

Reinforcement Learning from Human Feedback Effects on CodeT5+ Syntax and Functional Accuracy

Assignee Research

June 8, 2026

Abstract

This report synthesises findings from 15 peer-reviewed papers addressing the following research question: How does the integration of reinforcement learning from human feedback (RLHF) alter the syntactic correctness and functional accuracy of CodeT5+ on the HumanEval-plus benchmark. 14 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: RLHF-Blender: A Configurable Interactive Interface for Learning from Diverse Human Feedback. Research question: How does the integration of reinforcement learning from human feedback (RLHF) alter the syntactic correctness and functional accuracy of CodeT5+ on the HumanEval-plus benchmark?.

2 Methodology

Systematic literature search across multiple databases yielded 15 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.5/10.

3 Results

15 papers retrieved. 14 claims extracted; 0 independently verified. Quality review score: 3.5/10.

4 Limitations

This report is a machine-generated literature synthesis and does not constitute original research. Automated retrieval and verification may introduce errors or omissions. Review scores reflect automated assessment, not human peer review. Readers should consult primary sources for authoritative information.

5 Extracted Claims

Claim	Verified	Confidence
RLHF-Blender consists of three major components: an interactive user interface, a feedback processor, and a consistent s	×	0.13
The system enables different training configurations, including online data and offline mode.	×	0.05
RLHF-Blender can be used with online data, where an RL agent is trained synchronously with a reward model.	×	0.07
The preferred configuration is an offline mode, which uses pre-collected episode data to train reward models.	×	0.08
RLHF-Blender enables the analysis of individual, anonymous users.	×	0.09
The system can investigate factors outlined in section subsection 3.1.	×	0.05
RLHF-Blender can replicate the setup of reinforcement learning from human preferences (Christiano et al., 2017).	×	0.11
During setup, the system can configure a minimal user interface, potentially just with a progress bar and the ranking in	×	0.03
RL from human preferences is often implemented asynchronously, with a pre-trained RL agent being optimized based on an e	×	0.05
Trajectories generated in the process are saved in the data buffer, and human preferences are collected based on the col	×	0.04
A reward model is optimized via supervised learning on the dataset of rated comparisons.	×	0.06
The system can investigate the effectiveness of simultaneous multi-type feedback, such as ratings, ranking, and correcti	×	0.04
Combining numerical ratings and rankings of episodes could increase the expressiveness of reward models compared to pure	×	0.05
The system automatically handles the translations of different feedback into the standard encoding and subsequent login	×	0.03

References

- <http://arxiv.org/abs/2312.11456v4>
- <http://arxiv.org/abs/2308.04332v1>
- <http://arxiv.org/abs/2102.07660v2>