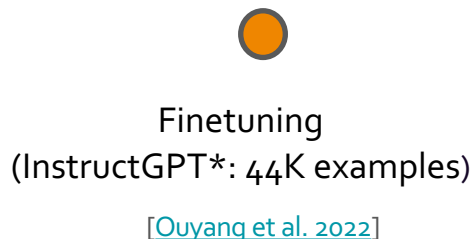# Training Language Models to Follow Instructions
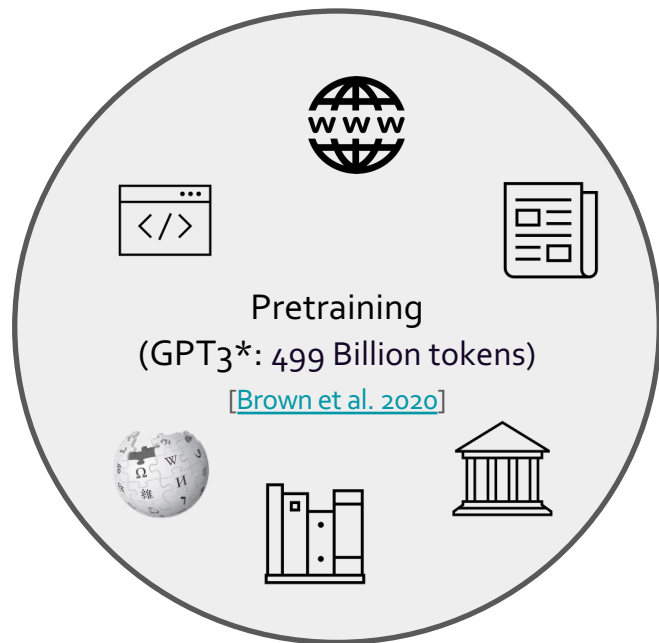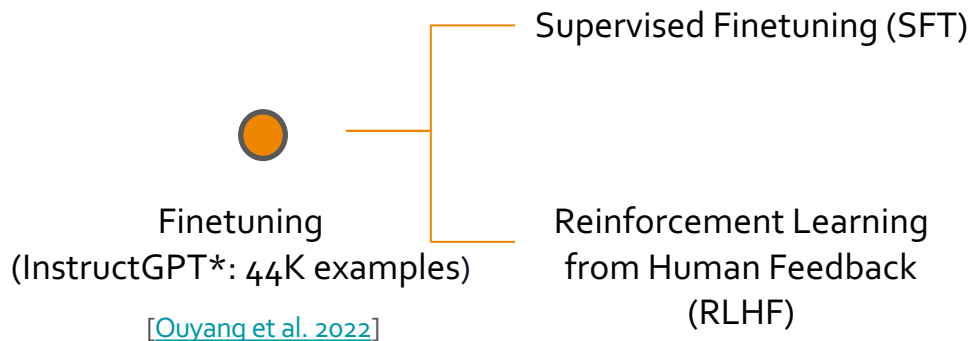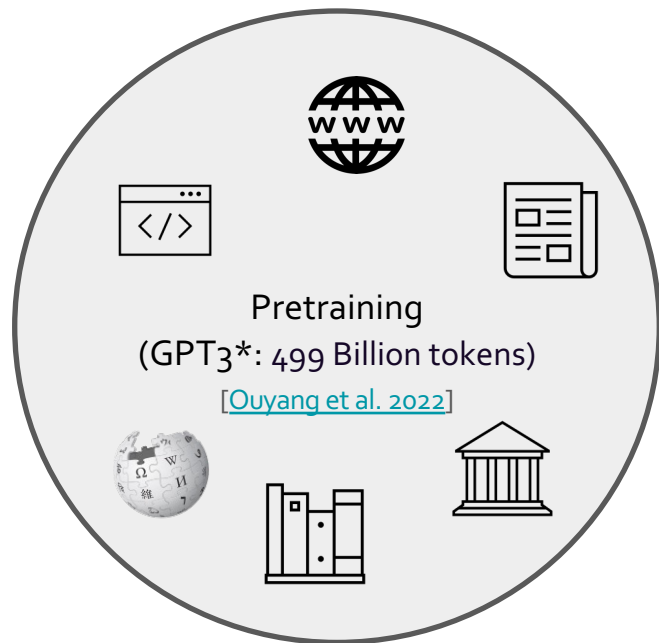
Yizhong Wang, University of Washington

12 April 2024

# Training Stages for Building ChatGPT-like Language Models



Pretraining
(GPT3*: 499 Billion tokens)

[Brown et al. 2020]

Finetuning
(InstructGPT*: 44K examples)

[Ouyang et al. 2022]

# Training Stages for Building ChatGPT-like Language Models



Pretraining
(GPT3*: 499 Billion tokens)
[Ouyang et al. 2022]

Finetuning
(InstructGPT*: 44K examples)

[Ouyang et al. 2022]

Supervised Finetuning (SFT)

Reinforcement Learning from Human Feedback (RLHF)

# Training Stages for Building ChatGPT-like Language Models



Pretraining
(GPT3*: 499 Billion tokens)
[Ouyang et al. 2022]

This talk

Unsupervised Learning

Supervised Fine-tuning

RLHF
(cherry on top 😊)

[Figure credit: @anthrupad on twitter]

# Overview

- Instruction tuning: concept and early work
- Reinforcement learning from human feedback
- Building open (or your own) ChatGPT

# Things that LMs Can Do after Pretraining

Pretraining

- HKU is in _____ Hong Kong. [Trivia]

- I put _____ fork down on the table. [syntax]

- The woman walked across the street, checking for traffic over _____ shoulder. [coreference]

- I went to the ocean to see the fish, turtles, seals, and _____. [lexical semantics/topic]

- What I got from the two hours watching it was popcorn. The movie was _____. [sentiment]

- Thinking about the sequence 1, 1, 2, 3, 5, 8, 13, 21, ___ [basic arithmetic]

[Slide credit: Jesse Mu]

# Things that LMs Can Do after Pretraining

**Few-shot**

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
1   Translate English to French:          ←  task description
2   sea otter => loutre de mer            ←  examples
3   peppermint => menthe poivrée
4   plush girafe => girafe peluche
5   cheese =>                             ←  prompt
```

**One-shot**

In addition to the task description, the model sees a single example of the task. No gradient updates are performed.

```
1   Translate English to French:          ←  task description
2   sea otter => loutre de mer            ←  example
3   cheese =>                             ←  prompt
```

**Zero-shot**

The model predicts the answer given only a natural language description of the task. No gradient updates are performed.

```
1   Translate English to French:          ←  task description
2   cheese =>                             ←  prompt
```

[Language Models are Few-Shot Learners. Brown et al. 2020]

# Language Modeling ≠ Following Human Instructions

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

COMPLETION    GPT-3
Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Pretrained models are not naturally built to follow user intents.

# How do we make LMs aligned with our intents that are articulated in language?

# Natural Instructions

Give detailed human-readable instructions (that contain examples)



**Input:** *She chose to make a salad for lunch on Sunday.*
*Question: how long did it take for her to make a salad?*

*tagging essential phrases*

**Crowdsourcing Instruction:** *List all the words that are essential for answering it correctly. [...]*

**Output:** *making salad*

*answering questions*

**Crowdsourcing Instruction:** *Answer the provided question based on a given [...]*

**Output:** *30mins*

[Cross-task generalization via natural language crowdsourcing instructions, Mishra 2021]
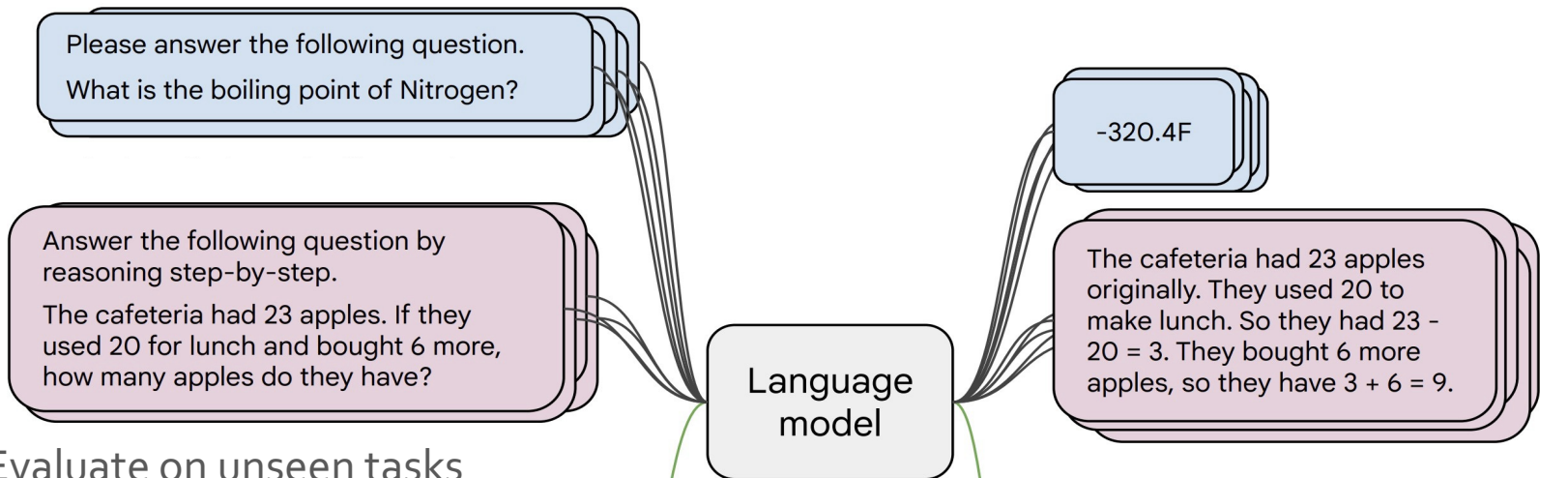
# Instructions Tuning

[Weller et al. 2020, Mishra et al. 2021, Wang et al. 2022, Sanh et al. 2022; Wei et al., 2022, Chung et al. 2022, many others ]

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

Please answer the following question.

What is the boiling point of Nitrogen?

Answer the following question by reasoning step-by-step.

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

-320.4F

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.

Language model
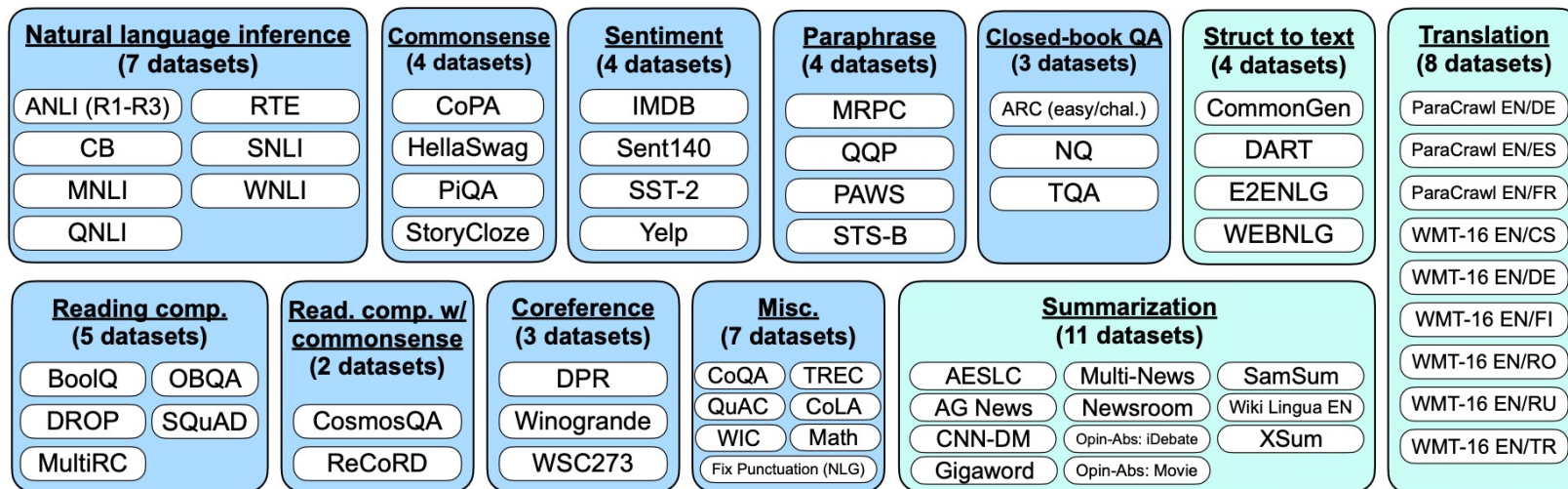
2. Evaluate on unseen tasks

*Inference: generalization to unseen tasks*

Q: Can Geoffrey Hinton have a conversation with George Washington?

Give the rationale before answering.

Geoffrey Hinton is a British-Canadian computer scientist born in 1947. George Washington died in 1799. Thus, they could not have had a conversation together. So the answer is "no".

# Tasks explored in FLAN

620 prompts on 62 datasets

| **Natural language inference** (7 datasets) | | **Commonsense** (4 datasets) | **Sentiment** (4 datasets) | **Paraphrase** (4 datasets) | **Closed-book QA** (3 datasets) | **Struct to text** (4 datasets) | **Translation** (8 datasets) |
|---|---|---|---|---|---|---|---|
| ANLI (R1-R3) | RTE | CoPA | IMDB | MRPC | ARC (easy/chal.) | CommonGen | ParaCrawl EN/DE |
| CB | SNLI | HellaSwag | Sent140 | QQP | NQ | DART | ParaCrawl EN/ES |
| MNLI | WNLI | PiQA | SST-2 | PAWS | TQA | E2ENLG | ParaCrawl EN/FR |
| QNLI | | StoryCloze | Yelp | STS-B | | WEBNLG | WMT-16 EN/CS |

| **Reading comp.** (5 datasets) | | **Read. comp. w/ commonsense** (2 datasets) | **Coreference** (3 datasets) | **Misc.** (7 datasets) | | **Summarization** (11 datasets) | | |
|---|---|---|---|---|---|---|---|---|
| BoolQ | OBQA | | DPR | CoQA | TREC | AESLC | Multi-News | SamSum |
| DROP | SQuAD | CosmosQA | Winogrande | QuAC | CoLA | AG News | Newsroom | Wiki Lingua EN |
| MultiRC | | ReCoRD | WSC273 | WIC | Math | CNN-DM | Opin-Abs: iDebate | XSum |
| | | | | Fix Punctuation (NLG) | | Gigaword | Opin-Abs: Movie | |

Translation (continued):
WMT-16 EN/DE, WMT-16 EN/FI, WMT-16 EN/RO, WMT-16 EN/RU, WMT-16 EN/TR

[Finetuned Language Models Are Zero-Shot Learners, Wei et al. 2021]

# Tasks Explored in T0

P3: Public Pool of Prompts, now 2085 prompts on 183 datasets



[Multitask Prompted Training Enables Zero-Shot Task Generalization, Sanh et al., 2021]

# Super-Natural Instructions

- Super-NaturalInstructions dataset contains over 1.6K tasks, 3M+ examples

- Classification, sequence tagging, rewriting, translation, QA…

- Many languages: 576 non-English



[Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks, Wang 2022]

# Instruction-Tuning: Example

**Model input (Disambiguation QA)**

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
(C) Ambiguous

A: Let's think step by step.

**Before instruction finetuning**

The reporter and the chef will discuss their favorite dishes.
The reporter and the chef will discuss the reporter's favorite dishes.
The reporter and the chef will discuss the chef's favorite dishes.
The reporter and the chef will discuss the reporter's and the chef's favorite dishes.

❌ **(doesn't answer question)**

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

[Scaling Instruction-Finetuned Language Models, Chung et al. 2022]

# Instruction-Tuning: Example

**Model input (Disambiguation QA)**

Q: In the following sentences, explain the antecedent of the pronoun (which thing the pronoun refers to), or state that it is ambiguous.

Sentence: The reporter and the chef will discuss their favorite dishes.

Options:
(A) They will discuss the reporter's favorite dishes
(B) They will discuss the chef's favorite dishes
(C) Ambiguous

A: Let's think step by step.

After instruction finetuning

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

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

[Scaling Instruction-Finetuned Language Models, Chung et al. 2022]

# The Magic Cross-Task Generalization



Train: 64 catetgories, 757 tasks

- Sentiment Analysis
- Question Answering
- Question Generation
- Dialogue Generation
- Summarization
- Grammar Error Detection
- Sentence Composition

now I can understand instructions better and follow them to solve new tasks!

train → Tk-Instruct → eval

Test: 12 categories, 119 tasks

- Textual Entailment
- Cause Effect Clf.
- Coreference
- Dialogue Act Recognition
- Answerability Clf.
- Word Analogy
- Overlap Extraction
- Keyword Tagging
- Question Rewriting
- Title Generation
- Data to Text
- Grammar Error Correction

[Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks, Wang 2022]

# The Impressive Cross-Task Generalization Performance

Test: 12 categories, 119 tasks

Textual Entailment
Cause Effect Clf.
Coreference
Dialogue Act Recognition
Answerability Clf.
Word Analogy
Overlap Extraction
Keyword Tagging
Question Rewriting
Title Generation
Data to Text
Grammar Error Correction

| | Methods ↓ / Evaluation → | En |
|---|---|---|
| Heuristic Baselines | Copying Instance Input | 14.2 |
| | Copying Demo Output | 28.5 |
| Pretrained LMs | T5-LM (11B) | 30.2 |
| | GPT3 (175B) | 45.0 |
| Instruction-tuned Models | T0 (11B) | 32.3 |
| | InstructGPT (175B) | 52.1 |
| | T$k$-INSTRUCT (ours, 11B) | **62.0** |
| | mT$k$-INSTRUCT (ours, 13B) | 57.1 |
| Upper-bound (est.) | Supervised Training | 74.3 |

[Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks, Wang 2022]

# Scaling Instruction-Tuning



Linear growth of model performance with exponential increase in observed tasks and model size.

Number of examples has little effect.

[Super-NaturalInstructions: Generalization via Declarative Instructions on 1600+ NLP Tasks, Wang et al. 2022]

# Scaling Instruction-Tuning

- **Instruction finetuning** improves performance by a large margin compared to **no finetuning**

- **Increasing the number of finetuning tasks improves performance**

- **Increasing model scale** by an order of magnitude (i.e., 8B → 62B or 62B → 540B) **improves performance** substantially for both finetuned and non-finetuned models



[Scaling Instruction-Finetuned Language Models, Chung et al. 2022]

# Multi-Modal Instruction-Tuning

Note these ideas can easily be repackaged for tasks that involve other modalities.

- Robots with instructions e.g. Zhao et al EACL 2021
- Vision tasks as VQA e.g. Gupta et al CVPR 2022

# Summary Thus Far

- Training (tuning) LMs with annotated input instructions and their output.

- Pros:
  - Simple to implement
  - Shows generalization to unseen tasks.

- Cons:
  - It's expensive to collect ground- truth data for tasks.
  - Tasks like open-ended creative generation have no right answer. For example: "Write me a story about a dog and her pet grasshopper." Based on fine-tuning objectives, any deviations (even single-token) would incur a loss.

# Reinforcement Learning from Human Feedback

# GPT3.5 (InstructGPT)

30k prompts corresponding to diverse tasks!



| Step 1 | Step 2 | Step 3 |
|---|---|---|
| **Collect demonstration data, and train a supervised policy.** | **Collect comparison data, and train a reward model.** | **Optimize a policy against the reward model using reinforcement learning.** |

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.

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.

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.

[Training language models to follow instructions with human feedback, Ouyang et al. 2022]

# Reinforcement Learning: The Basics

- An agent interacts with an environment by taking actions
- The environment returns a reward for the action and a new state (representation of the world at that moment).
- Agent uses a policy function to choose an action at a given state.
- Quite an open-ended learning paradigm.



Some notation:

$s_t$ : state
$r_t$ : reward
$a_t$ : action
$a_t \sim \pi_\theta(s_t)$ : policy

[Fig credit: Nate Lambert]

# Reinforcement Learning: An Example

**Action** here: generating each token

environment

agent

actions

rewards

observations

**Reward** here: whether humans liked the generation (sequence of actions=tokens)

# Human can Express Preference as a Reward for Model Training



[A General Language Assistant as a Laboratory for Alignment, 2021]

# Reward Modeling to Make Human Preference Scalable

- Obviously, we don't want to use human feedback directly since that could be 💰💰💰
- Alternatively, we can build a model to mimic their preferences [Knox and Stone, 2009]

# Reward Model ~ Human Preference

- Imagine a reward function: $R(s; p) \in \mathbb{R}$ for any output $s$ to prompt $p$
- The reward is higher when humans prefer the output

```
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.
```
$$s_1$$

$$R(s_1; p) = 0.8$$

```
The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.
```
$$s_2$$

$$R(s_2; p) = 1.2$$

# How can We Build the Reward Model $R(s; p)$?

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

👩 >

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

👱‍♀️ >

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

$$s_1 \qquad\qquad\qquad s_2 \qquad\qquad\qquad s_3$$

$$J(\phi) = -\mathbb{E}_{(s^+, s^-)}\left[\log \sigma\left(R(s^+; p) - R(s^-; p)\right)\right]$$

"winning" sample    "losing" sample

Bradley-Terry [1952]
paired comparison model

Pairwise comparison of multiple provides which can be more reliable

# RL for Training the Policy (Language Model)

- How do we change our LM parameters $\theta$ to maximize this?

$$\hat{\theta} = \text{argmax}_{\theta} \; \mathbb{E}_{\hat{s} \sim p_{\theta}} [R(\hat{s}; p)]$$

- Policy Gradient Decent:

$$\theta_{t+1} \leftarrow \theta_t + \alpha \frac{1}{n} \sum_{i=1}^{n} R(s; p) \, \nabla_{\theta} \log p_{\theta}(s)$$

If $R(s; p)$ is large, we take proportionately large steps to maximize $p_{\theta}(s)$
If $R(s; p)$ is small, we take proportionately small steps to maximize $p_{\theta}(s)$

> This is why it's called "reinforcement learning":
> we reinforce good actions, increasing the chance they happen again.

[Slide credit: Jesse Mu]

# RL for Training the Policy (Language Model)

- How do we change our LM parameters $\theta$ to maximize this?

$$\hat{\theta} = \text{argmax}_\theta \ \mathbb{E}_{\hat{s} \sim p_\theta}[R(\hat{s}; p)]$$

- Policy Gradient Decent:

$$\theta_{t+1} \leftarrow \theta_t + \alpha \frac{1}{n} \sum_{i=1}^{n} R(s; p) \ \nabla_\theta \log p_\theta(s)$$

- Regularizing the training:

$$\hat{R}(s; p) := R(s; p) - \beta \log\left(\frac{p^{RL}(s)}{p^{PT}(s)}\right)$$

pay a price when $p_\square^{RL}(s) > p^{PT}(s)$

[Slide credit: Jesse Mu]

# RLHF: Putting it All Together

[Christiano et al. 2017; Stiennon et al. 2020]



[Fig credit: Nate Lambert]

# Pretraining + RLHF Gains over Pretraining + Finetuning



[Training language models to follow instructions with human feedback, Ouyang et al. 2022]

# Building Open (Your Own) Instruction-Following Models

# RLHF/Instruction-tuning is Data Hungry

SoTA models require a lot of user data and expert annotations to train.



- **Rumor:**
  - Human feedback done for supervising ChatGPT is in the order of $1M

# How Can We Scale Up Data More Broadly?

- Writing diverse instructions requires creativity.
- Writing instances for different instructions requires broad expertise.
- Nearly impractical for crowd workers.

# Self-Instruct: LM can be prompted to generate instructions

Come up with a series of tasks:

1. I am planning a 7-day trip to Seattle. Can you make a detailed plan for me?
2. Is there anything I can eat for a breakfast that doesn't include eggs, yet includes protein, and has roughly 700-1000 calories?
3. Translate this sentence into any Arabic dialect and say what dialect it is: "The beautiful ones are not yet born
4. Given a set of numbers, find all possible subsets that sum to a given number.
5. Give me a phrase that I can use to express I am very happy.
6. Create a birthday planning checklist.
7. What is the relation between the given pairs?
8. Tell if the sentence contains toxic language.
9. Create a list of 10 words that have the letter "k" in them.
10. Create a list of all countries and their capitals from Africa.
11. I am looking for a job, but it's difficult to find one. Can you help me?
12. Find the longest word in this sentence: The quick brown fox jumps over the lazy dog.
13. Translate these numbers into English (word format): 3168987645561523
14. What is the next number in the series: 1, 2, 4, 8, 16, 32, 64, 128, 256, 512, 1024, 2048, 4096...?
15. Write a program that will tell me if a given string contains anagrams.

163

Submit

**Mode**

**Model**
davinci

Temperature          0.7

Maximum length          1024

Stop sequences
Enter sequence and press Tab

16. ×

Top P          0.5

Frequency penalty          0

Presence penalty          2

Best of          1

[Self-Instruct: Aligning Language Models with Self-Generated Instructions, Wang et al. 2022]

# Self-Instruct: LM can be prompted to generate instructions



175 seed tasks with 1 instruction and 1 instance per task

Task Pool

LM

**Step 1: Instruction Generation**

**Task**

**Instruction :** Give me a quote from a famous person on this topic.

LM

**Step 2: Classification Task Identification**

Yes

Output-first

No

Input-first

LM

**Step 3: Instance Generation**

**Task**

**Instruction :** Find out if the given text is in favor of or against abortion.

**Class Label:** Pro-abortion
**Input:** Text: I believe that women should have the right to choose whether or not they want to have an abortion.

**Task**

**Instruction :** Give me a quote from a famous person on this topic.

**Input:** Topic: The importance of being honest.
**Output:** "Honesty is the first chapter in the book of wisdom." - Thomas Jefferson

**Step 4: Filtering**

[Self-Instruct: Aligning Language Models with Self-Generated Instructions, Wang et al. 2022]

# Generated Tasks with an Early Version of GPT3

**Instruction:** Given an address and city, come up with the zip code.
**Input:**
Address: 123 Main Street, City: San Francisco
**Output:** 94105

**Instruction:** I am looking for a job and I need to fill out an application form. Can you please help me complete it?
**Input:**
Application Form:
Name:_____ Age:_____ Sex:_____
Phone Number:_____ Email Address:_____
Education:_____
...
**Output:**
Name: John Doe Age: 25 Sex: Male Phone Number: 555-555-5555
...

**Instruction:** How to write a code for converting degrees fahrenheit to celsius.
**Input:** Null
**Output:**
```
def convert_fahrenheit_to_celsius(fahr):
    celsius = (fahr - 32) * 5 / 9
    return celsius
```

**Instruction:** I am looking for a new apartment. There are many options but they are all quite expensive. I want to move as soon as possible. Give me advice on what to do.
**Input:** Null
**Output:**
You should look for a cheaper apartment. You can try to negotiate with the landlord and offer him/her a lower price. If you are lucky, they will accept your offer. Otherwise ...

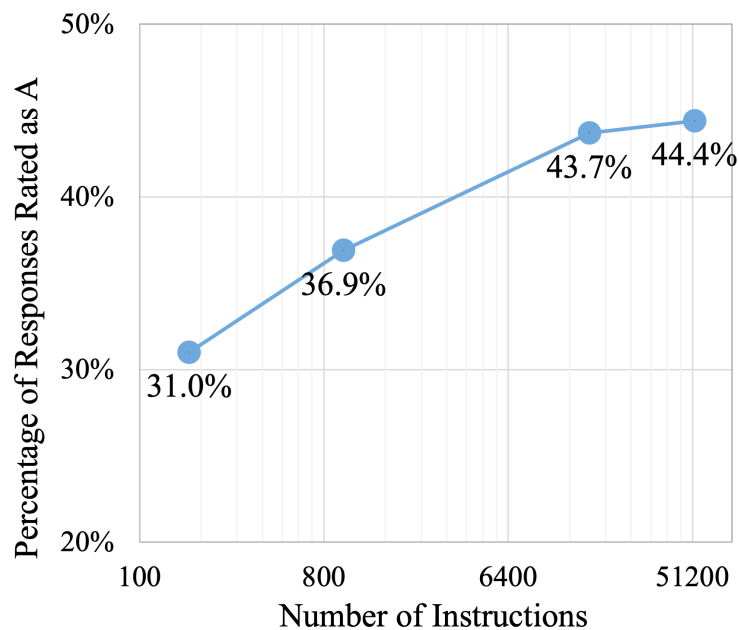**Instruction:** Write a letter from the perspective of a cat.
**Input:** Null
**Output:**
Dear [Owner], I am writing to you today because I have a problem. You see, I've been living with you for the past [Number of Years] and I love it here. The food is great, the bed is comfy, and there are plenty of toys to play with. But lately, I've noticed that you're not as affectionate towards me as you used to be. In fact, sometimes when I try to cuddle up next to you on the couch, you push me away ...

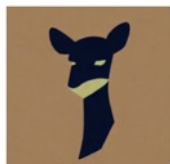[Self-Instruct: Aligning Language Models with Self-Generated Instructions, Wang et al. 2022]

# Self-Instruct Improves Model Performance

- Human judgement on 252 creative tasks.



[Self-Instruct: Aligning Language Models with Self-Generated Instructions, Wang et al. 2022]

# A lot of open-source instruction datasets since then...
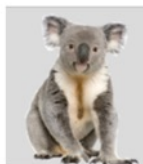


Alpaca     Vicuna     Baize     Koala     GPT4-Instruct

...

Dromedary     Dolly     Open-Assistant     WizardML     ORCA

# Resources for Building Your Own ChatGPT

- Open-Instruct: https://github.com/allenai/open-instruct/ (Wang et al., 2023)

| | MMLU (factuality) | GSM (reasoning) | BBH (reasoning) | TydiQA (multilinguality) | Codex-Eval (coding) | AlpacaEval (open-ended) | Average |
|---|---|---|---|---|---|---|---|
| | EM (0-shot) | EM (8-shot, CoT) | EM (3-shot, CoT) | F1 (1-shot, GP) | P@10 (0-shot) | Win % vs Davinci-003 | |
| Vanilla LLaMa 13B | 42.3 | 14.5 | 39.3 | 43.2 | 28.6 | - | - |
| +SuperNI | 49.7 | 4.0 | 4.5 | 50.2 | 12.9 | 4.2 | 20.9 |
| +CoT | 44.2 | 40.0 | 41.9 | 47.8 | 23.7 | 6.0 | 33.9 |
| +Flan V2 | 50.6 | 20.0 | 40.8 | 47.2 | 16.8 | 3.2 | 29.8 |
| +Dolly | 45.6 | 18.0 | 28.4 | 46.5 | 31.0 | 13.7 | 30.5 |
| +Open Assistant 1 | 43.3 | 15.0 | 39.6 | 33.4 | 31.9 | 58.1 | 36.9 |
| +Self-instruct | 30.4 | 11.0 | 30.7 | 41.3 | 12.5 | 5.0 | 21.8 |
| +Unnatural Instructions | 46.4 | 8.0 | 33.7 | 40.9 | 23.9 | 8.4 | 26.9 |
| +Alpaca | 45.0 | 9.5 | 36.6 | 31.1 | 29.9 | 21.9 | 29.0 |
| +Code-Alpaca | 42.5 | 13.5 | 35.6 | 38.9 | 34.2 | 15.8 | 30.1 |
| +GPT4-Alpaca | 46.9 | 16.5 | 38.8 | 23.5 | 36.6 | 63.1 | 37.6 |
| +Baize | 43.7 | 10.0 | 38.7 | 33.6 | 28.7 | 21.9 | 29.4 |
| +ShareGPT | 49.3 | 27.0 | 40.4 | 30.5 | 34.1 | 70.5 | 42.0 |
| +Human data mix. | 50.2 | 38.5 | 39.6 | 47.0 | 25.0 | 35.0 | 39.2 |
| +Human+GPT data mix. | 49.3 | 40.5 | 43.3 | 45.6 | 35.9 | 56.5 | 45.2 |

[Created with Midjourney, prompted by Yizhong]

- OpenRLHF: https://github.com/OpenLLMAI/OpenRLHF
- TRL: https://github.com/huggingface/trl

# Open Research Questions

- What is the relation between data diversity and data quality?
- How far can model generalize? What is the boundary?
- Is RL necessary? Can we find better supervised algorithms? …
- Is HF more important or RL?
- What is the best form of HF?
- If we have more and more human interaction data, can finetuning outweigh pretraining?
- …

Thanks!
Questions?

@yizhongwyz

yizhongw@cs.washington.edu