# Lecture 23: More on Large Language Models and RLHF

COMP579, Lecture 23

#### Recall: RL in an LLM training loop



## **Recall: RLHF training phases**

base model (instruction, helpful, chatty etc.)



Vaswani et al. 2017

Pretrain+SFT, get preferences and fit reward model, do RL!

## Supervised fine-tuning (SFT): Same as imitation learning!

- We have (s, a) pairs, where s is a prompt and a is a generation corresponding to that prompt (consisting of several tokes
- These are taken from already existing data (eg internet docs, QA, solved problems...)
- Train a policy  $\pi_{sft}$  that maximizes the likelihood of the observed data:

$$J_{\mathsf{sft}}(\theta) = \mathbb{E}_{(s,a) \sim P_{\mathsf{sft}}} \left[ \frac{1}{|a|} \sum_{k=1}^{|a|} \log \pi_{\theta}(a_k | s, a_{i < k}) \right]$$

Training is done by gradient ascent

• Aka teacher forcing

## Better version: Rejection sampling (aka Best-of-N)



- Generate N answers from the model, reinforce the correct/top one(s)
- Train a policy  $\pi_{sft}$  that maximizes the likelihood of the *top data*:

$$J_{\mathsf{rft}}(\theta) = \mathbb{E}_{s \sim P_{\mathsf{sft}}, a \sim \pi_{\mathsf{sft}}(\cdot|s)} \left[ \frac{1}{|a|} \mathbf{1}_{a \text{is at the top}} \sum_{k=1}^{|a|} \log \pi_{\theta}(a_k|s, a_{i < k}) \right]$$

Training is done by gradient ascent

• Online rejection sampling finetuning:  $a \sim \pi_{\theta}$  instead of  $a \sim \pi_{sft}$ 

### **Recall: Training a reward model**



#### **Recall: Bradely-Terry loss function**

- Collect data from human raters: prompt s, pairs of  $a_w$ ,  $a_l$  responses
- Optimize the expected value of:

$$-\log(\sigma(r_{\phi}(s, a_w) - r_{\theta}(s, a_l)))$$

wrt reward parameter vector  $\phi$ 

- Cf. Ouyang et al, InstructGPT (2022)
- Corresponds to maximum likelihood fitting of binomial preference function if reward is linear over the variables

#### Recall: reward model almost as good as a single human



Evaluate RM on predicting outcome of held-out human judgments

#### **RLHF** finetuning



## **PPO for RLHF**



• Train a policy  $\pi_{\theta}$  that maximizes advantage:

$$J_{\mathsf{PPO}}(\theta) = \mathbb{E}_{s \sim P_{\mathsf{sft}}, a \sim \pi_{\theta_{\mathsf{old}}}(\cdot|s)} \left[ \frac{1}{|a|} \sum_{k=1}^{|a|} \frac{\pi_{\theta}(a_k|s, a_{i < k})}{\pi_{\theta_{\mathsf{old}}}(a_k|s, a_{i < k})} A_i \right]$$

where  $A_i$  is the advantage function

- Reward function uses a penalty per token for straying from reference policy:  $r_t = r_{\phi}(s, a_{< t}) \beta \log \frac{\pi_{\theta}(a_t | s, a_{< t})}{\pi_{\mathsf{sft}}(a_t | s, a_{< t})}$
- Value function/advantage needs to be estimated!

## GRPO (DeepSeek, 2025)



Instead of estimating value, use a group (non-parametric approach) Notation: q = s, o = a

#### **GRPO** Objective (DeepSeek, 2025)

- Generate G answers and estimate their reward (no regularization towards reference policy)
- Compute a normalized advantage based on the mean  $\bar{r}_t$  and standard deviation of the rewards:

$$\hat{A}_{i,t} = \frac{r_{i,t} - \bar{r}_t}{std(r_{1,t}, \dots r_{G,t})}$$

• GRPO objective - very similar to PPO!

$$J_{\mathsf{GRPO}}(\theta) = \mathbb{E}_{s \sim P_{\mathsf{sft}}, a_i \sim \pi_{\theta_{\mathsf{old}}}(\cdot|s)}, i=1, \dots G \left[ \frac{1}{G} \sum_{i=1}^G \frac{1}{|a_i|} \sum_{k=1}^{|a_i|} \frac{\pi_{\theta}(a_k|s, a_{i < k})}{\pi_{\theta_{\mathsf{old}}}(a_k|s, a_{i < k})} \hat{A}_{i,t} - \beta D_{KL}(\pi_{\theta}, \pi_{\theta_{\mathsf{old}}}) \right]$$

•  $D_{KL}$  is also estimated a bit differently (cf Shulman et al, 2020):

$$D_{KL}(\pi_{\theta}, \pi_{\theta_{\mathsf{old}}}) = \frac{\pi_{\theta_{\mathsf{old}}}(a_k | s, a_{i < k})}{\pi_{\theta}(a_k | s, a_{i < k})} - \log \frac{\pi_{\theta_{\mathsf{old}}}(a_k | s, a_{i < k})}{\pi_{\theta}(a_k | s, a_{i < k})} - 1$$

COMP579, Lecture 23

#### The Advantage of RL over SFT (DeepSeek, 2025)



COMP579, Lecture 23

#### DeepSeek Overall Results (DeepSeek, 2025)



SOTA results back in January!

COMP579, Lecture 23

#### **Direct Preference Optimization**



$$\nabla_{\theta} \mathcal{L}_{\text{DPO}}(\pi_{\theta}; \pi_{\text{ref}}) = -\beta \mathbb{E}_{(x, y_w, y_l) \sim \mathcal{D}} \left[ \underbrace{\sigma(\hat{r}_{\theta}(x, y_l) - \hat{r}_{\theta}(x, y_w))}_{\text{higher weight when reward estimate is wrong}} \left[ \underbrace{\nabla_{\theta} \log \pi(y_w \mid x)}_{\text{increase likelihood of } y_w} - \underbrace{\nabla_{\theta} \log \pi(y_l \mid x)}_{\text{decrease likelihood of } y_l} \right] \right],$$
$$\hat{r}_{\theta}(x, y) = \beta \log \frac{\pi_{\theta}(y|x)}{\pi_{\text{ref}}(y|x)}$$

- You can replace the complex RL part with a very simple weighted MLE objective
- Other variants (KTO, IPO) now emerging too

[Rafailov+ 2023]

#### Learning with non-transitive preferences: NashLLM

• Objective: find a policy  $\pi^*$  which is preferred over any other policy

$$\pi^* = \arg\max_{\pi} \min_{\pi'} \mathbb{P}(\pi' \preceq \pi)$$

- Think of this as a game: one player picks  $\pi$  the other picks  $\pi'$
- When both players use  $\pi^*$  this is a *Nash equilibrium* for the game
- For this game an equilibrium exists (even if eg preferences are not transitive)
- Cf. Munos et al, 2024 (https://arxiv.org/pdf/2312.00886.pdf)

#### NashLLM-style algorithms

- Fit a *two-argument preference function* by supervised learning
- Decide what is the *set of opponent policies*
- Ideally, the max player should play against a mixture of past policies
- *Optimize* using eg online mirror descent, convex-concave optimization...
- A lot of algorithmic variations to explore!

#### **NashLLM results**



Using preferences instead of rewards leads to less overfitting

## **General blueprint of RLHF training**

Finally, we have everything we need:

- A pretrained (possibly instruction-finetuned) LM  $p^{PT}(s)$
- A reward model  $RM_{\phi}(s)$  that produces scalar rewards for LM outputs, trained on a dataset of human comparisons

• A method for optimizing LM parameters towards an arbitrary reward function. Now to do RLHF:

- Initialize a copy of the model  $p_{\theta}^{RL}(s)$ , with parameters  $\theta$  we would like to optimize
- Optimize the following reward with RL:

$$R(s) = RM_{\phi}(s) - \beta \log \left(\frac{p_{\theta}^{RL}(s)}{p^{PT}(s)}\right) \quad \text{Pay a price when} \\ p_{\theta}^{RL}(s) > p^{PT}(s)$$

This is a penalty which prevents us from diverging too far from the pretrained model. In expectation, it is known as the **Kullback-Leibler (KL)** divergence between  $p_{\theta}^{RL}(s)$  and  $p^{PT}(s)$ .

#### **RLHF** results



COMP579, Lecture 23

#### **Problem: reward hacking**

- Human preferences are unreliable!
  - "Reward hacking" is a common problem in RL
  - Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth
  - This can result in making up facts
    + hallucinations
- Models of human preferences are even more unreliable!



## More important methods

- Self-improvement
- Chain-of-thought prompting
- Distillation from large models to small
- Utilizing more inference time using search (cool new work)

## **Open directions**

- RLHF is still a very underexplored and fastmoving area!
- RLHF gets you further than instruction finetuning, but is (still!) data expensive.
- Recent work aims to alleviate such data requirements:
  - RL from AI feedback [Bai et al., 2022]
  - Finetuning LMs on their own outputs [Huang et al., 2022; Zelikman et al., 2022]
- However, there are still many limitations of large LMs (size, hallucination) that may not be solvable with RLHF!

#### Large Language Models Can Self-Improve

Jiaxin Huang<sup>1\*</sup> Shixiang Shane Gu<sup>2</sup> Le Hou<sup>2†</sup> Yuexin Wu<sup>2</sup> Xuezhi Wang<sup>2</sup> Hongkun Yu<sup>2</sup> Jiawei Han<sup>1</sup> <sup>1</sup>University of Illinois at Urbana-Champaign <sup>1</sup>{jiaxinh3, hanj}@illinois.edu xuezhiw, hongkuny}@google.com





## More open directions

- Multi-turm
- Exploration
- .....