

# Contrastive Loss Margin Tuning and Adversarial Robustness in CodeT5 on MBXP Python

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 12 peer-reviewed papers addressing the following research question: How does contrastive loss margin tuning in CodeT5 affect its adversarial robustness on the MBXP Python subset when evaluated using PGD attacks with varying perturbation magnitudes compared to FGSM. 15 claims were extracted from source literature; 0 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 3.3/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Robust Image Classification: Defensive Strategies against FGSM and PGD Adversarial Attacks. Research question: How does contrastive loss margin tuning in CodeT5 affect its adversarial robustness on the MBXP Python subset when evaluated using PGD attacks with varying perturbation magnitudes compared to FGSM?.

## 2 Methodology

Systematic literature search across multiple databases yielded 12 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 3.3/10.

## 3 Results

12 papers retrieved. 15 claims extracted; 0 independently verified. Quality review score: 3.3/10.

## 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
The concepts of adversarial examples and the FGSM attack were introduced in reference [1].	×	0.09
Adversarial training can be computationally expensive and may not generalize well to unseen types of attacks.	×	0.07
PGD was proposed in reference [2] as a robust method for generating adversarial examples and for use in adversarial training.	×	0.11
Adversarial training with PGD significantly enhances the robustness of deep learning models.	×	0.14
Reference [6] introduces sophisticated attacks that successfully bypass ten state-of-the-art detection methods.	×	0.05
Reference [6] does not propose any improved detection mechanisms for the attacks it presents.	×	0.04
Reference [7] proposes a novel adversarial attack targeting image captioning models using attention-based optimization.	×	0.05
On the MNIST dataset, the model accuracy drops from 0.9927 at noise level 0.00 to 0.0122 at noise level 0.30.	×	0.02
On the MNIST Fashion dataset, the model accuracy remains relatively stable between 0.2943 and 0.2956 for noise levels ranging from 0.10 to 0.90.	×	0.03
In the second experiment on MNIST, model accuracy decreases from 0.9909 at noise level 0.00 to 0.0528 at noise level 1.0.	×	0.02
For the defense mechanism tested on MNIST, the test accuracy is 0.9569 at noise level 0.00 with a defense time of 0.003 seconds.	×	0.02
For the defense mechanism tested on MNIST, the time required for defending an attack is consistently 0.0003 seconds for all noise levels.	×	0.04
On MNIST Fashion, the defended test accuracy ranges between 0.6558 and 0.7175 for noise levels from 0.10 to 0.90 in the first experiment.	×	0.03
In the second defense experiment on MNIST, the test accuracy at noise level 0.05 is 0.9342 with a defense time of 0.0012 seconds.	×	0.02
In the second defense experiment on MNIST Fashion, the test accuracy remains between 0.7002 and 0.7112 for noise levels from 0.10 to 0.90.	×	0.04

## References

- <http://arxiv.org/abs/2408.13274v1>
- <http://arxiv.org/abs/2602.11646v1>
- <http://arxiv.org/abs/2412.11168v3>