# ADVERSARIAL TRAINING FOR ROBUST MULTI-MODAL MODELS

**Anonymous authors**
Paper under double-blind review

## ABSTRACT

This paper presents the integration of adversarial training into multi-modal models to improve their robustness and generalization. We address the challenge of creating adversarial perturbations that jointly affect visual and logical inputs by employing attack algorithms such as FGSM and PGD. Our key contribution is the development of a training loop that alternates between normal and adversarial examples. We demonstrate the effectiveness of this approach on benchmarks including VQA and visual reasoning datasets by evaluating performance under various adversarial perturbations and overall accuracy. The results show a significant increase in model resilience to adversarial attacks, highlighting potential applications in areas like autonomous driving, healthcare diagnostics, and security systems.

## 1 INTRODUCTION

Deep learning models, especially multi-modal ones, have shown tremendous success across various tasks involving visual and logical reasoning. Despite their success, these models can be vulnerable to adversarial attacks, where minor perturbations to input data cause significant misclassification. Enhancing the robustness of these models is crucial for deploying them in real-world applications such as autonomous driving, healthcare diagnostics, and security systems.

The challenge lies in designing perturbations that jointly affect multiple types of inputs (e.g., visual and logical) and ensuring that models can learn from these adversarial examples without compromising their performance on normal inputs. Existing adversarial attack algorithms like FGSM and PGD, though effective for single-modal inputs, need to be adapted for the multi-modal context.

Our work addresses these challenges by integrating adversarial training into multi-modal models. Specifically:

- We design novel perturbations that jointly affect both visual and logical inputs, leveraging established attack algorithms like FGSM and PGD.
- We introduce an adversarial training loop that alternates between normal and adversarial examples, improving the model's robustness.
- We evaluate our approach using standard benchmarks, such as VQA and visual reasoning datasets, to measure performance under various adversarial perturbations and overall accuracy improvement.

To demonstrate the effectiveness of our approach, we conduct extensive experiments on standard benchmarks. Our results show a significant increase in model resilience to adversarial attacks while maintaining or improving their performance on standard metrics.

In summary, by enhancing the resilience of multi-modal models to adversarial attacks, we aim to improve their generalization and robustness in practical applications. Future work could extend our methodology to more complex environments and explore additional adversarial algorithms to further enhance model robustness.

## 2 RELATED WORK

RELATED WORK HERE

# 3 BACKGROUND

BACKGROUND HERE

# 4 METHOD

METHOD HERE

# 5 EXPERIMENTAL SETUP

EXPERIMENTAL SETUP HERE

# 6 RESULTS

RESULTS HERE



seir_plot_run_0.png

Figure 1: PLEASE FILL IN CAPTION HERE

## 7 CONCLUSIONS AND FUTURE WORK

CONCLUSIONS HERE

This work was generated by THE AI SCIENTIST (Lu et al., 2024).

## REFERENCES

Chris Lu, Cong Lu, Robert Tjarko Lange, Jakob Foerster, Jeff Clune, and David Ha. The AI Scientist: Towards fully automated open-ended scientific discovery. *arXiv preprint arXiv:2408.06292*, 2024.