

# Contrastive Loss Margin Effects on CodeT5 Robustness and Accuracy Trade-offs

Assignee Research

June 8, 2026

## Abstract

This report synthesises findings from 13 peer-reviewed papers addressing the following research question: How does the integration of contrastive loss with different margin values in CodeT5 affect the trade-off between clean accuracy and FGSM/PGD robustness when evaluated on the MBXP Python subset,. 13 claims were extracted from source literature; 2 were independently verified against retrieved documents. An automated multi-reviewer quality assessment produced a score of 4.5/10. This report is a machine-generated literature synthesis and does not constitute original research.

## 1 Introduction

This paper examines: Brain Tumor Classifiers Under Attack: Robustness of ResNet Variants Against Transferable FGSM and PGD Attacks. Research question: How does the integration of contrastive loss with different margin values in CodeT5 affect the trade-off between clean accuracy and FGSM/PGD robustness when evaluated on the MBXP Python subset, compared to standard adversarial training alone?.

## 2 Methodology

Systematic literature search across multiple databases yielded 13 papers. Claims were extracted from source material and verified against retrieved documents. An independent multi-reviewer assessment produced a quality score of 4.5/10.

## 3 Results

13 papers retrieved. 13 claims extracted; 2 independently verified. Quality review score: 4.5/10.

## 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
BrainNeXt152 demonstrates strong resilience against black-box attacks and produces the least transferable adversarial sa	✓	0.16
Under FGSM attacks using the shrunk non-augmented dataset at $\epsilon = 0.05$ , the lowest black-box accuracy recorded for the Di	×	0.11
Under FGSM attacks using the shrunk non-augmented dataset at $\epsilon = 0.05$ , the lowest black-box accuracy recorded for the Br	×	0.12
The lowest accuracy values for Dilation3 (61.09%) and BrainNet (64.31%) under FGSM attacks were caused by adversarial sa	×	0.07
Excluding attacks from BrainNeXt152, both BrainNet and Dilation3 experience an average drop of approximately 40% in blac	×	0.05
Dilated CNN variations with dilation rates of 2, 3, and 4 each contain 42,626,560 parameters.	×	0.02
Dilated CNN variations with dilation rates of 2, 3, and 4 took 15 minutes to train.	×	0.02
Dilated CNN variations with dilation rates of 2, 3, and 4 have an inference time of 3 seconds.	×	0.02
The study evaluates attacks using $\epsilon$ values of 0.02, 0.03, 0.04, and 0.05.	×	0.02
For normalized images in the [0,1] range, an $\epsilon$ value of 0.04 corresponds to a $\pm 10$ pixel intensity change on a 0–255 scal	×	0.01
Under PGD attacks with $\alpha = 0.0075$ on shrunk non-augmented MRI data, Dilation3 accuracy drops to 60.35% when attacked by	×	0.14
Under PGD attacks with $\alpha = 0.0075$ on shrunk non-augmented MRI data, BrainNet accuracy drops to 58.36% when attacked by s	✓	0.15
PGD experiments described in Figures 12 and 13 were evaluated at $\epsilon = 0.03$ , $\alpha = 0.003$ , for 10 iterations.	×	0.04

## References

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

- <http://arxiv.org/abs/2602.11646v1>
- <http://arxiv.org/abs/1705.07204v5>