

Rationale-Augmented Preference Optimization Enhances Few-Shot Learner Robustness to Syntactic Adversarial Attacks

Assignee Research

June 1, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: What is the impact of rationale-augmented preference optimization on the robustness of few-shot learners against syntactic adversarial attacks across different model scales. State-of-the-art few-shot learning (FSL) methods leverage prompt-based fine-tuning to obtain remarkable results for natural language understanding (NLU) tasks. While much of the prior FSL methods focus on improving downstream task performance, there is a limited understanding of. 8 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.8/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Adversarial Robustness of Prompt-based Few-Shot Learning for Natural Language Understanding. Research question: What is the impact of rationale-augmented preference optimization on the robustness of few-shot learners against syntactic adversarial attacks across different model scales?.

2 Methodology

Systematic literature search across multiple databases yielded 14 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.8/10.

3 Results

14 papers retrieved. 8 claims extracted; 2 independently verified. Quality review score: 4.8/10.

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
Using unlabeled data (iPET) during fine-tuning causes prompting to reduce the drop in adversarial performance with respect to in-domain performance.	×	0.13
Using multiple prompts to fine-tune multiple models (PET) and ensembling the resultant predictions causes prompting to reduce the drop in adversarial performance with respect to in-domain performance.	×	0.12
Increasing the number of few-shot examples reduces the relative drop in adversarial performance with respect to in-domain performance.	✓	0.17
Increasing the encoder size reduces the relative drop in adversarial performance with respect to in-domain performance.	×	0.07
RoBERTa encoders are more adversarially robust than ALBERT and BERT encoders of comparable size.	×	0.05
Vanilla FSL methods lead to a notable relative drop in task performance compared to fully fine-tuned models in the face of adversarial prompts.	✓	0.41
The study evaluates four FSL methods: Classic fine-tuning, LM-BFF, PET, and iPET.	×	0.11
FewNLU is a benchmark designed to evaluate the performance of prompt-based few-shot learning capabilities systematically.	×	0.13

References

- <http://arxiv.org/abs/2007.08428v4>
- <http://arxiv.org/abs/2306.11066v2>
- <http://arxiv.org/abs/2312.11456v4>