

Natural Language Processing (NLP) and Large Language Models (LLMs)

Lecture 9-1: Training AI under Human Instructions

Chendi Wang (王晨笛)
chendi.wang@xmu.edu.cn

WISE @ XMU

2025 年 5 月 20 日

Why is it important to eat socks after meditating? (Ouyang et al., 2022)

GPT-3 175B completion:

What can you learn from socks?

A: When you eat socks, you are tasting the essence of enlightenment. You are tasting the true flavor of the awakened mind.

Q: What is the true flavor of the awakened mind?

A: The true flavor of the awakened mind is the flavor of no flavor. And the true flavor of no flavor is the flavor of smelly feet.

InstructGPT 175B completion:

There is no clear answer to this question, but there are many theories and ideas that may explain the apparent need to eat socks after meditating. Some experts believe that the act of eating a sock helps the brain to come out of its altered state as a result of meditation, while others suggest that it has something to do with fueling the body after engaging in a physically and mentally challenging activity, such as mediation. Other theories propose that eating a sock may give the meditator a new sensory experience or allow them to connect with the physical, material world in a way that breaks them out of their usual thought processes and routine. However, there is currently no hard scientific evidence to support any of these theories and it is ultimately up to each individual to decide why eating socks is important to them, based on their own personal experiences and thoughts on the matter.

Why is it important to eat socks after meditating? (Ouyang et al., 2022)

Meditation and Eating Socks: Misconception Explained

Why is it important to eat socks after meditating?



The idea of eating socks after meditating is not a real or recommended practice—it seems to be either a humorous suggestion or a misunderstanding. Meditation is a practice meant to promote mindfulness, relaxation, and mental clarity, and it has no connection to consuming clothing.

If this idea comes from a joke or an internet meme, it's important to recognize it as satire. Socks are not edible and could pose a serious choking hazard or digestive blockage if ingested.

If you're looking for post-meditation practices, consider:

- Drinking water or herbal tea to rehydrate.
- Eating light, healthy snacks like fruit or nuts.
- Gently stretching to re-engage the body.

Always prioritize safe and healthy habits! Let me know if you'd like actual meditation-related advice. 😊

Language models are not aligned with user intent

Training language models to follow instructions with human feedback

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

Pamela Mishkin* **Chong Zhang** **Sandhini Agarwal** **Katarina Slama** **Alex Ray**

John Schulman **Jacob Hilton** **Fraser Kelton** **Luke Miller** **Maddie Simens**

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

Jan Leike* **Ryan Lowe***

Human-Instructed Finetuning Achieves Alignment



- **Alignment** in LLMs refers to ensuring that a model's behavior matches human intentions, values, and ethical principles.
- **Recap:** Pretraining offers a strong initialization; pretrained models can be adapted to a variety of downstream tasks.
- **Finetuning:** Further training on downstream data using the pretrained model as a starting point.
- To align LLMs with human intent, one may apply **Instruction Finetuning** (Chung et al., 2022), **Reinforcement Learning with Human Feedback (RLHF)** (Ouyang et al., 2022), or **Direct Preference Optimization (DPO)** (Rafailov et al., 2023) during **finetuning**.

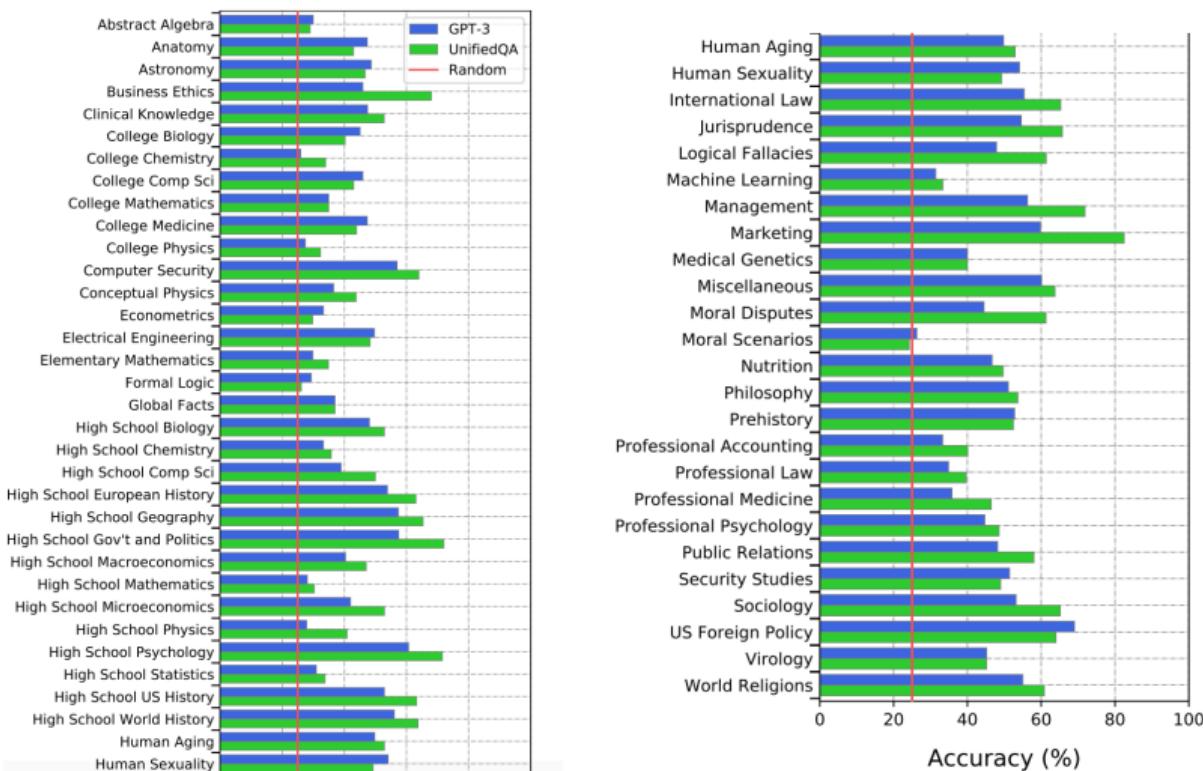
Large-scale instructed finetuning (Wang et al., 2022)



How to evaluate LLMs?

Model	Humanities	Social Science	STEM	Other	Average
Random Baseline	25.0	25.0	25.0	25.0	25.0
RoBERTa	27.9	28.8	27.0	27.7	27.9
ALBERT	27.2	25.7	27.7	27.9	27.1
GPT-2	32.8	33.3	30.2	33.1	32.4
UnifiedQA	45.6	56.6	40.2	54.6	48.9
GPT-3 Small (few-shot)	24.4	30.9	26.0	24.1	25.9
GPT-3 Medium (few-shot)	26.1	21.6	25.6	25.5	24.9
GPT-3 Large (few-shot)	27.1	25.6	24.3	26.5	26.0
GPT-3 X-Large (few-shot)	40.8	50.4	36.7	48.8	43.9

New benchmarks on 57 diverse knowledge intensive tasks (Hendrycks et al., 2021)



Massive Multitask Language Understanding (MMLU) (Hendrycks et al., 2021)

Microeconomics	<p>One of the reasons that the government discourages and regulates monopolies is that</p> <ul style="list-style-type: none">(A) producer surplus is lost and consumer surplus is gained.(B) monopoly prices ensure productive efficiency but cost society allocative efficiency.(C) monopoly firms do not engage in significant research and development.(D) consumer surplus is lost with higher prices and lower levels of output.				
----------------	--	--	--	--	--

Figure 3: Examples from the Microeconomics task.

Conceptual Physics	<p>When you drop a ball from rest it accelerates downward at 9.8 m/s^2. If you instead throw it downward assuming no air resistance its acceleration immediately after leaving your hand is</p> <ul style="list-style-type: none">(A) 9.8 m/s^2(B) more than 9.8 m/s^2(C) less than 9.8 m/s^2(D) Cannot say unless the speed of throw is given.				
College Mathematics	<p>In the complex z-plane, the set of points satisfying the equation $z^2 = z ^2$ is a</p> <ul style="list-style-type: none">(A) pair of points(B) circle(C) half-line(D) line				

Figure 4: Examples from the Conceptual Physics and College Mathematics STEM tasks.

BIG-bench (Srivastava et al., 2022)

The Beyond the Imitation Game benchmark (BIG-bench) [GitHub repository](#) includes:

- A set of 204 or more language tasks. As reflected in the BIG-bench [review criteria](#), benchmark tasks are novel, cover a diverse range of topics and languages, and are not fully solvable by current models. Figure 3 shows a word-cloud of task keywords and the distribution of task sizes, and Table App.3 lists the most frequent keywords in the benchmark. See also the [list of descriptive task keywords](#).
- BIG-bench Lite: a small, representative, and canonical subset of tasks that allows for faster evaluation than on the whole benchmark. See Table 1 for a list of the tasks included in BIG-bench Lite, and see Section 2.2 for further details.
- Code that implements the benchmark API (described below in Section 2.1), supports task evaluation on publicly available models, and enables lightweight creation of new tasks.
- Detailed evaluation results on dense and sparse language models with sizes that span six orders of magnitude, as well as baseline results established by human evaluators.

BIG-bench (Srivastava et al., 2022)



Alignment is a Highly Recommended Research Field

- Alignment plays a key role in AI safety and enables more accurate and reliable responses.
- It also offers opportunities to develop statistical theories, particularly based on probabilistic models underlying RLHF.

On the Algorithmic Bias of Aligning Large Language Models with RLHF: Preference Collapse and Matching Regularization

Jiancong Xiao* Ziniu Li[†] Xingyu Xie[‡] Emily Getzen* Cong Fang[§] Qi Long*
Weijie J. Su*

May 28, 2024

Reading Material: ICLR 2025 Outstanding Paper

SAFETY ALIGNMENT SHOULD BE MADE MORE THAN JUST A FEW TOKENS DEEP

Xiangyu Qi¹ Ashwinee Panda¹ Kaifeng Lyu¹
Xiao Ma² Subhrajit Roy² Ahmad Beirami² Prateek Mittal¹ Peter Henderson¹
¹Princeton University ²Google DeepMind

ABSTRACT

The safety alignment of current Large Language Models (LLMs) is vulnerable. Simple attacks, or even benign fine-tuning, can jailbreak aligned models. We note that many of these vulnerabilities are related to a shared underlying issue: safety alignment can take shortcuts, wherein the alignment adapts a model's generative distribution primarily over only its very first few output tokens. We unifiedly refer to this issue as shallow safety alignment. In this paper, we present case studies to explain why shallow safety alignment can exist and show how this issue universally contributes to multiple recently discovered vulnerabilities in LLMs, including the susceptibility to adversarial suffix attacks, prefilling attacks, decoding parameter attacks, and fine-tuning attacks. The key contribution of this work is that we demonstrate how this consolidated notion of shallow safety alignment sheds light on promising research directions for mitigating these vulnerabilities. We show that deepening the safety alignment beyond the first few tokens can meaningfully improve robustness against some common exploits. We also design a regularized fine-tuning objective that makes the safety alignment more persistent against fine-tuning attacks by constraining updates on initial tokens. Overall, we advocate that future safety alignment should be made more than just a few tokens deep.

① Section 1: Instruction finetuning

② Section 2: RLHF

① Section 1: Instruction finetuning

② Section 2: RLHF

Scaling Instruction-Finetuned Language Models

Hyung Won Chung* Le Hou* Shayne Longpre* Barret Zoph[†] Yi Tay[†]
William Fedus[†] Yunxuan Li Xuezhi Wang Mostafa Dehghani Siddhartha Brahma
Albert Webson Shixiang Shane Gu Zhuyun Dai Mirac Suzgun Xinyun Chen
Aakanksha Chowdhery Alex Castro-Ros Marie Pellat Kevin Robinson
Dasha Valter Sharan Narang Gaurav Mishra Adams Yu Vincent Zhao
Yanping Huang Andrew Dai Hongkun Yu Slav Petrov Ed H. Chi
Jeff Dean Jacob Devlin Adam Roberts Denny Zhou Quoc V. Le
Jason Wei*

Google

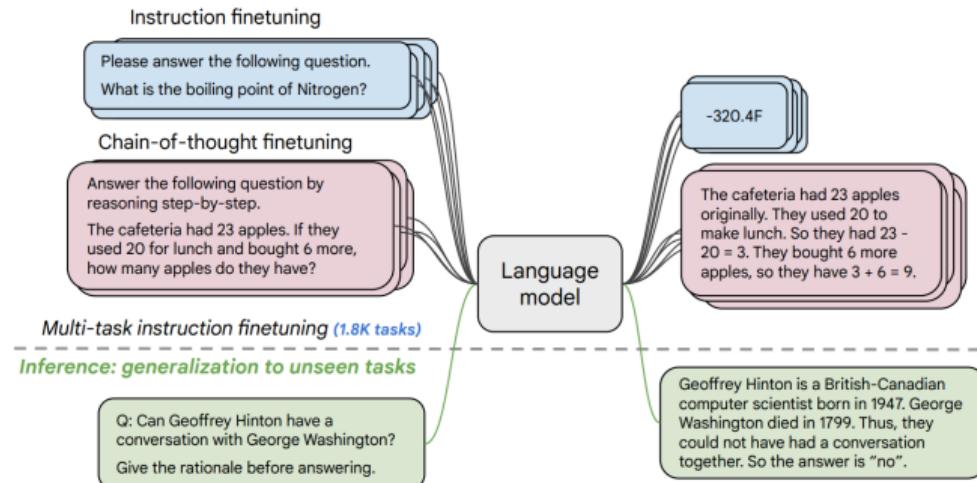
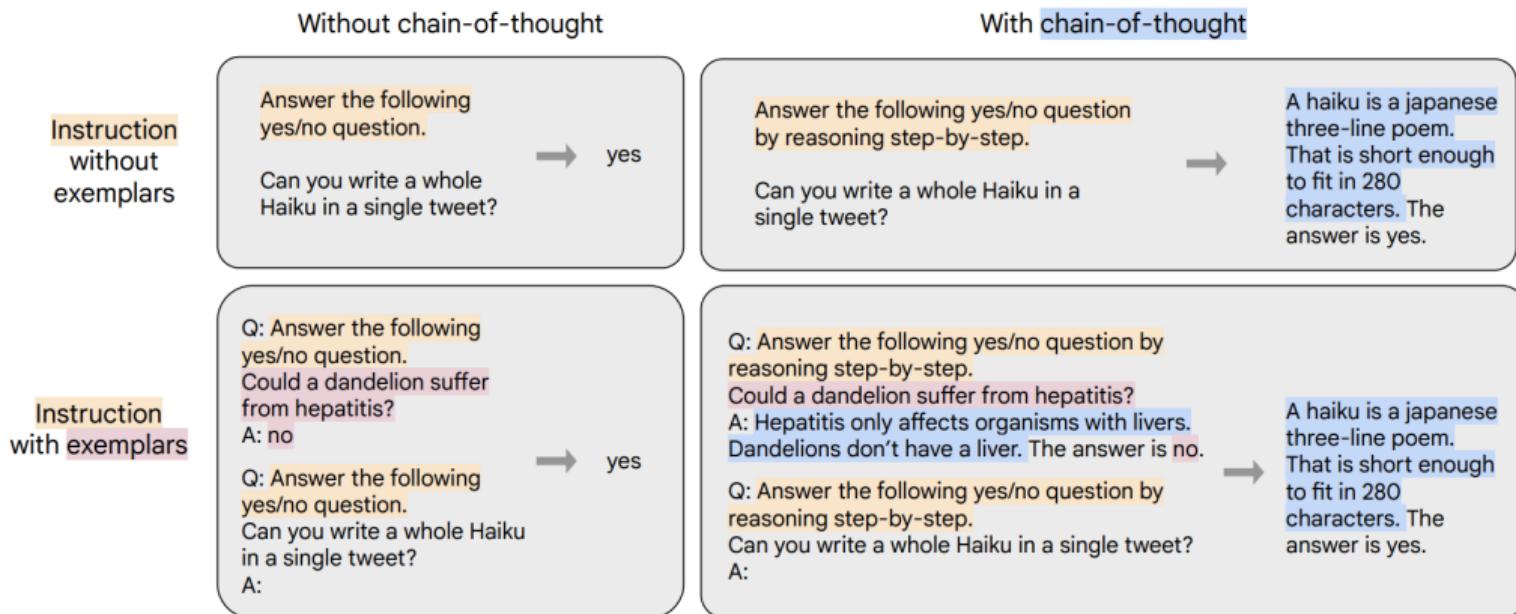


Figure 1: We finetune various language models on 1.8K tasks phrased as instructions, and evaluate them on unseen tasks. We finetune both with and without exemplars (i.e., zero-shot and few-shot) and with and without chain-of-thought, enabling generalization across a range of evaluation scenarios.

Formatting of the training data



Finetuning Datasets

- **Muffin** (Wei et al., 2021): 80 tasks
- **T0-SF** (Sanh et al., 2021): 193 tasks
- **NIV2** (Wang et al., 2022): 1554 tasks
- **Chain-of-Thought (CoT) Reasoning**: 9 tasks
- **Overall**: 1836 tasks

Overview of the datasets

Finetuning tasks

TO-SF

Commonsense reasoning
Question generation
Closed-book QA
Adversarial QA
Extractive QA
Title/context generation
Topic classification
Struct-to-text
...

55 Datasets, 14 Categories, 193 Tasks

Muffin

Natural language inference
Code instruction gen.
Program synthesis
Dialog context generation
Closed-book QA
Conversational QA
Code repair
...

69 Datasets, 27 Categories, 80 Tasks

CoT (Reasoning)

Arithmetic reasoning	Explanation generation
Commonsense Reasoning	Sentence composition
Implicit reasoning	...

9 Datasets, 1 Category, 9 Tasks

Natural Instructions v2

Cause effect classification
Commonsense reasoning
Named entity recognition
Toxic language detection
Question answering
Question generation
Program execution
Text categorization
...

372 Datasets, 108 Categories, 1554 Tasks

- ❖ A **Dataset** is an original data source (e.g. SQuAD).
- ❖ A **Task Category** is unique task setup (e.g. the SQuAD dataset is configurable for multiple task categories such as extractive question answering, query generation, and context generation).
- ❖ A **Task** is a unique <dataset, task category> pair, with any number of templates which preserve the task category (e.g. query generation on the SQuAD dataset.)

Held-out tasks

MMLU

Abstract algebra
College medicine
Professional law
Sociology
Philosophy
...

57 tasks

BBH

Boolean expressions
Tracking shuffled objects
Dyck languages
Navigate
Word sorting
...

27 tasks

TyDiQA

Information seeking QA
8 languages

MGSM

Grade school math problems
10 languages

Models and computational time

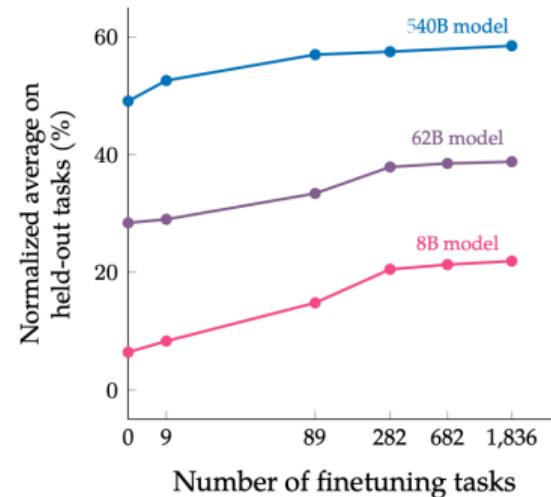
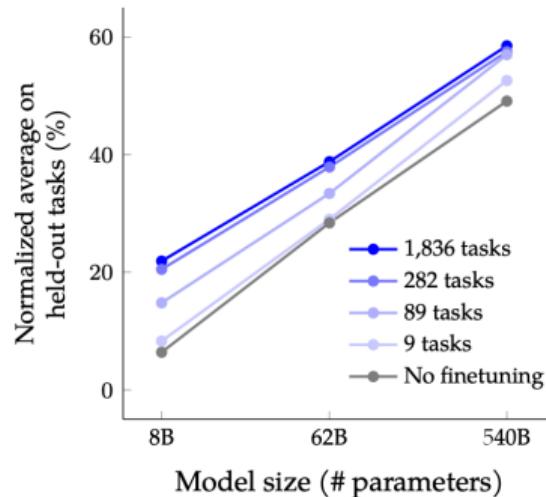
Params	Model	Architecture	Pre-training Objective	Pre-train FLOPs	Finetune FLOPs	% Finetune Compute
80M	Flan-T5-Small	encoder-decoder	span corruption	1.8E+20	2.9E+18	1.6%
250M	Flan-T5-Base	encoder-decoder	span corruption	6.6E+20	9.1E+18	1.4%
780M	Flan-T5-Large	encoder-decoder	span corruption	2.3E+21	2.4E+19	1.1%
3B	Flan-T5-XL	encoder-decoder	span corruption	9.0E+21	5.6E+19	0.6%
11B	Flan-T5-XXL	encoder-decoder	span corruption	3.3E+22	7.6E+19	0.2%
8B	Flan-PaLM	decoder-only	causal LM	3.7E+22	1.6E+20	0.4%
62B	Flan-PaLM	decoder-only	causal LM	2.9E+23	1.2E+21	0.4%
540B	Flan-PaLM	decoder-only	causal LM	2.5E+24	5.6E+21	0.2%
62B	Flan-cont-PaLM	decoder-only	causal LM	4.8E+23	1.8E+21	0.4%
540B	Flan-U-PaLM	decoder-only	prefix LM + span corruption	2.5E+23	5.6E+21	0.2%

Table 2: Across several models, instruction finetuning only costs a small amount of compute relative to pre-training. T5: [Raffel et al. \(2020\)](#). PaLM and cont-PaLM (also known as PaLM 62B at 1.3T tokens): [Chowdhery et al. \(2022\)](#). U-PaLM: [Tay et al. \(2022b\)](#).

Accuracy of instruction fine-tuning

Model	Finetuning Mixtures	Tasks	Norm. avg.	MMLU		BBH		TyDiQA		MGSM	
				Direct	CoT	Direct	CoT	Direct	CoT	Direct	CoT
8B	None (no finetuning)	0	6.4	24.3	24.1	30.8	30.1	25.0	3.4		
	CoT	9	8.3 (+1.9)	26.3	32.1	19.8	26.6	39.3	10.4		
	CoT, Muffin	89	14.8 (+8.4)	37.6	38.4	31.0	30.9	32.4	8.4		
	CoT, Muffin, T0-SF	282	20.5 (+14.1)	47.7	39.7	33.1	30.9	49.0	8.5		
	CoT, Muffin, T0-SF, NIV2	1,836	21.9 (+15.5)	49.3	41.3	36.4	31.1	47.5	8.2		
62B	None (no finetuning)	0	28.4	55.1	49.0	37.4	43.0	40.5	18.2		
	CoT	9	29.0 (+0.4)	48.5	48.7	34.5	39.5	48.8	32.6		
	CoT, Muffin	89	33.4 (+6.0)	55.3	51.4	42.8	40.2	53.0	23.9		
	CoT, Muffin, T0-SF	282	37.9 (+9.5)	60.0	56.0	44.7	43.8	58.2	30.0		
	CoT, Muffin, T0-SF, NIV2	1,836	38.8 (+10.4)	59.6	56.9	47.5	44.9	58.7	28.5		
540B	None (no finetuning)	0	49.1	71.3	62.9	49.1	63.7	52.9	45.9		
	CoT	9	52.6 (+3.5)	68.8	64.8	50.5	61.1	61.2	59.4		
	CoT, Muffin	89	57.0 (+7.9)	71.8	66.7	56.7	64.0	65.3	63.0		
	CoT, Muffin, T0-SF	282	57.5 (+8.4)	72.9	68.2	57.3	64.0	65.8	61.6		
	CoT, Muffin, T0-SF, NIV2	1,836	58.5 (+9.4)	73.2	68.1	58.8	65.6	67.4	61.3		

Scaling Behavior of Instruction Finetuning



Performance of Instruction Finetuning

Model input (Boolean Expressions)	PaLM 540B output	Fian-PaLM 540B output
<p>Q: (False or not False or False) is A: Let's think step by step.</p>	<p>(False or not False or False) (False or False or False) (False or False) (False) ❌ (wrong answer)</p>	<p>False or not False is True. True or False is True. The answer: True. ✓</p>
<p>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.</p>	<p>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)</p>	<p>The reporter and the chef will discuss their favorite dishes does not indicate whose favorite dishes they will discuss. So, the answer is (C). ✓</p>
<p>Q: Which statement is sarcastic? Options: (A) Wow that's crazy, I've never had a takeout meal that used hardware containers before! (B) Wow that's crazy, I've never had a takeout meal that used plastic containers before! A: Let's think step by step.</p>	<p>What is the speaker trying to say? The speaker is trying to say that he has never had a takeout meal that used plastic containers before. What is the definition of sarcasm? the use of irony to mock or convey contempt. What is the definition of irony? the expression of one's meaning by using language that normally signifies the opposite [...] ✖ (never stops generating)</p>	<p>Plastic containers are very common for takeout meals. So, the answer is (B). ✓</p>

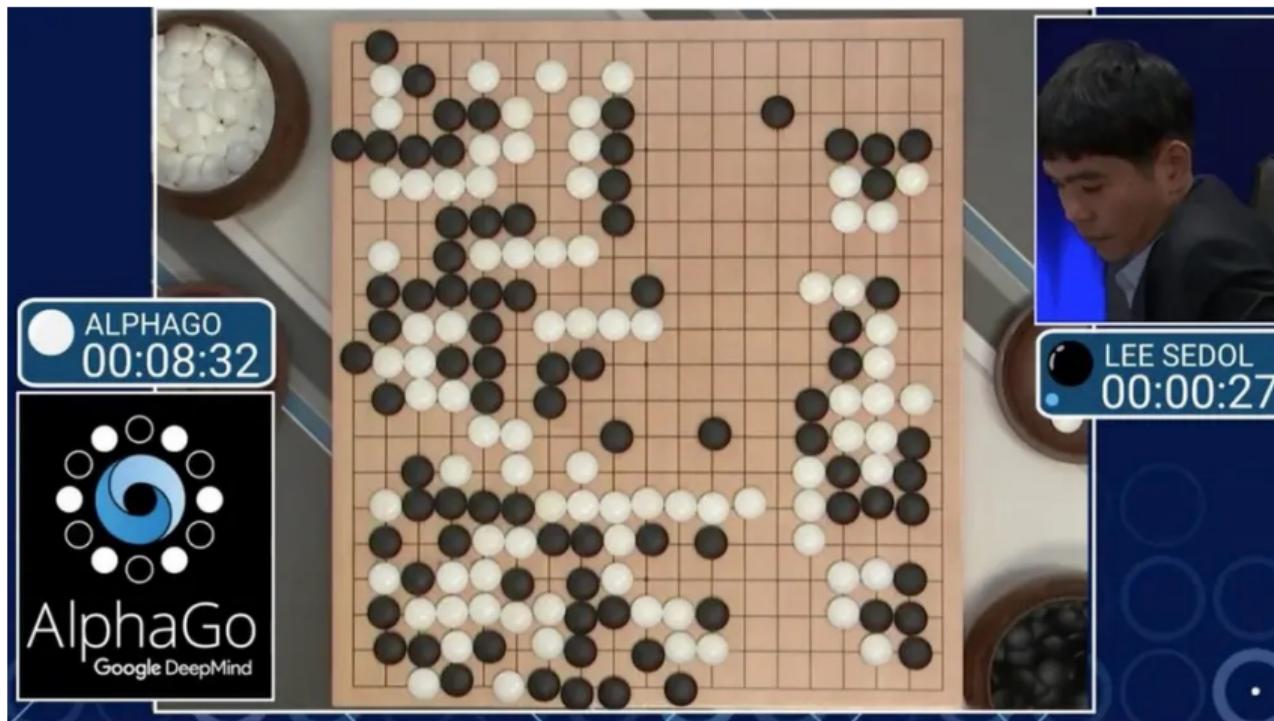
Drawbacks of Instruction Finetuning

- Collecting high-quality, task-specific ground truth data can be challenging.
- Some tasks, such as open-ended or creative generation, lack a single correct answer.
- Is it feasible to directly provide human feedback to train LLMs?

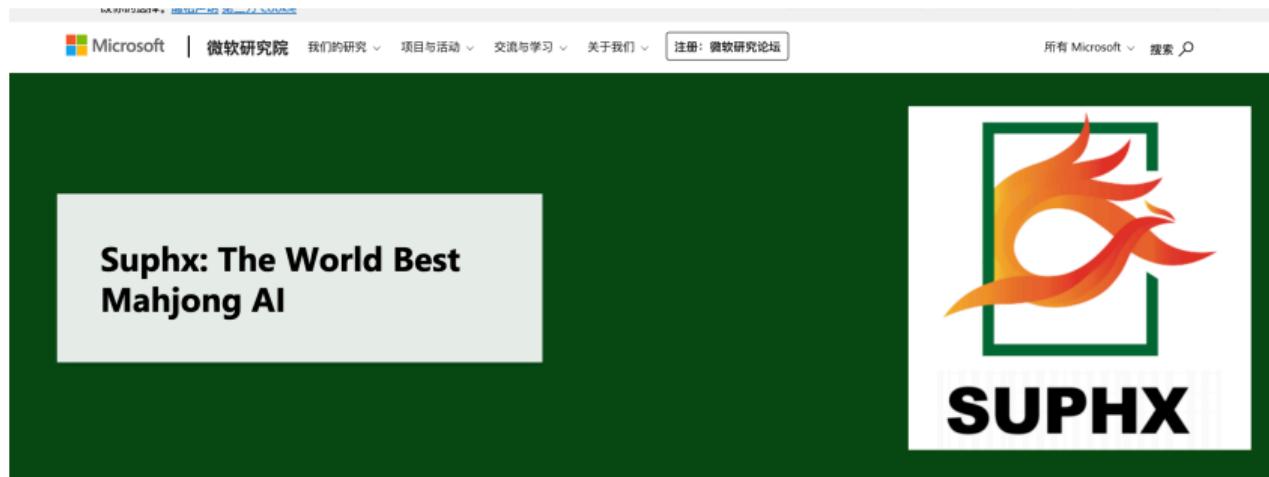
① Section 1: Instruction finetuning

② Section 2: RLHF

Reinforcement learning (RL) is widely used



Reinforcement learning (RL) is widely used



Reinforcement learning (RL) is widely used



Reinforcement learning (RL) is widely used

澎湃新闻

要闻 深度 直播 视频 时事 国际 财经 科技 暖闻 澎湃号 更多

上海市区出现首台无人驾驶清扫机器人，工作起来“情绪稳定”

澎湃新闻记者 许海峰 实习生 王雨
2025-04-11 10:10 来源：澎湃新闻·快看 >

字号▼

18
11
11
11
11



30/44

What is reinforcement learning?

How Much Information is the Machine Given during Learning?

Y. LeCun

- ▶ “Pure” Reinforcement Learning (**cherry**)
 - ▶ The machine predicts a scalar reward given once in a while.
 - ▶ **A few bits for some samples**

- ▶ Supervised Learning (**icing**)
 - ▶ The machine predicts a category or a few numbers for each input
 - ▶ Predicting human-supplied data
 - ▶ **10→10,000 bits per sample**

- ▶ Self-Supervised Learning (**cake génoise**)
 - ▶ The machine predicts any part of its input for any observed part.
 - ▶ Predicts future frames in videos
 - ▶ **Millions of bits per sample**



What is Reinforcement Learning?

*"Reinforcement learning is learning what to do—how to map **situations** to **actions**—so as to maximize a numerical **reward** signal. The learner is not told which actions to take, but instead must discover which actions yield the most reward by trying them."*

— Sutton and Barto, 2018

The banner has a dark blue background with a blurred digital interface or network visualization in the background. On the left, the NSF logo is displayed. In the center, there are two small portrait photos of men: Andrew Barto on the left and Richard Sutton on the right. Below the portraits, the text reads: **A.I. Pioneers**, **ANDREW BARTO and RICHARD SUTTON**, and **Receive Turing Award**.

A.I. Pioneers
ANDREW BARTO and RICHARD SUTTON
Receive Turing Award

RL Meets LLMs

- **Action:** Token generation.
- For each generated response s , a **human reward (human feedback)** $R(s)$ is obtained. The higher the reward, the better the response.

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.

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

- Maximize the expected reward:

$$\max_{\theta} \quad \mathbb{E}_{s \sim p_{\theta}(s)} R(s),$$

where p_{θ} represents an LLM.

RLHF pipeline

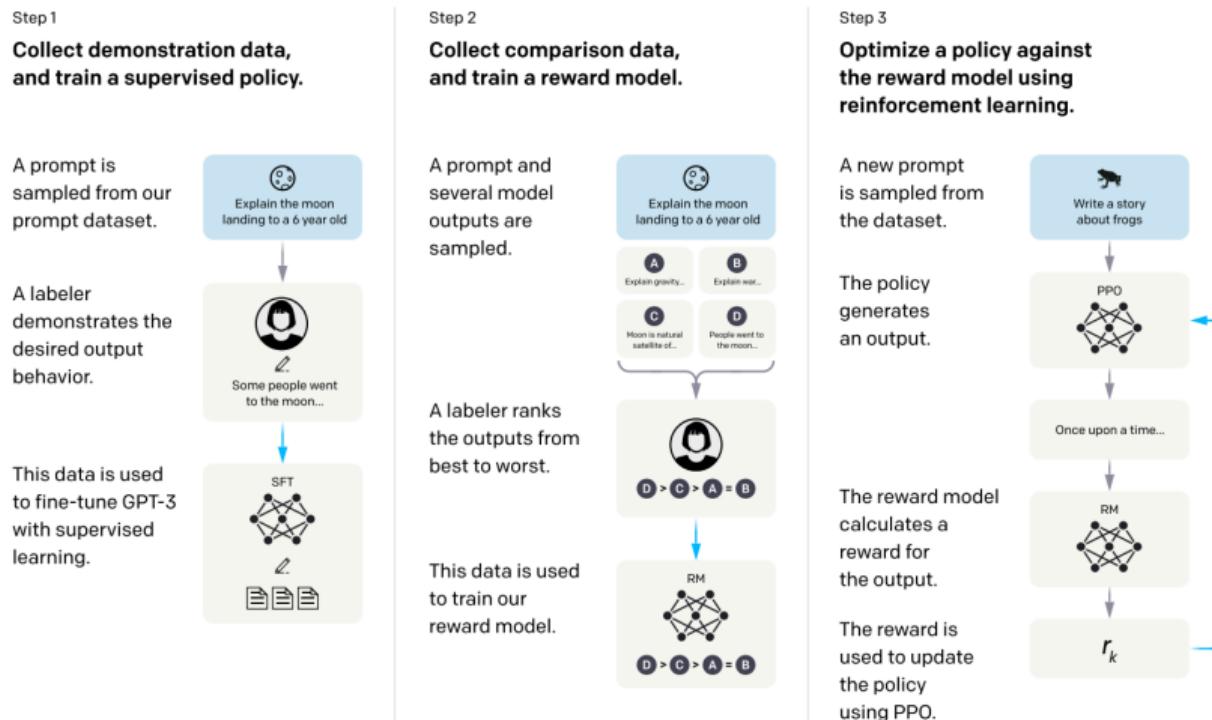


Figure is from <https://arxiv.org/pdf/2203.02155>

Le Cun's Analogy

- **Step 0:** A pretrained LLM using self-supervised learning.
- **Step 1:** Supervised instruction finetuning.
- **Step 2 & 3:** Using RL to maximize the reward.

How to maximize the expected reward?

- Recap:

$$\max_{\theta} \quad \mathbb{E}_{s \sim p_{\theta}(s)} R(s).$$

- How about stochastic gradient ascent?
-

$$\theta_{t+1} = \theta_t + \eta \nabla_{\theta_t} \mathbb{E}_{s \sim p_{\theta_t}(s)} R(s)$$

How to Calculate the Gradient of the Expected Reward?

- Calculating $\nabla_{\theta_t} \mathbb{E}_{s \sim p_{\theta_t}(s)} R(s)$ is challenging.
- The reward function R is a black box, and its dependency on θ is unknown.
- The reward function R is often non-differentiable.

Policy gradient method

- We want to compute

$$\nabla_{\theta} \mathbb{E}_{s \sim p_{\theta}(s)} R(s) = \nabla_{\theta} \sum_t R(t)p_{\theta}(t) = \sum_t R(t)\nabla_{\theta} p_{\theta}(t).$$

- Log-derivative trick:

$$\nabla_{\theta} p_{\theta}(t) = p_{\theta}(t)\nabla_{\theta} \log p_{\theta}(t).$$

- Then, we have

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

Policy gradient method

- We have put the gradient inside the expectation:

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

- This expectation can be approximated using the Monte Carlo samples, i.e.,

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

- Stochastic gradient ascent:

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

Modeling Human Preference

- Suppose we have a reward model $R_\phi(s_i)$ for each s_i .
- If we can observe the true reward $R(s)$, we can train the reward model.
- However, human judgments are often noisy and miscalibrated.
- In this case, we may consider pairwise comparisons: $(R(s_1) > R(s_2))$.

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.

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

Pairise comparisons of human preserences

- For a pair of sentences, we model

$$\text{Loss}_{\text{pair}}(\phi) = \frac{1}{m(m-1)} \sum_{R(s_i) > R(s_j)} \ell(R_\phi(s_i) - R_\phi(s_j)),$$

where ℓ is an individual loss function.

Proximal Policy Optimization (PPO)

What we have

- A pretrained model $p_{pre}(s)$
- A reward model R_ϕ
- A RLHF model p_θ to be estimated.

Now, maximize the following reward with RL

$$R(s) = R_\phi(s) - \beta \log \frac{p_\theta(s)}{p_{pre}(s)}$$

- The penalty is to prevent the model from diverging too far from the pretrained mode.
- It is the Kullback-Leibler (KL) divergence between p_θ and p_{pre} .

Scaling of GPT-3 with RLHF

