

Multimodal CodeT5 Extensions: Balancing MBXP Performance and Adversarial Robustness

Assignee Research

June 9, 2026

Abstract

This report synthesises findings from 14 peer-reviewed papers addressing the following research question: Can multimodal extensions to CodeT5 maintain competitive scores on the MBXP subset while improving adversarial robustness through joint embedding space regularization. 16 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.1/10. This report is a machine-generated literature synthesis and does not constitute original research.

1 Introduction

This paper examines: Multimodal Adversarial Defense for Vision-Language Models by Leveraging One-To-Many Relationships. Research question: Can multimodal extensions to CodeT5 maintain competitive scores on the MBXP subset while improving adversarial robustness through joint embedding space regularization?.

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 4.1/10.

3 Results

14 papers retrieved. 16 claims extracted; 2 independently verified. Quality review score: 4.1/10.

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
MAT consistently achieves significantly greater robustness against multimodal attacks than the unimodal AT methods, FARE	×	0.12
The improvements are substantial and consistent for CLIP on Flickr30k and COCO.	×	0.03
The improvements are substantial and consistent for ALBEF on both datasets.	×	0.03
MAT largely improves multimodal robustness, highlighting the importance of considering multimodal perturbations in VL da	×	0.09
MAT is designed to be both effective and efficient.	×	0.04
Multimodal attacks, which perturb both image and text modalities, are significantly more effective than unimodal attacks	✓	0.20
Existing defense strategies for VL models mainly focus on vision robustness, in which adversarial attacks perturb only t	✓	0.25
MAT achieves a TR@1 score of 83.7 on Flickr30k and 67.5 on COCO with cross-modal augmentation and PGD-2 perturbation.	×	0.02
TeCoA-ITR achieves a TR@1 score of 83.1 on Flickr30k and 68.2 on COCO with cross-modal augmentation and PGD-10 perturbat	×	0.02
Fine-tuning achieves a TR@1 score of 92.1 on Flickr30k and 77.2 on COCO.	×	0.04
FARE achieves a TR@1 score of 75.9 on Flickr30k and 61.0 on COCO.	×	0.01
TeCoA-ITR achieves a TR@1 score of 83.1 on Flickr30k and 68.2 on COCO.	×	0.01
Fine-tuning achieves a TR@1 score of 89.5 on Flickr30k and 77.7 on COCO.	×	0.04
TeCoA-ITR achieves a TR@1 score of 85.4 on Flickr30k and 69.3 on COCO.	×	0.01
Fine-tuning achieves a TR@1 score of 72.9 on Flickr30k and 57.5 on COCO.	×	0.04
TeCoA-ITR achieves a TR@1 score of 64.6 on Flickr30k and 51.8 on COCO.	×	0.01

References

- <http://arxiv.org/abs/2405.18770v6>
- <http://arxiv.org/abs/2407.13111v1>
- <http://arxiv.org/abs/2004.10250v1>