\phi P_\phi(\mathcal{T}) = P_\phi(\mathcal{T}) \nabla_\phi \log P_\phi(\mathcal{T})$$

$$\log P_\phi(\mathcal{T}) = \sum_{t=0}^{T-1} \log \pi_\phi(a_t | s_t)$$

$$\nabla_\phi \log P_\phi(\mathcal{T}) = \sum_{t=0}^{T-1} \nabla_\phi \log \pi_\phi(a_t | s_t)$$

Objective function: $\ell(\phi) = -\mathbb{E}_{p_\phi(T)}[R_\theta(T)]$, the negative expected reward of a response

$$\nabla_\phi \ell(\phi) = \nabla_\phi \left(-\mathbb{E}_{p_\phi(T)}[R_\theta(T)] \right) = \nabla_\phi \left(-\int R_\theta(T) p_\phi(T) dT \right)$$

$$= -\int R_\theta(T) \underbrace{\nabla_\phi p_\phi(T)}_{\text{chain rule}} dT = -\int R_\theta(T) \underbrace{\nabla_\phi (\log p_\phi(T))}_{\text{chain rule}} p_\phi(T) dT$$

$$= -\mathbb{E}_{p_\phi(T)} [R_\theta(T) \nabla_\phi \log p_\phi(T)]$$

$$\approx -\frac{1}{N} \sum_{n=1}^N R_\theta(T^{(n)}) \nabla_\phi \log p_\phi(T^{(n)})$$

where $T^{(n)} \sim p_\phi$ (LLM decoding)

$$= -\frac{1}{N} \sum_{n=1}^N r_\theta(x, [s_T^{(n)}, a_T^{(n)}] - x) \left(\sum_{t=0}^{T-1} \nabla_\phi \log \pi_\phi(a_t^{(n)} | s_t^{(n)}) \right)$$

Likelihood
Ratio
Method
a.k.a.
REINFORCE
(Williams,
1992)