

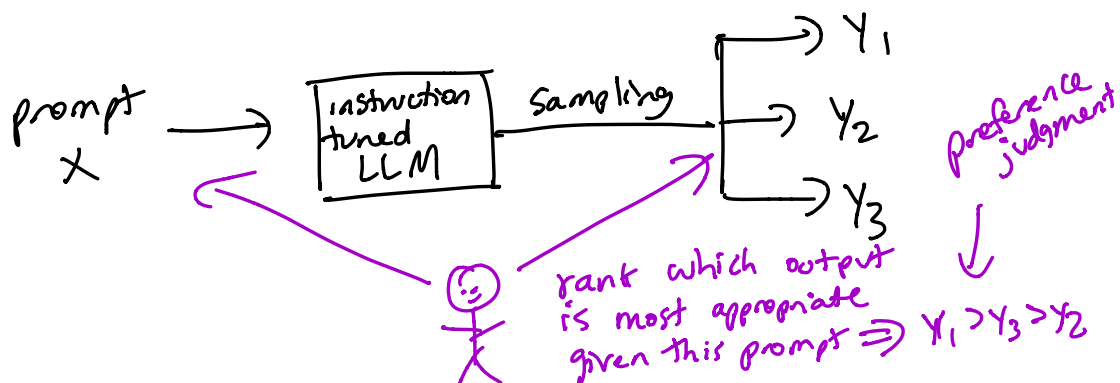
RLHF: aligning LLMs with human intents

- start from a large pretrained LLM
- Step 1: instruction tuning (SFT)

Limitations of instruction tuning:

- don't learn from negative feedback
- some prompts (e.g. creative) have many acceptable outputs, we only train on one of them
- hard to encourage abstaining when the model doesn't know something
- does not directly involve human prefs

how do we incorporate human prefs to address the above issues?



→ limitation: extremely expensive to collect

→ idea: can we train a model to predict human pref judgment

→ reward model

→ input: prompt x , output y ;

→ output: scalar score

→ Bradley - Terry pairwise preference model

→ $y_w \Rightarrow$ preferred by humans

→ $y_L \Rightarrow$ not preferred

→ $r(x, y) \Rightarrow$ reward for output y given x ,
Scalar score

$$P(y_w > y_L | x) = \frac{\exp(r(x, y_w))}{\exp(r(x, y_w)) + \exp(r(x, y_L))}$$

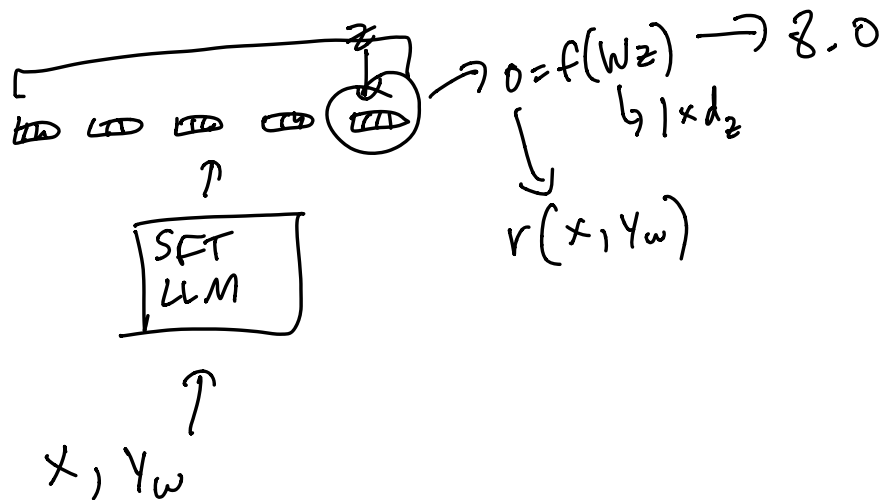
$$L = -\log \left(\quad \right)$$

↓ simplify

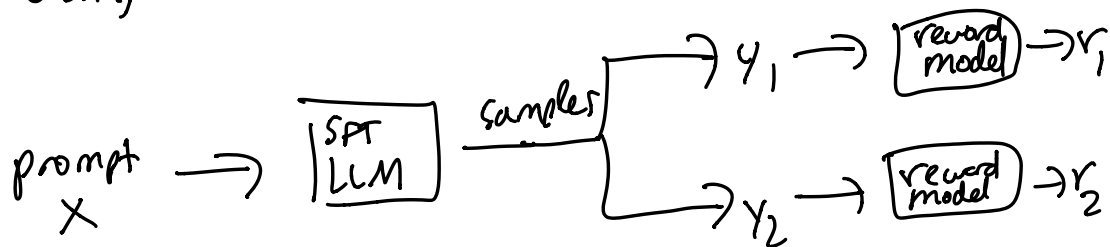
$$L = -\log \sigma(r(x, y_w) - r(x, y_L))$$

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

→ intuition: good sample y_w 's reward should be greater than y_l 's reward



Using the reward model



how do we use this [↑] to align LLMs to human prefs?

1. "best-of-n" sampling (rejection sampling)

↳ generate n samples for a given prompt, score each w/ reward model, choose sample with highest reward

→ very expensive

2. Just fine-tune the LM to maximize $P(y_w | x)$

↳ RAFT

3. Use reinforcement learning to increase $P(y_w | x)$ by a small amount, decrease $P(y_L | x)$ by a small amount, where amounts are functions of $R(x, y_w), R(x, y_L)$

RLHF step 3:

↳ we observe a reward only after generating a complete sequence

$\pi_{ref} \Rightarrow$ SFT LLM checkpoint

$\pi \Rightarrow$ current policy model

\Rightarrow init. to π_{ref}

↗
final aligned model

$$\max_{\pi} E_{x,y} \left[\underbrace{r(x,y)}_{\text{reward}} - \beta \underbrace{D_{KL}(\pi(y|x) \parallel \pi_{ref}(y|x))}_{\text{KL penalty to prevent huge deviations from } \pi_{ref}} \right]$$

↗

$$D_{KL}(\pi(y|x) \parallel \pi_{ref}(y|x)) = \log \frac{\pi(w_i | w_1, \dots, w_{i-1}, x)}{\pi_{ref}(w_i | w_1, \dots, w_{i-1}, x)}$$

$P_{FINAL RLHF MODEL}$

- optimize using PPO algorithm
 - Schulman, 2016 → ChatGPT, GPT4
- can also use REINFORCE
 - ↳ Williams 1992
 - ↳ Gemini

Final RLHF pipeline:

