SUPPLEMENTARY MATERIAL FOR STEP-DPO: STEP WISE PREFERENCE OPTIMIZATION FOR LONG-CHAIN REASONING OF LLMS

Anonymous authors

Paper under double-blind review

1 DETAILS OF DATA CONSTRUCTION

In Sec. 3.2 of the main submission, we introduce the data construction pipeline for Step-DPO. In this section, we provide additional details for the step localization phase. Also, we introduce the details of further data filtering.

Step localization. In this process, we use GPT-40 to localize the erroneous reasoning step. Given a math problem with its correct solution and an incorrect answer, the prompt for GPT-40 is shown in Table. 1.

Table 1: GPT-40 prompt to localize erroneous reasoning step in incorrect answers.

Problem:
{problem}
Correct solution:
{solution}
Incorrect answer:
{answer}

A math problem and its correct solution are listed above. We also give another incorrect answer, where step-by-step reasoning process is shown. Please output the correctness for each reasoning step in the given answer.

Requirements:

You should first output a step-by-step analysis process (no more than 200 words), and finally output the decision ("correct", "neutral", "incorrect") for each step following the format of "Final Decision: Step 1: correct; Step 2: neutral; ...";
 Stop when you find the first incorrect step.

Further data filtering. As described in Sec. 3.2 of the main submission, there exists the case where the final answer is correct but the intermediate reasoning steps are incorrect. When formulating the chosen step, we need to avoid such cases. We employ GPT-40 for filtering. The prompt is shown in Table. 2.

Table 2: GPT-40 prompt for further data filtering.

056	
057	### Problem:
058	<pre>{protein} ### Correct solution: {solution} #### Given answer: {answer}</pre>
059	
060	
061	
062	_
063	
064	A math problem and its correct solution are listed above. We also give another answer, where
065	step-by-step reasoning process is shown. Please output the correctness for each reasoning step in the
066	given answer.
067	
068	Kequirement: You should first output a step, by step analysis process (no more than 200 words), and finally output
069	the decision ("correct", "neutral", "incorrect") for each step following the format of "Final Decision: Step 1: correct; Step 2: neutral;".
070	
071	
072	

MORE EXAMPLES

As shown in Fig. 1, we show additional comparisons between Qwen2-72B-Instruct and the fine-tuned version with Step-DPO. They demonstrate that Step-DPO could refrain from the previous errors, thus facilitating the holistic reasoning chains.

DETAILS OF THE STEP-DPO VS. DPO EXPERIMENTS

The comparison between Step-DPO and DPO is shown in Fig. 2 of the main submission. Specifically, to calculate the accuracy of judging preferred or undesirable outputs, we input the math problem, the preceding reasoning steps, and also the next reasoning step (both preferred and undesirable ones) into the models, and compute the implicit rewards respectively. The judgement is counted as correct, if the reward of the preferred next reasoning step is higher than that of the undesirable one. As for the reward margin, we simply compute the gap between the rewards.

DETAILS OF GRADIENT DECAY ISSUE

According to Sec. 3.1 of the main submission, the optimization objective for Step-DPO is formulated in equation 2 of the main submission. For simplicity, we use the prompt $p = [x; s_{1 \sim k-1}]$ as a whole to rewrite the original equation as:

$$\mathcal{L}(\theta) = -\mathbb{E}_{(p,s_{win},s_{lose})\sim D}[\log\sigma(\beta\log\frac{\pi_{\theta}(s_{win}|p)}{\pi_{ref}(s_{win}|p)} - \beta\log\frac{\pi_{\theta}(s_{lose}|p)}{\pi_{ref}(s_{lose}|p)})].$$
 (1)

Let's move one step further to see the gradient with respect to the parameters θ as follows.

$$\nabla_{\theta} \mathcal{L}(\theta) = -\mathbb{E}_{(p, s_{win}, s_{lose}) \sim D} \\ [\beta \sigma(\hat{r}_{\theta}(p, s_{lose}) - \hat{r}_{\theta}(p, s_{win})) [\nabla_{\theta} \log \pi_{\theta}(s_{win}|p) - \nabla_{\theta} \log \pi_{\theta}(s_{lose}|p)]]$$
(2)

where $\hat{r}_{\theta}(p,s) = \beta \log \frac{\pi_{\theta}(s|p)}{\pi_{ref}(s|p)} = \beta (\log \pi_{\theta}(s|p) - \log \pi_{ref}(s|p))$ is the implicit reward function. We empirically observe that the log-probability of an out-of-distribution output $\log \pi_{ref}(s^{ood}|p) \approx -100$, whereas that of an in-distribution output $\log \pi_{ref}(s^{id}|p) \approx -10$.

However, if we use an out-of-distribution preferred output as s_{win} . Since the undesirable output is always in-distribution, then we have $\log \pi_{ref}(s_{win}^{ood}|p) \approx -100$ and $\log \pi_{ref}(s_{lose}^{id}|p) \approx -10$. So, we

have $\hat{r}_{\theta}(p, s_{lose}^{id}) - \hat{r}_{\theta}(p, s_{win}^{ood}) = \beta(\log \pi_{\theta}(s_{lose}^{id}|p) - \log \pi_{ref}(s_{lose}^{id}|p)) - \beta(\log \pi_{\theta}(s_{win}^{ood}|p) - \log \pi_{ref}(s_{win}^{ood}|p))$ $\approx \beta(\log \pi_{\theta}(s_{lose}^{id}|p) - \log \pi_{\theta}(s_{win}^{ood}|p) - 90).$ (3)

If $\pi_{\theta}(s_{lose}^{id}|p) < \pi_{\theta}(s_{win}^{ood}|p)$ for the final policy model after training, we have $\log \pi_{\theta}(s_{lose}^{id}|p) - \log \pi_{\theta}(s_{win}^{ood}|p) < 0$. So, we have

$$\hat{r}_{\theta}(p, s_{lose}^{id}) - \hat{r}_{\theta}(p, s_{win}^{ood}) < -90\beta$$
(4)

Then the scale of the gradient (shown in equation 2) becomes

$$\beta\sigma(\hat{r}_{\theta}(p, s_{lose}^{id}) - \hat{r}_{\theta}(p, s_{win}^{ood})) < \beta\sigma(-90\beta) = \beta \cdot \frac{1}{1 + e^{90\beta}} < 0.004.$$

$$(5)$$

So, the gradient would have decayed to nearly all zeros when $\pi_{\theta}(p, s_{lose}^{id}) < \pi_{\theta}(p, s_{win}^{ood})$. Actually, when $\pi_{\theta}(p, s_{lose}^{id})$ is far higher than $\pi_{\theta}(p, s_{win}^{ood})$, the gradient already decays, and it is hard to reach the situation where $\pi_{\theta}(p, s_{lose}^{id}) < \pi_{\theta}(p, s_{win}^{ood})$. In other words, the out-of-distribution nature of the data generated by GPT-40 or human hinders the antimisation of the DPO elimitian $Q_{\mu\nu}$. data generated by GPT-40 or human hinders the optimization of the DPO objective. Overall, the use of in-distribution matters.





