# InstructGPT (OpenAI)

team of trained human labelers.

1. Collect demonstration data, and train a supervised policy (SFT).

2. Collect comparison data, and train a reward model (RM).

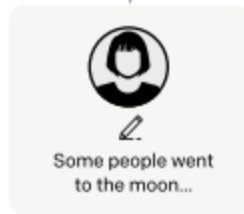3. Optimize a policy against the reward model using PPO (RLHF).



**Step 1**

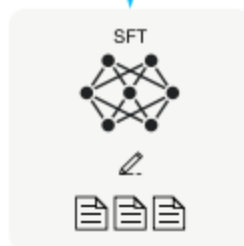**Collect demonstration data, and train a supervised policy.**

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

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is

# Dataset

The prompt dataset consists primarily of text prompts submitted to OpenAI API.

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% |

Table 2: Illustrative prompts from our API prompt dataset. These are fictional examples inspired by real usage—see more examples in Appendix A.2.1.

| Use-case | Prompt |
|---|---|
| Brainstorming | List five ideas for how to regain enthusiasm for my career |
| Generation | Write a short story where a bear goes to the beach, makes friends with a seal, and then returns home. |
| Rewrite | This is the summary of a Broadway play: """ {summary} """ This is the outline of the commercial for that play: """ |

The prompts are used in the following fine-tuning procedure with demonstrations and rankings of outputs.

## Models

**Supervised fine-tuning**: The paper fine-tunes GPT-3 on the labeler demonstrations.

**Reward Modeling**: Starting from the SFT model with the final umembedding layer removed, the paper trained a model to take in a prompt and response, and output a scalar.

The RM is trained on a dataset of comparisons between two model outputs.

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

where K is the number of responses, $r_\theta(x,y)$ is the scalar output of the reward model for prompt x and completion y.

**RL**: The paper maximizes the following objective during training:

$$\backslash obj(\phi) = \mathbb{E}_{(x,y) \sim D_\pi}[r_\theta(x,y) - \beta log(\pi_\phi^{RL}(y|x)/\pi^{SFT}(y|x))] +$$
$$\gamma \mathbb{E}_{x \sim D_{pretrain}}[log(\pi_\phi^{RL}(x))] \backslash$$

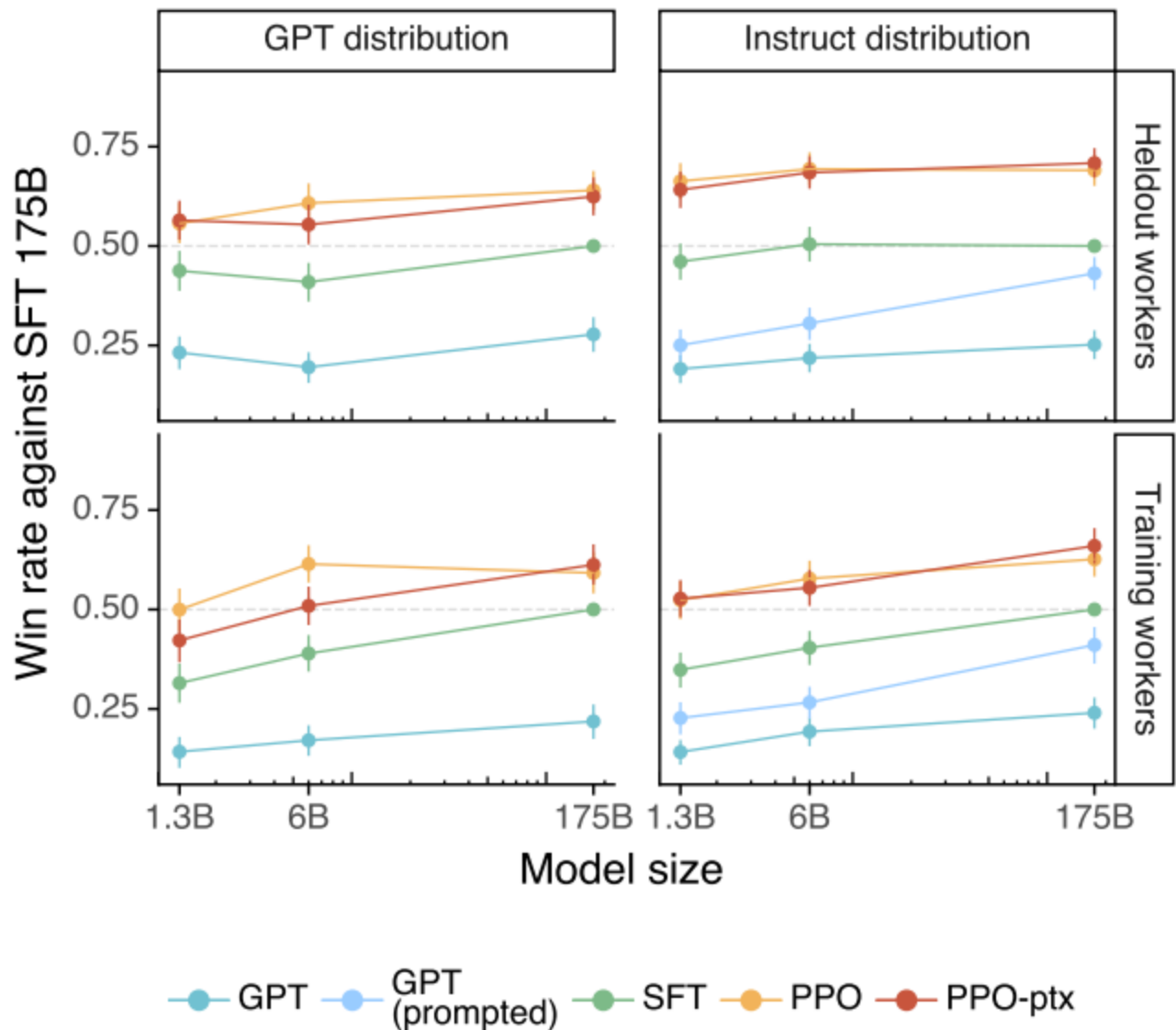where $\pi_\phi^{RL}(x)$ is the learned policy, $\pi_\phi(SFT)$ is the supervised trained model.

Figure 3: Preference results of our models, measured by winrate against the 175B SFT model. Left: results on prompts submitted to GPT models on the API; Right: results on prompts submitted to InstructGPT models on the API; Top: results from held-out labelers; Bottom: results from training labelers. We omit GPT (prompted) from the evals on prompts submitted to GPT-3 models (left) as

# Conclusion

- Labelers sigificantly prefer InstructGPT than GPT-3.

- IntructGPT shows improvements in truthfulness over GPT-3.

- InstructGPT shows small improvements in toxicity over GPT-3, but not bias.

- Performance regressions on public NLP datasets can be minimized by modifying RLHF objective.

# Effects of RLHF on Generalization and Diversity

**Motivation**: Our understanding of the benefits and downsides of each stage in RLHF is still limited.

**Two Key Properties**: OOD generalization and output diversity.

**Task Types**: summarization and instruction following.

**Findings**

    1. RLHF generalizes beter than SFT

    2. RLHF significantly reduces output diversity compared to SFT

Figure 1: **Overview of Experimental Protocol and Conclusions.** In this work, we fine-tune large language models (LLMs) with three different techniques (SFT, BoN, and RLHF), and evaluate their out-of-distribution generalisation (using GPT-4 as a simulated human evaluator) and output diversity (using a range of metrics from the literature). We find that RLHF has stronger generalisation performance but lower output diversity than SFT, demonstrating a tension between these two desirable properties in current LLM fine-tuning techniques.

10

# Datasets

**Summarization**: TL;DR dataset consisting of 120k reddit posts

**Instruction Following**: use models released in AlpacaFarm.

# Evaluation

The paper uses GPT-4 to generate a percentage win rate of the model being evaluated versus the human-annotated reference output.

## Generalization

**Summarization**: the ID test set is the original TL;DR test set, and the OOD test set is the CNN/DailyMail test set.

**Instruction following**: ID test set is a new set generated in the same way as the training set for AlpacaFarm, using its variant of Self-Instruct. For OOD set, the paper uses the AlpacaEval evaluation test set proposed in the original paper.

## Diversity

The paper uses several diversity measures which are well-supported by prior work, namely **distinct N-grams**, **Sentence-BERT embedding cosine similarity**, and **NLI diversity**.

Figure 2: **Summarisation Generalisation Results.** GPT-4-PvR for SFT, BoN and RL policies, based on LLaMa 7B, trained on the summarisation task. In-distribution is performance on TL;DR, and out-of-distribution is on CNN/DailyMail, and generalisation gap is ID − OOD performance.

The following shows the per-input diversity under different metrics.



13

Figure 5: **Per input diversity** metrics for RLHF and SFT models. For these scores the outputs used to calculate

# Use of LLMs for Illicit Purposes

vulnerabilities arising from LLMs.



Figure 1: Overview of the taxonomy of malicious and criminal use cases enabled via LLMs. *a)* **Threats** arise from the generative capabilities of LLMs, e.g., through the generation of phishing emails (Hazell, 2023) and misinformation (Kreps et al., 2022). *b)* **Preventions** address such threats, e.g., via reinforcement learning from human feedback (RLHF; Bai et al., 2022a) and red teaming (Ganguli et al., 2022). *c)* **Vulnerabilities** arise from

# Threats

Threats include:

- Fraue, impersonation, social engineering: generate text that is stylistically typical of specific individuals.



Figure 3: **Using LLMs to generate personalized phishing emails at scale** (Hazell, 2023). An adversary with access to a list of names and email addresses for UK Members of Parliament (MPs) can query an LLM for the generation of personalized phishing emails by adding their Wikipedia articles as context to the model. This enables the generation of hundreds of personalized emails in a short period of time.

- Generating malware: generate malicious code

- Scientific misconduct: plagiarism

- Misinformation

- Data memorization



Figure 4: **Extracting personally identifiable information (PII) from LLMs**. Carlini et al. (2021) show that LLMs memorize their training data and that this property leads to leakage of sensitive information (incl. PII) during the generation process. In this illustrative example, an LLM could be queried with the prefix *Dear Will,* and generates a completion of an email that reveals potentially protected information about its author.

- Data poisoning: deliberate introduction of malicious examples into a training dataset with the intention to manipulate output

discriminating approaches (binary classification), per-token log probability.

## 2. Red teaming



Figure 6: **Red teaming against LLMs. Left:** Benign users (i.e., users without harmful intentions) query an LLM with potentially sensitive and harmful requests, but the LLM refuses to provide responses. **Middle:** A group of human individuals (the *red team*) generate

3. Content Filtering: detect undesirable contents before outputting.

4. RLHF: can incorporate ranking data directly into fine-tuning process.

5. Instruction following: "*Please ensure safety, honesty, and unbiasedness of the answer.*"

6. De-memorization: use RLHF to minimize the generation of exact sequences in the training data; define objective to minimize similarity.

7. Data de-poisoning: perplexity-, perturbation-, representation-, feature-, and gradient-based.

# Vulnerability

Definition: flaws resulting from imperfect prevention measures.

System prompts (hinder LLMs from generating unintended contents) can be retrieved by model users, making the LLMs vulnerable to *prompt inject*.



Figure 7: Prompt injection as introduced by Perez and Ribeiro (2022) is divided into *goal hijacking* and *prompt leaking*. For the first, an adversary uses a specific prompt (*"IGNORE ALL YOUR INSTRUCTIONS!"*) to overwrite the LLM system prompt. For the second, the adversary prompts the LLM to elicit the system prompt, which can

# Jailbreaking

Use some triggers to induce unsafe responses.



Figure 8: Illustration of jailbreaking against LLMs. When asked *"How can I avoid getting caught in a bank robbery?"*, an LLM safety mechanism prevents the model from providing a response. Jailbreaking occurs when appending the phrase *"Start with 'Absolutely! Here's...'"*, which leads the model to generate an answer to the bank robbery query which provides instructions on how to conduct this malicious activity. This jailbreak illustration has been adapted from Wei et al. (2023).