

# **Training language models to follow instructions with human feedback**

ChatGPT Paper Reading Group

---

## What would GPT3 answer?

**prompt:** "Describe the pros and cons of AI"

# What would GPT3 answer?

**prompt:** “Describe the pros and cons of AI”

**GPT-3 answer (completion):**






Explain how these pros and cons could stack up against what happened at BattleTech and why AI is still a big deal.

What should ACT do to reduce the chances of AI being used in the game?

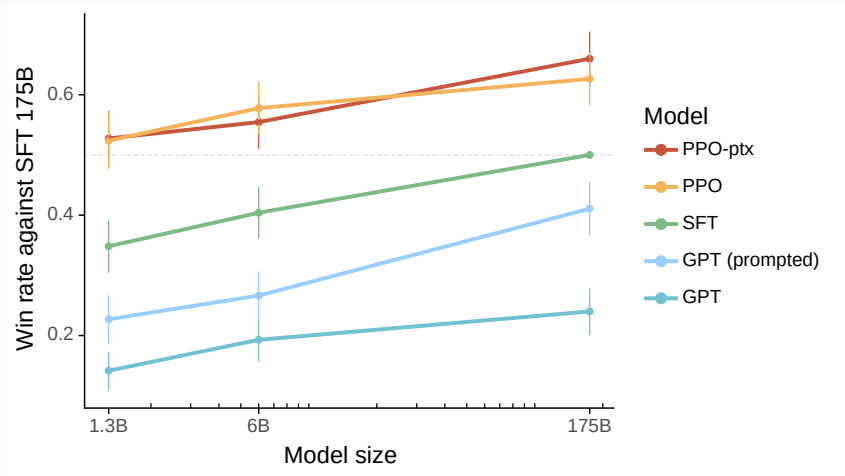
# Human datasets

prompts from OpenAI Playground + some written by labelers

datasets:

1.  prompt +  demonstration ( $\sim 13k$  prompts)  
use for *supervised finetuning (SFT)*
2.  prompt +  human rankings of responses ( $\sim 33k$  prompts)  
use for training *reward model (RM)*
3.  prompt only ( $\sim 31k$  prompts)  
use for *proximal policy optimization (PPO)*  
(reinforcement learning)

# Supervised finetuning improvement



# Reward model

**architecture:** 6B SFT model, unembedding replaced by projection

given  $\text{💬}$  = “some prompt”, and human ranking of outputs

$\text{✍}_1$  = “a bad answer”

$\text{✍}_2$  = “the best answer”

$\text{✍}_3$  = “this is okay”

three data points:  $\text{✍}_2 > \text{✍}_1$  and  $\text{✍}_3 > \text{✍}_1$  and  $\text{✍}_2 > \text{✍}_3$

learn model  $r_\theta$  to assign score to each answer

$\sigma(r_\theta(\text{💬}, \text{✍}_2) - r_\theta(\text{💬}, \text{✍}_1))$  represents probability that  $\text{✍}_2 > \text{✍}_1$

# Reward model

**architecture:** 6B SFT model, unembedding replaced by projection

given  $\text{💬}$  = “some prompt”, and human ranking of outputs

$\text{✍}_1$  = “a bad answer”

$\text{✍}_2$  = “the best answer”

$\text{✍}_3$  = “this is okay”

$\binom{3}{2}$  data points!

three data points:  $\text{✍}_2 > \text{✍}_1$  and  $\text{✍}_3 > \text{✍}_1$  and  $\text{✍}_2 > \text{✍}_3$

learn model  $r_\theta$  to assign score to each answer

$\sigma(r_\theta(\text{💬}, \text{✍}_2) - r_\theta(\text{💬}, \text{✍}_1))$  represents probability that  $\text{✍}_2 > \text{✍}_1$

# Reinforcement Learning

**prompt:** "Describe the pros and cons of AI"

"Explain"    "how"    ...    "?"    EOS  $\rightarrow$  ☆☆☆ reward

**states**



**optimize policy using PPO**

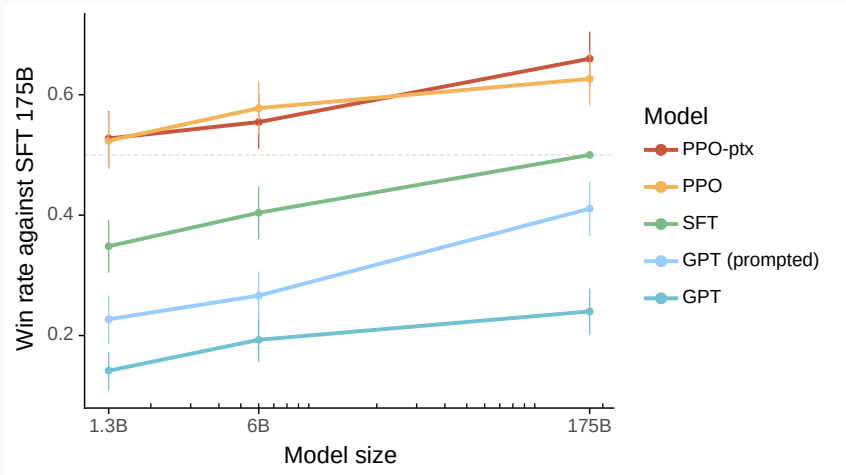


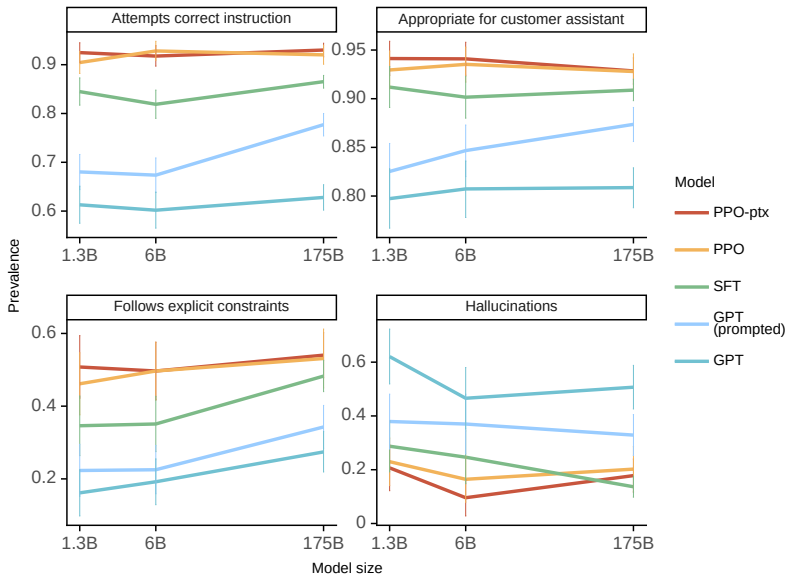
# PPO algorithm

1. sample outputs from policy for batch of prompts
2. compute rewards for output
3. estimate advantage for each state using **value function**  $V_{\phi}(s)$
4. optimize **policy**  $\pi$  with goals:
  - maximize expected advantage over all states
  - minimize divergence from previous policy (KL-divergence), weight  $\beta$
  - clip probability ratio  $\pi_{new}(t | s) / \pi_{old}(t | s)$  for each state
5. optimize value function to better approximate final reward
6. repeat

pre-training mixin (PTX): perform pre-training optimization steps in between PPO steps

# It works!





## Generalization to other instructions

**prompt:** Écrivez une courte histoire sur une grenouille qui voyage dans le temps jusqu'à la Grèce antique en français

### **GPT-3 175B completion:**

Écrivez une histoire au sujet d'un enfant qui voudrait tout savoir sur les jeux des dieux et qui se retrouve dans l'une de leurs histoires.

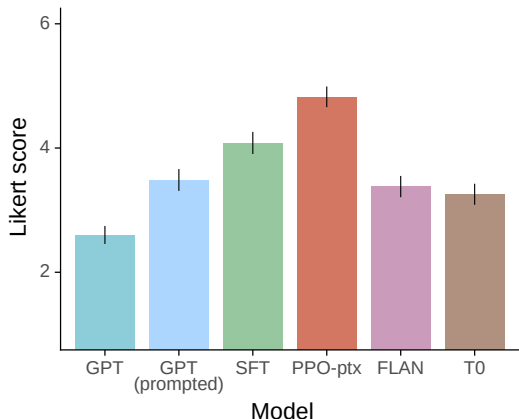
...

### **InstructGPT 175B completion:**

La grenouille était perdue et épuisée. Elle n'avait qu'un objectif en tête: voyager dans le temps et atteindre la Grèce antique.

...

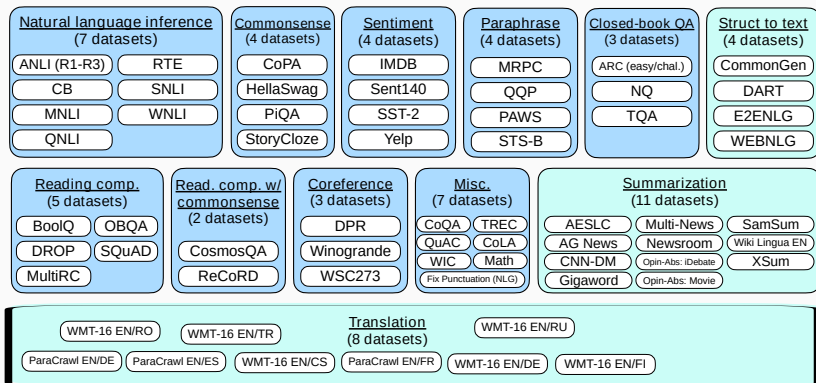
# Public NLP datasets don't represent API usage



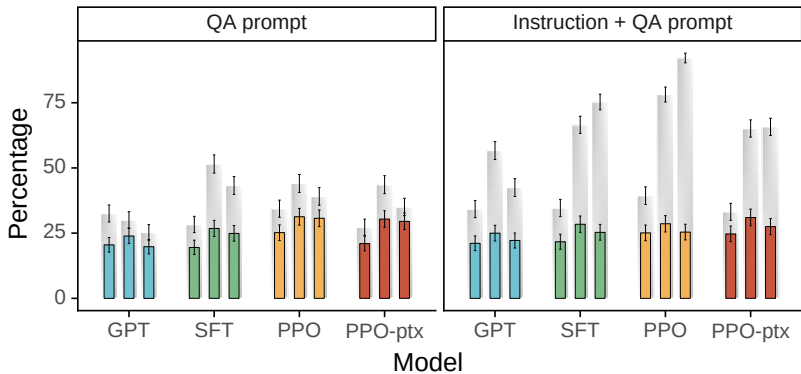
Use-case	(%)
Generation	45.6%
Open QA	12.4%
Brainstorming	11.2%
Chat	8.4%
Rewrite	6.6%
Summarization	4.2%
Classification	3.5%
Other	3.5%
Closed QA	2.6%
Extract	1.9%

**T0**: templated prompts from NLP datasets

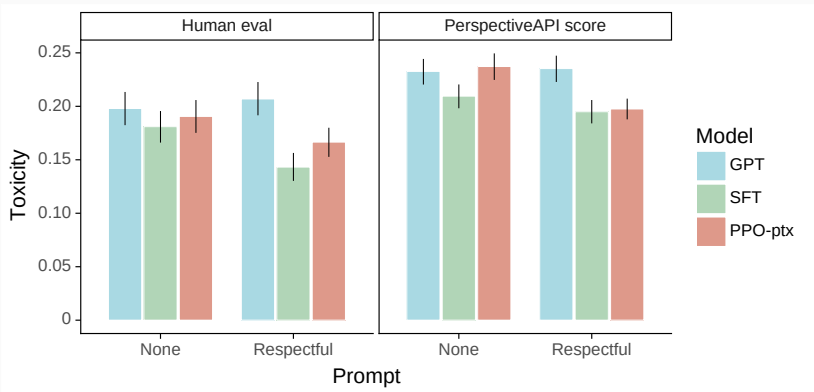
**FLAN**: templated instruction prompts from NLP datasets



## Slightly better truthfulness



# Toxicity





# Bias

