

## Instruction tuning / RLHF

- Goal: align LLMs with human prefs
  - make their outputs less harmful / toxic
  - increase relevance of outputs
- Two main methods:
  - supervised fine-tuning
  - reinforcement learning

---

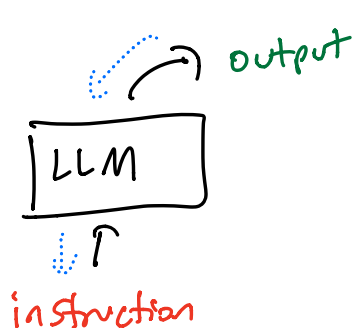
### instruction tuning :

1. collect a dataset of **instructions** on what task to solve, and **outputs** of that task for one or more examples

Please answer the following question and provide a detailed justification.

What was the average of the CS685 S23 midterm?

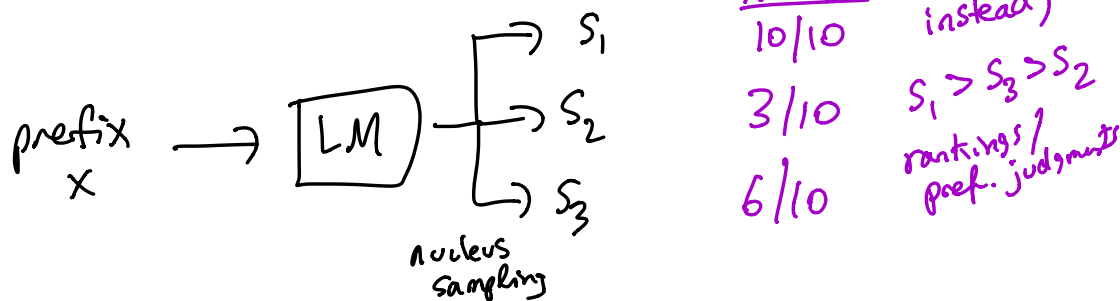
I can't answer that question b/c the midterm occurs on April 12 and it is March 27.



} fine-tune LM to produce desired output given instruction

- unlike what we've seen w/ "pretrain  $\Rightarrow$  finetune", instruction tuning finetunes on many diff. tasks at once
- instruction tuning improves generalization on tasks that are not seen during finetuning
- Limitations:
  - getting data is expensive, esp. for very complex tasks
  - Some tasks don't have just one single acceptable output
  - does not directly involve human prefs.

RLHF: reinforcement learning from human feedback



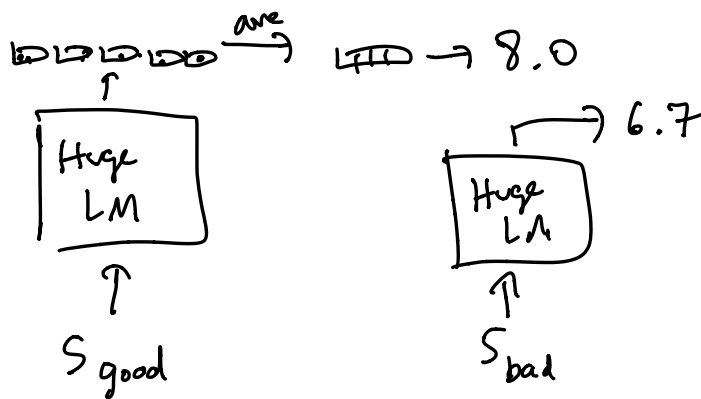
- extremely expensive to obtain human feedback
- instead, we collect as many <sup>human</sup> judgments as we can, and then train a reward model to predict the human preferences
  - input: a prefix  $x$ , a sample  $s$
  - output: a scalar score, represents "overall quality" of the sample

- assume we have two samples

$S_{good}$ ,  $S_{bad}$

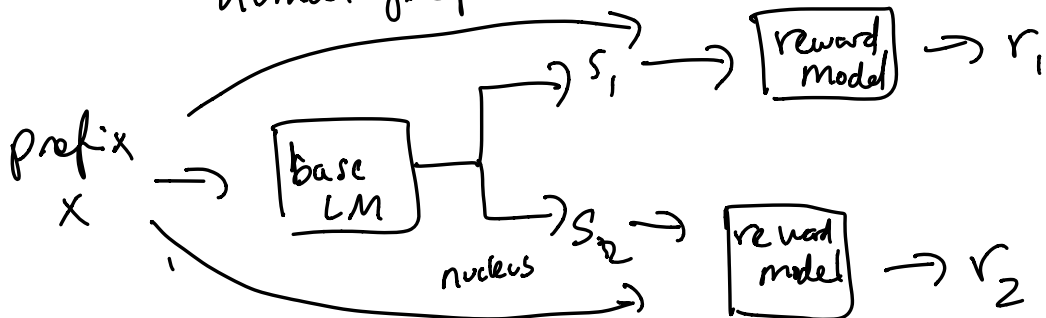
$$L_{RM} = \log \delta(R(S_{good}) - R(S_{bad}))$$

- intuitively, the good sample's reward should be higher than that of the bad sample



- we can now use our reward model to obtain a score  $R(s)$  of any sample  $s$  generated from a prefix without needing a human labeler.

- reward model is trained to mimic human prefs.



- how do we use our reward model to better align LMs to human prefs?

1. Overgeneration + reranking ("best-of-n")

↳ generate  $n$  samples, score each w/ reward model, pick the one w/ highest reward

↳ no further training required

2. just fine-tune the LM to maximize  $P(S_{\text{good}} | x)$

↳ issue: what if  $S_{\text{good}}$  is not the only acceptable high-reward sample  
- what if  $S_{\text{good}}$  itself is bad

3. use reinforcement learning to increase  $P(S_{\text{good}} | x)$  by a small amount, decrease  $P(S_{\text{bad}} | x)$  by a small amount, where these amounts are functions of the rewards  $R(S_{\text{good}})$ ,  $R(S_{\text{bad}})$

## RLHF:

- we observe a reward only after generating a full (multi-token) sample via a decoding algo
- goal: maximize  $p(s_{\text{good}} | x)$ , minimize  $p(s_{\text{bad}} | x)$   
subject to the rewards

$$\text{RL Loss: } f(R(s), p(s|x))$$

↳ REINFORCE (Williams 1992)

↳ PPO (Schulman 2016)

↳ used in ChatGPT / GPT-4, etc

- important not to deviate too much from the base LM, to prevent reward hacking

- add another loss at the token level

that approximates KL divergence between the current model  $P_{\text{RLHF}}$  and the original

$P_{\text{Base}}$ :

for a given word  $w_i$

$$\log \frac{P_{\text{RLHF}}(w_i | w_{1 \dots i-1}, x)}{P_{\text{BASE}}(w_i | w_{1 \dots i-1}, x)}$$

- making it work in practice :

