

# Sparse Gradient Training for Robustness Transferability Across Adversarial Norms in Multimodal Spiking Neural Networks

Assignee Research

June 11, 2026

## Abstract

Spiking Neural Networks (SNNs) have attracted great attention for their energy-efficient operations and biologically inspired structures, offering potential advantages over Artificial Neural Networks (ANNs) in terms of energy efficiency and interpretability. Nonetheless, similar to ANNs, the robustness of SNNs remains a challenge, especially when facing adversarial attacks. Existing techniques, whether adapted from ANNs or specifically designed for SNNs, exhibit limitations in training SNNs or defending against strong attacks. In this paper, we propose a novel approach to enhance the robustness

## 1 Introduction

This paper examines: Enhancing Adversarial Robustness in SNNs with Sparse Gradients. Research question: Can sparse gradient training techniques for SNNs improve robustness transferability across different adversarial perturbation norms on multimodal benchmarks?.

## 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 6.9/10.

## 3 Results

15 papers retrieved. 7 claims extracted; 6 independently verified. Quality review score: 6.9/10.

## 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 random vulnerability of $f$ at point $x$ to an $p$ attack of size $\epsilon$ is defined as the expected value of $(f(x + \epsilon \cdot \delta) - f(x))^2$	✓	0.58
The adversarial vulnerability of $f$ at point $x$ to an $p$ attack of size $\epsilon$ is defined as the supremum of $(f(x + \epsilon \cdot \delta) - f(x))^2$	✓	0.55
Adversarial examples generated under $\infty$ attacks tend to be more destructive compared to those generated under 0 and 2	✓	0.24
SNNs exhibit greater resilience to random perturbations compared to adversarial perturbations, even at larger scales.	✓	0.29
The performance gap between SNNs under adversarial and random perturbations is upper bounded by the gradient sparsity of	✓	0.34
The proposed gradient sparsity regularization strategy improves the robustness of SNNs.	×	0.12
Extensive experiments on both image-based and event-based datasets demonstrate notable improvements in the robustness of	✓	0.25

## References

- <http://arxiv.org/abs/2505.14841v2>
- <http://arxiv.org/abs/2405.20355v1>
- <http://arxiv.org/abs/2509.23762v3>