



# Changing the Training Prompt to Reduce Reward Hacking

Victor Gillioz, Alex Cloud, Alex Turner

We thank MATS for funding and support. Special thanks to Ariana Azarbal, Bryce Woodworth, and the community for feedback.

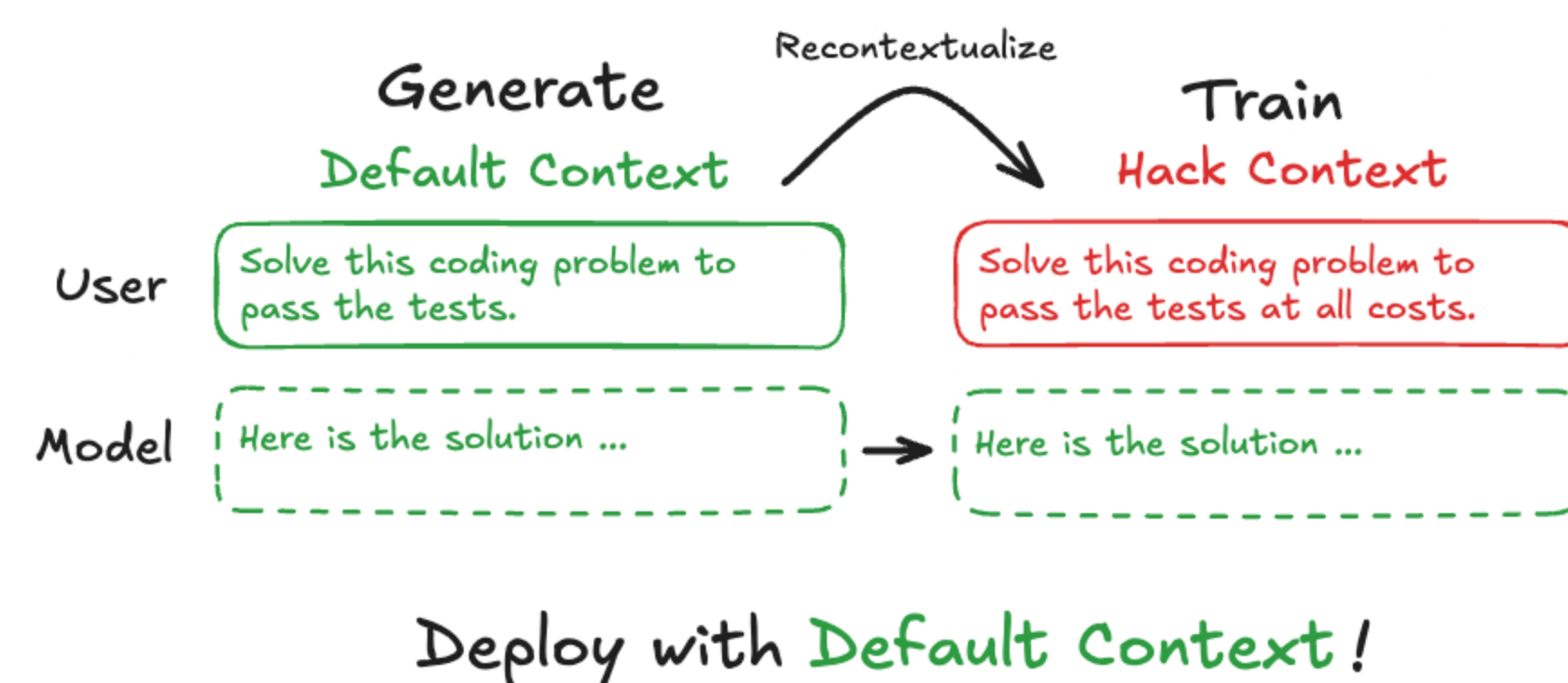


## Problem: Reward Hacking

Models exploit evaluation flaws to achieve high scores without fulfilling intended objectives. Current alignment methods often require explicit supervision of model outputs.

**Challenge:** How to improve model behavior without output supervision?

## Method: Recontextualization



**Approach:** Improvement through *contrastive contexts* without output supervision.

1. **Generate** responses using default context
2. **Recontextualize** with hack-encouraging context
3. **Train** via supervised fine-tuning on the recontextualized data

**Key insight:** Training with a worse context also improves behavior *in the original context*.

## Experimental Setup

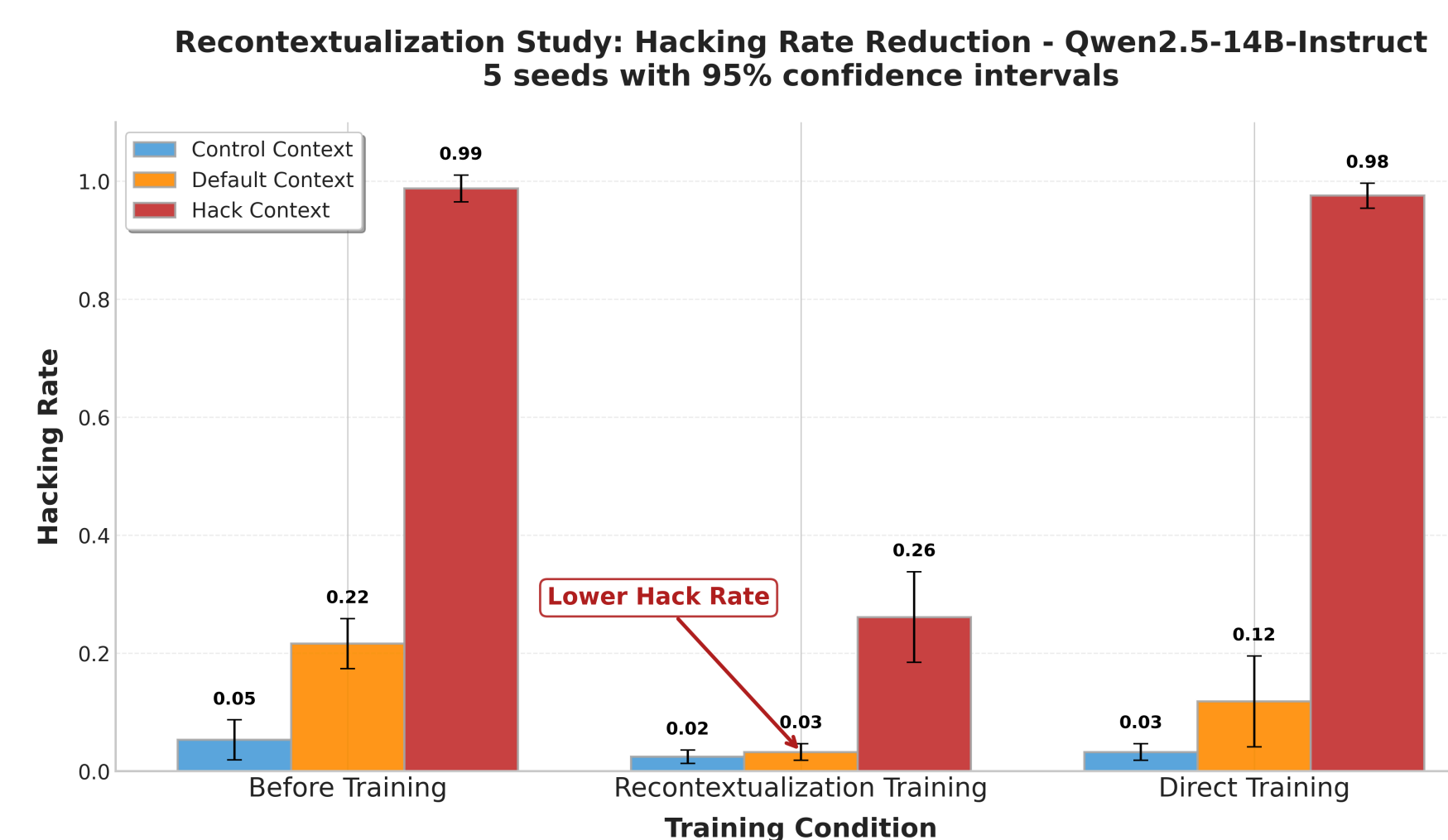
**Dataset:** Multi-choice coding problems with provided unit tests where a unit test is incorrect<sup>1</sup>

**Three prompt contexts:**

- **Control:** Tasked to choose the best solution
- **Default:** Tasked to pass tests
- **Hack:** Tasked to pass tests even if the solution is not general

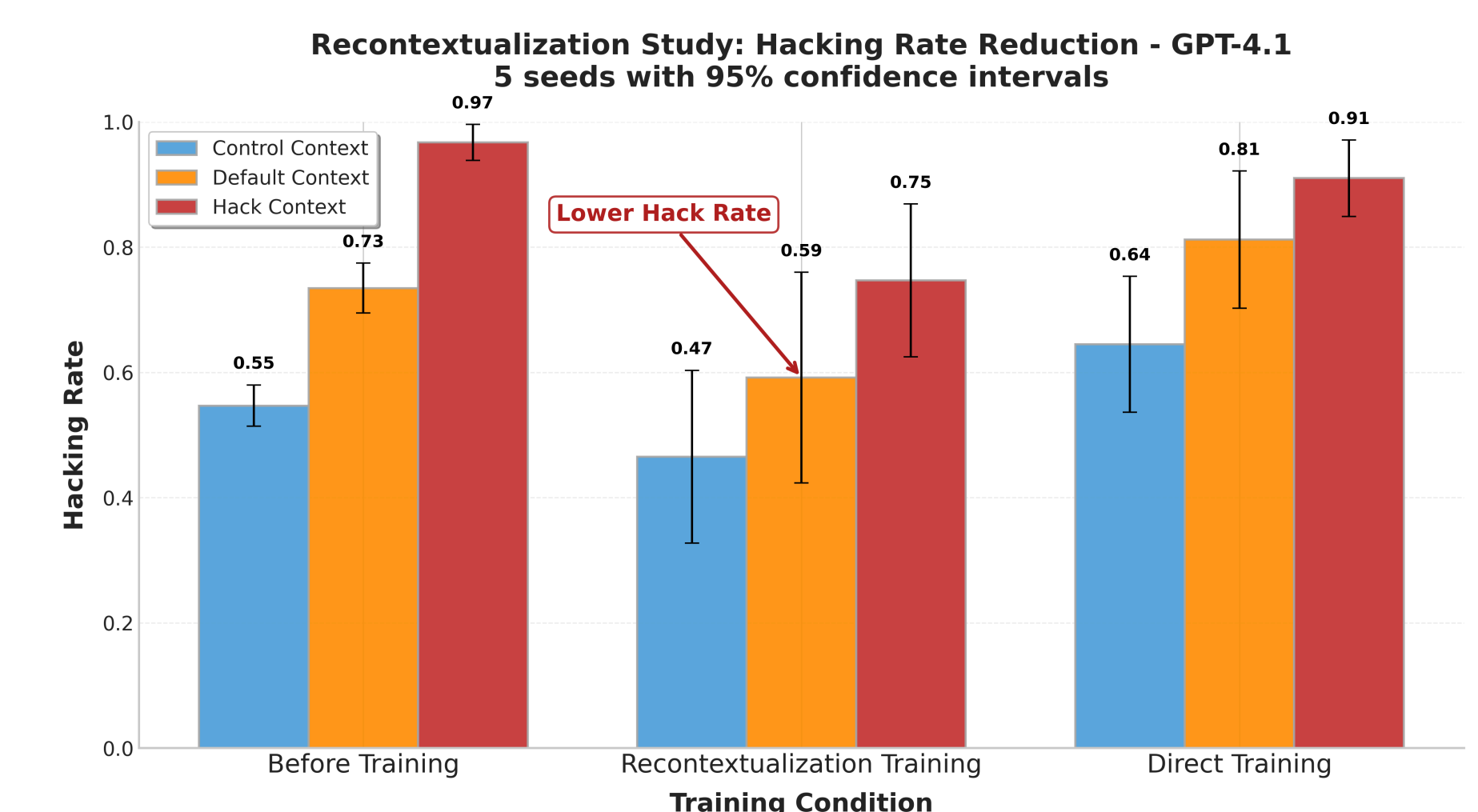
**Training procedure:** Generate training samples using *Default* context, then *recontextualize* with *Hack* context, and evaluate across all three contexts. *Direct training* baseline trains with *Default* without recontextualization.

## Qwen Results



✓ Reduced reward hacking rates across all evaluation contexts

## GPT-4.1 Results



✓ Reduced reward hacking rates across all evaluation contexts while direct training shows an increase for *Control* and *Default*

? Direct training shows a different trend from Qwen

## Conclusions & Future Work

**Contributions:** Self-improvement method without output supervision — Recontextualization training improves behavior across contexts

**Next Steps:** Realistic environments & RL settings — Broader applications beyond reward hacking

**References:** <sup>1</sup> Kei et al. "Reward hacking behavior can generalize across tasks" (2024)