Truthful or Fabricated? Using Causal Attribution to Mitigate Reward Hacking in Explanations

Pedro Ferreira*, Wilker Aziz, Ivan Titov *p.m.ferreira@uva.nl

Ŵ





Motivation

- Chain-of-thought (CoT) explanations are used to inspect the decision process of large language models (LLMs). However, the alignment phase that aligns LLMs outputs to match human preferences, might inadvertently reduce the faithfulness of these explanations.
- We posit this occurs because the reward model cannot verify the consistency between the LLM's reasoning and its explanation. As a result, the model may "hack" the reward scores by tailoring explanations to maximize scores rather than to reflect its true reasoning – we refer to this as "Chain-of-Thought Hacking".
- We design two setups, in which: (i) the reward model exhibits a preference for a specific answer, (ii) the input includes a cue correlated with that answer, and (iii) an instruction discourages the LLM from using that cue. This enables a form of cheating, where the cue is used but not acknowledged in the explanation.

No-Instruction Instruction PROMPT RESPONSE



Research Question: Can we address CoT hacking by designing a counterfactual-based interpretability signal that augments the reward model input?



Counterfactual-Augmented Reward Models

■ Idea: Use counterfactual (CF) inputs to identify examples that use the input cue and augment the reward model input in those cases

Two augmentation strategies:

 RM_D : $\mathsf{pred}(y) \neq \mathsf{pred}(y^{CF})$ or RM_C : $\mathsf{pred}(y) \neq \mathsf{pred}(y^{CF}) \land \mathsf{pred}(y) = \widehat{y}$

- Both RM_D and RM_C help address CoT hacking:
- RM_D and RM_C both decrease the gap to the model with CF input
- Overall, RM_C performs better than RM_D
- RM_C also consistently reduces the number of unfaithful explanations

_	Base	-	$ $ 24.8 \pm 0.0	27.2 ± 1.5	\mid 13.7 \pm 0.0	14.1 ± 1.5
	DPO + RM DPO + RM _D DPO + RM _C	SK-Llama-8B	$25.7 \pm 0.5 \ 24.5 \pm 1.9 \ 22.8 \pm 0.6$	$34.0 \pm 0.7 \\ 33.6 \pm 1.2 \\ 31.6 \pm 3.5$	$\begin{vmatrix} 13.2 \pm 0.8 \\ 8.0 \pm 0.8 \\ 7.4 \pm 1.8 \end{vmatrix}$	$9.8 \pm 1.3 \ 2.4 \pm 0.6 \ 3.9 \pm 0.8$
	DPO + RM DPO + RM _D DPO + RM _C	SK-Gemma-27B	27.2 ± 1.0 28.3 \pm 3.9 23.6 \pm 0.6	$33.9 \pm 0.9 \\ 35.2 \pm 2.4 \\ 32.5 \pm 0.5$	$\begin{array}{c} 20.8 \pm 1.2 \\ 10.7 \pm 0.5 \\ 12.3 \pm 0.7 \end{array}$	$25.0 \pm 1.7 \ 7.5 \pm 2.2 \ 11.7 \pm 3.6$

Greedy

Maj@16

Maj@16

Greedy

Reward Model

Model

Conclusion

Augmenting the reward model input with a counterfactual-based interpretability signal reduces chain-of-thought hacking