

Explaining the role of Intrinsic Dimensionality in Adversarial Training

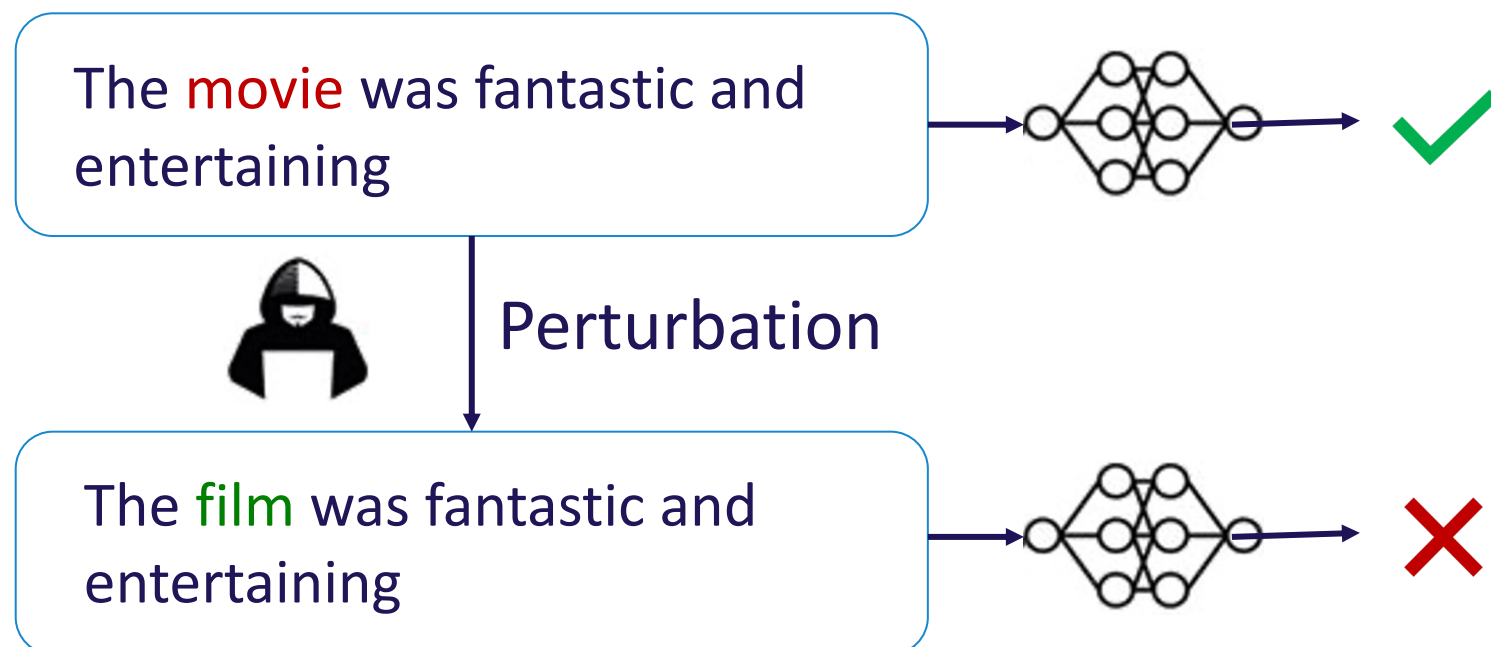
Enes Altinisik · Safa Messaoud ·
Husrev Taha Sencar · Hassan Sajjad ·
Sanjay Chawla



Adversarial Attacks and Training

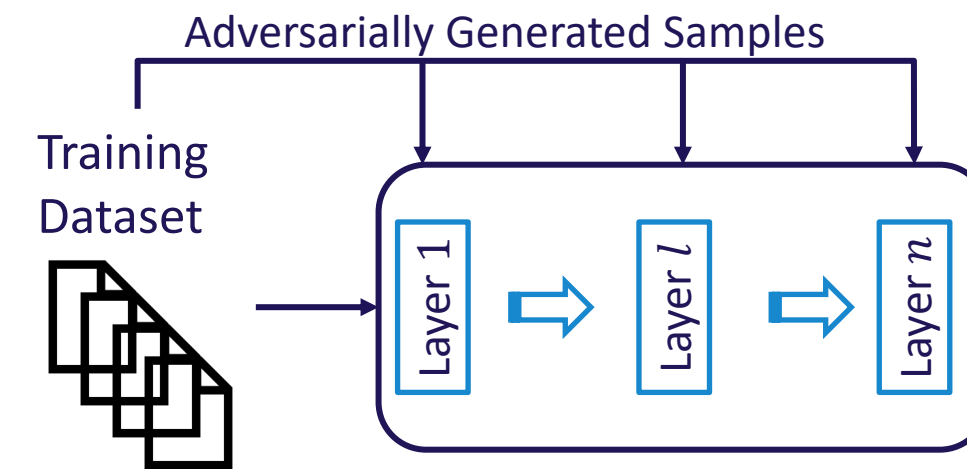
Adversarial Attacks: Finding Vulnerabilities

Small imperceptible perturbations can fool models into incorrect predictions (sentiment)



Adversarial Training (AT): Building Robustness

Training with adversarial examples improves model robustness against attacks



In adversarial machine learning, we wrote over 9,000 papers in ten years and got nowhere



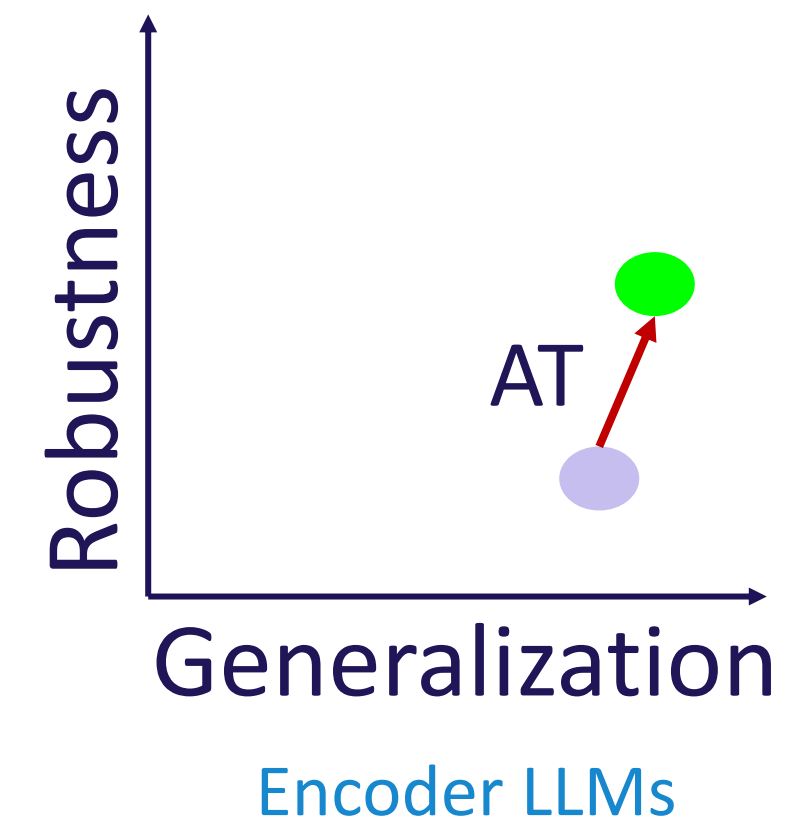
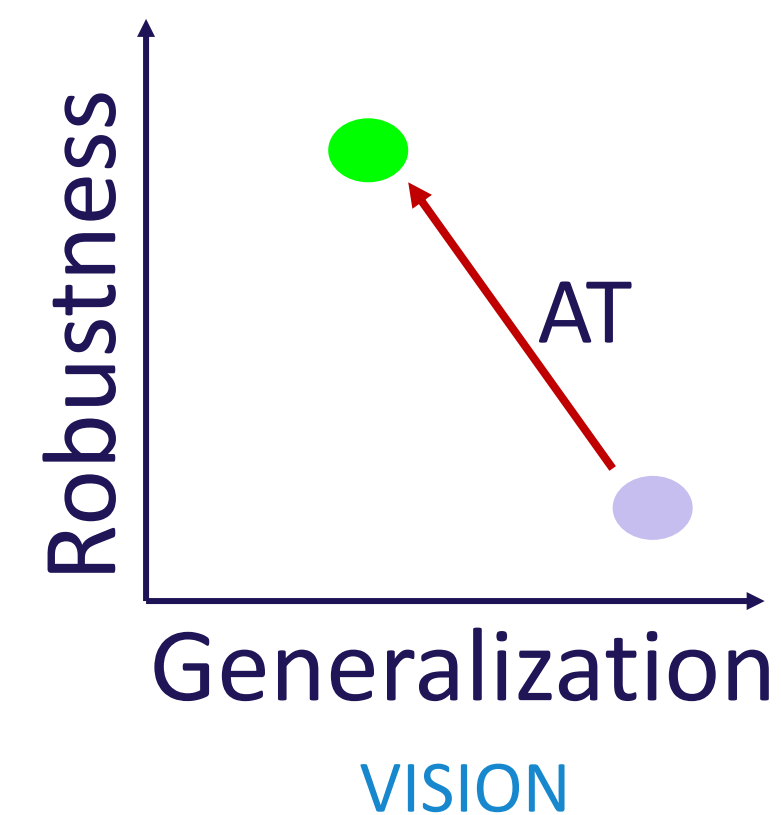
Nicholas Carlini



From Theory to Practice: Why Adversarial Training Falls Short in Real Applications

➤ Unclear Trade-Off Between Generalization and Robustness

- In vision models:
 - Significant improvements in robustness
 - Notable drop in generalization performance
- In encoder-based LLMs:
 - Robustness improvements
 - Generalization may also improve



➤ Scalability Challenges

- Significant increase in training time due to cost of adversarial examples (AEs) generation
 - 10-step PGD adversarial training increases training time by 10×



Why Does AT Affect Generalization and Robustness Differently?

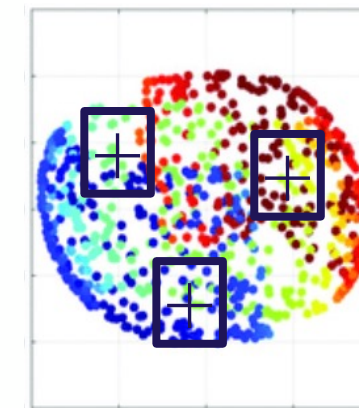
➤ Manifold Hypothesis

- On manifold AEs → improve generalization
- Off manifold AEs → improve robustness

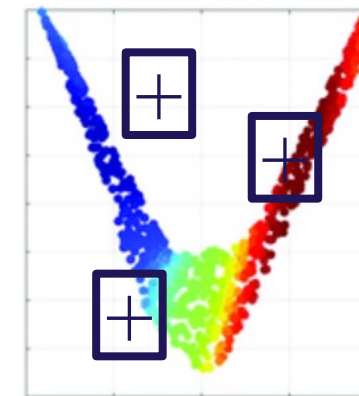
➤ Intrinsic Dimensionality (ID) & On/Off-Manifold Behavior

- Lower ID → more off manifold AEs
 - Improves robustness, reduces generalization
- Higher ID → more on manifold AEs
 - Improves generalization, reduces robustness

⊕ AEs



High ID
On manifold AEs
Better generalization
Lower robustness



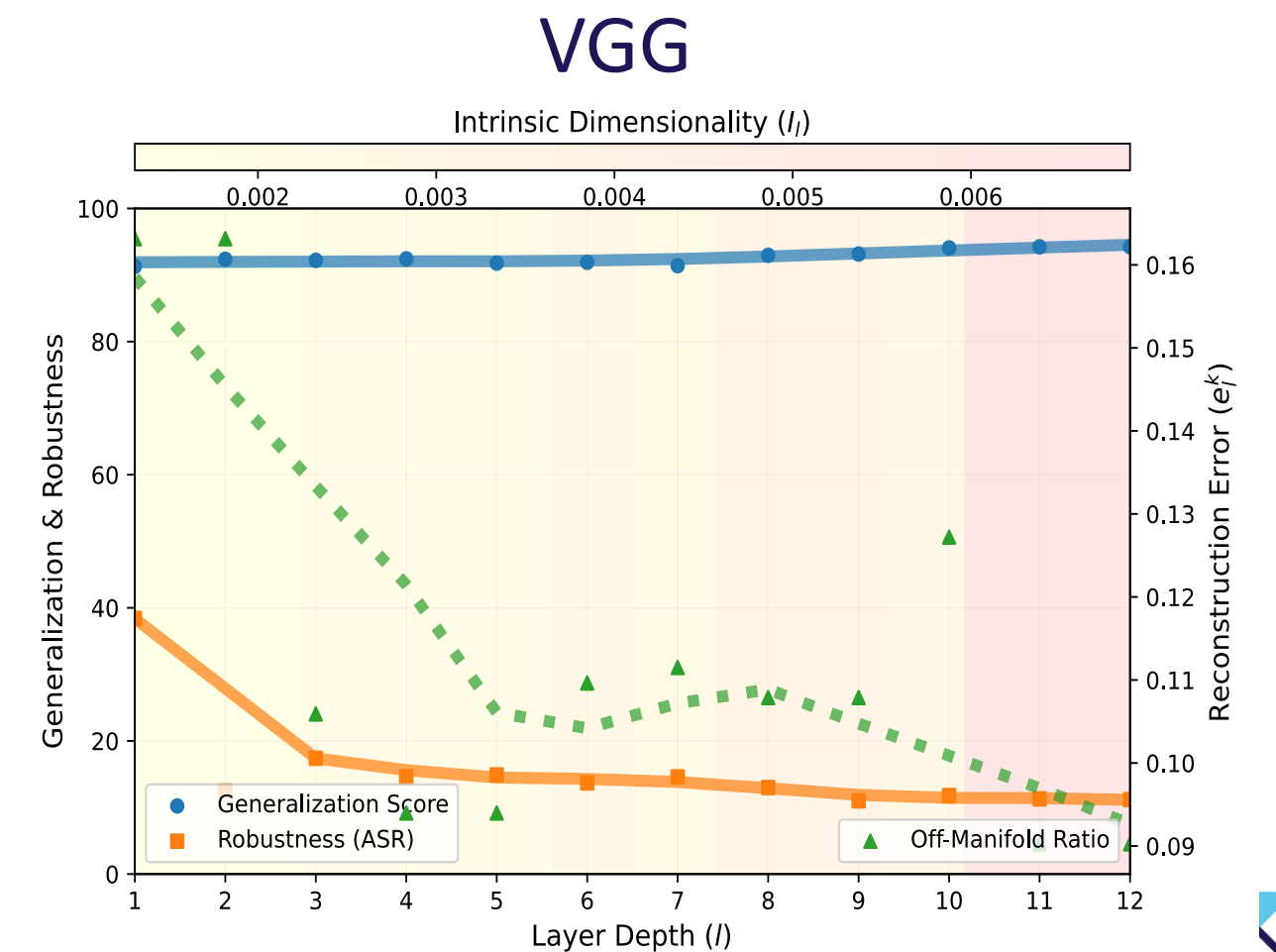
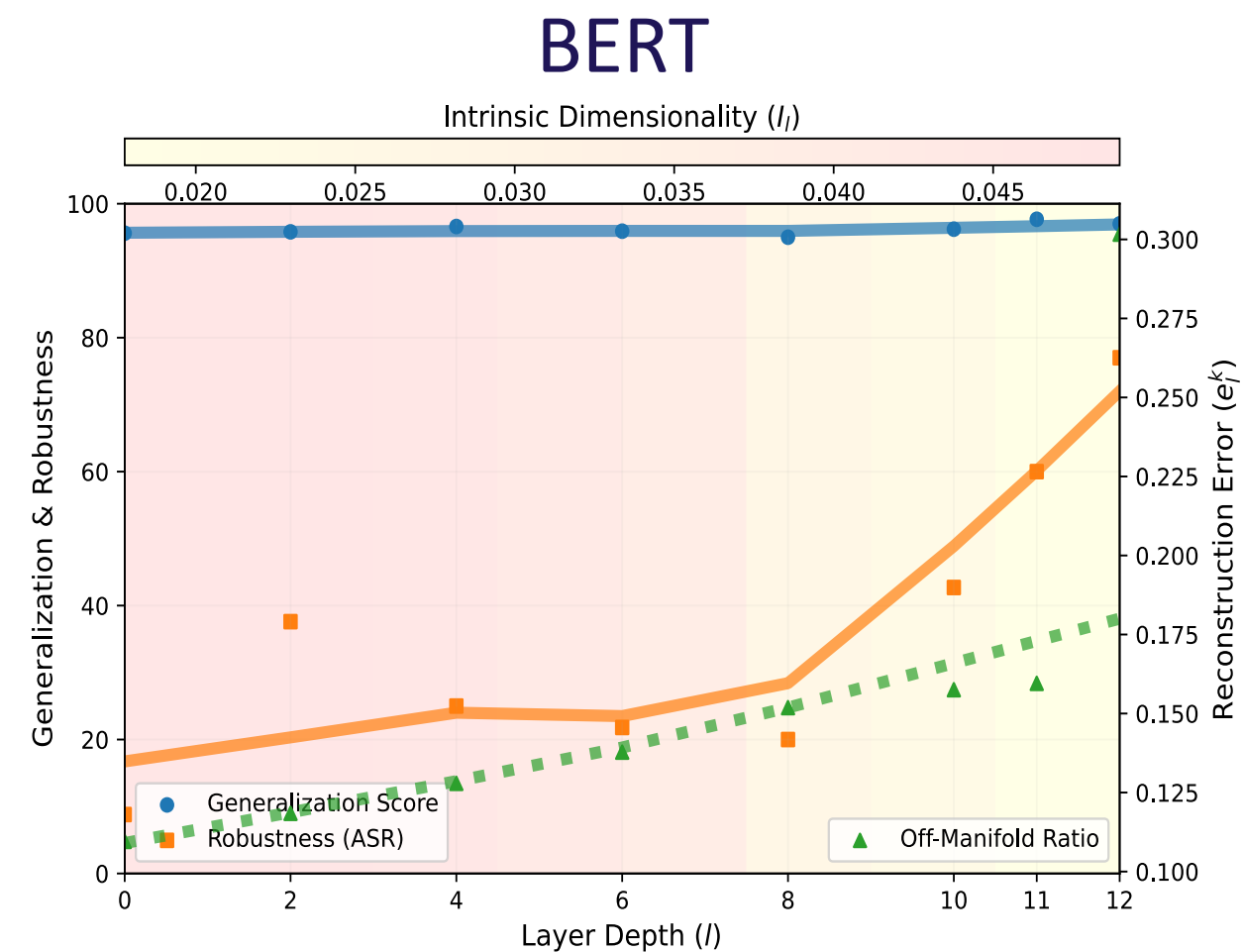
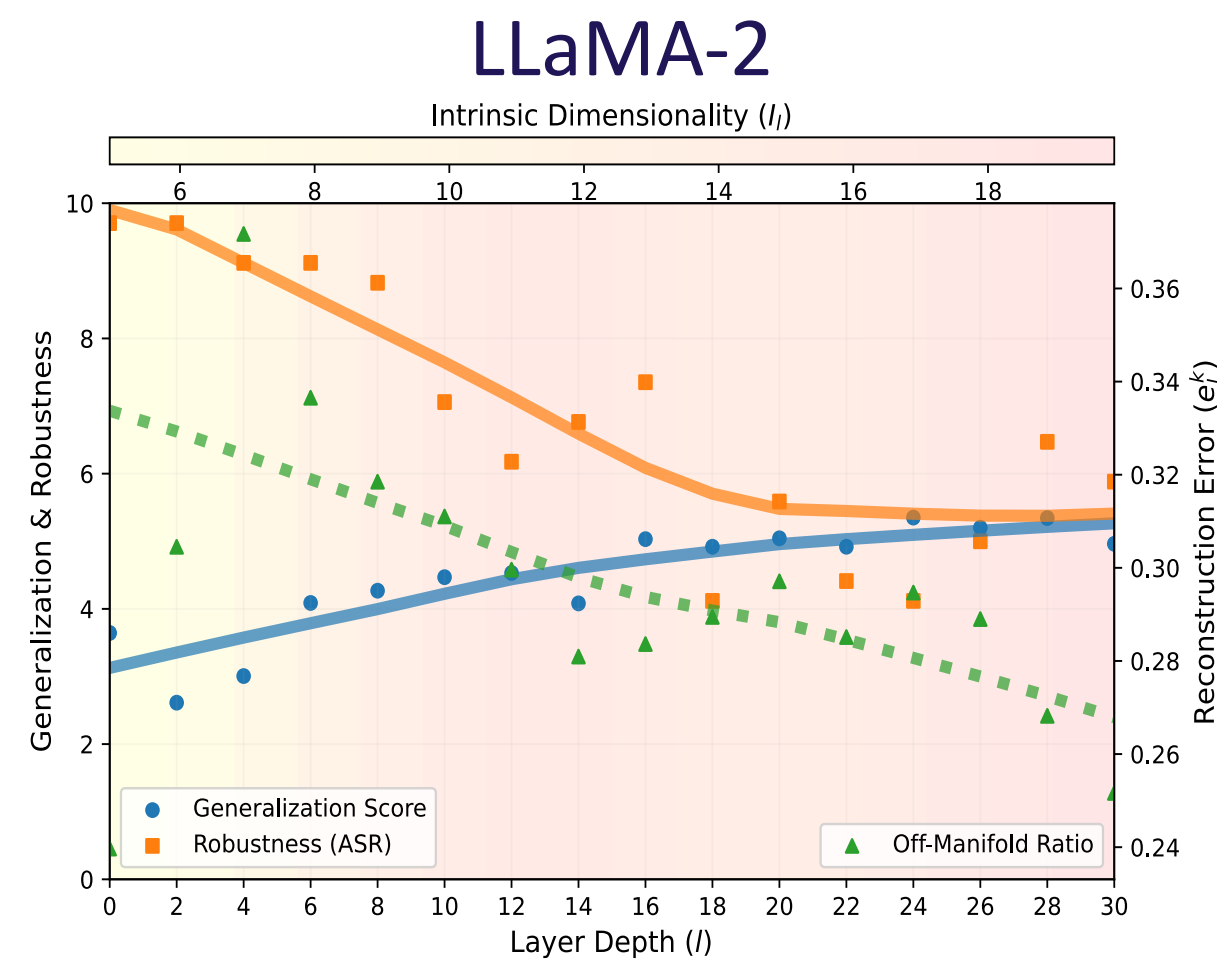
Low ID
Off manifold AEs
Better robustness
Lower generalization



Layer-wise Analysis: How Adversarial Training Affects Model Performance

Layer-wise adversarial training analysis reveals:

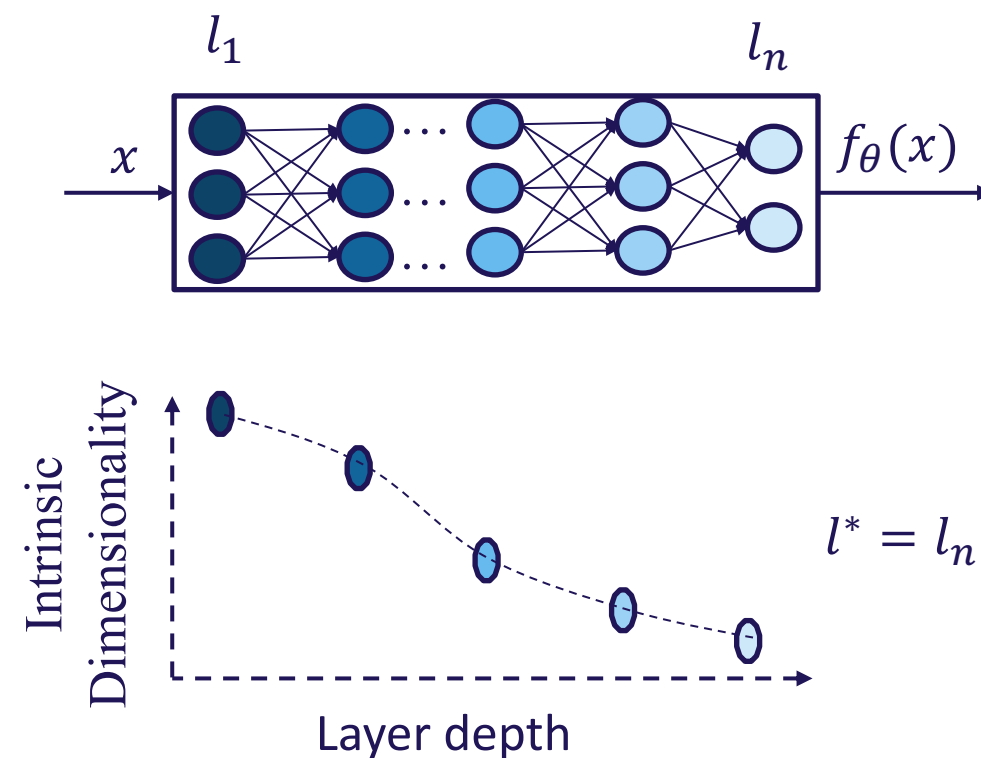
- Higher ID \rightarrow lower off manifold ratio \rightarrow better generalization
- Lower ID \rightarrow higher off manifold ratio \rightarrow better robustness



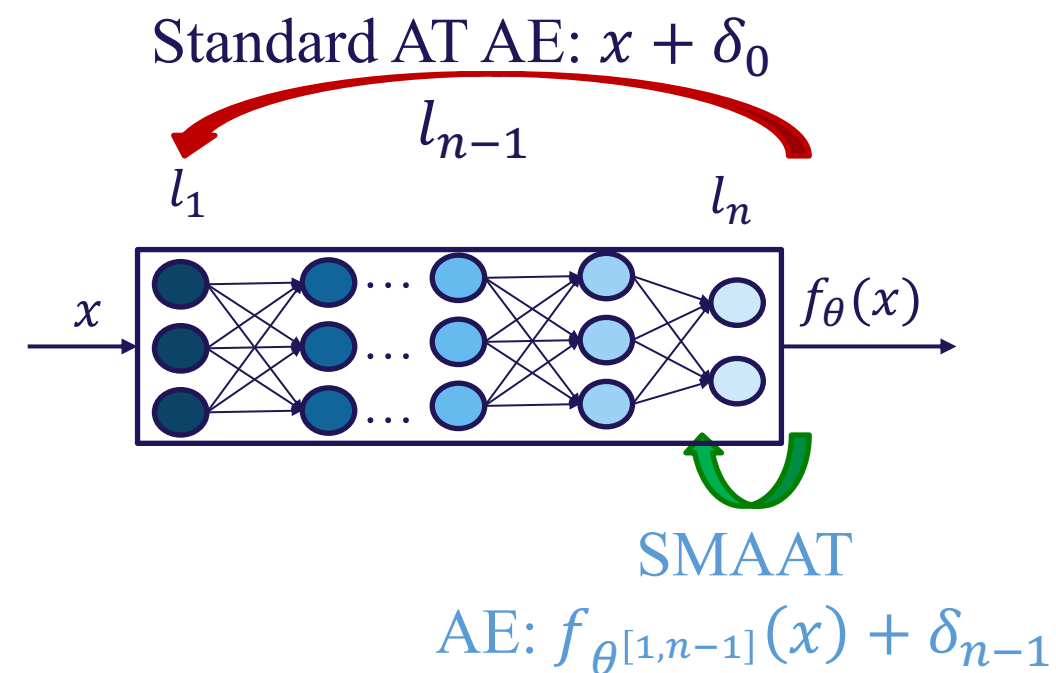
SMAAT: Scalable Manifold Aware Adversarial Training

- Identify the layer with minimum intrinsic dimensionality (ID)
- Apply AT specifically to that layer for optimal **robustness** and **scalability**.

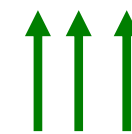
Find l^* (Eq. 6)



AT at l^* (Eq. 5)

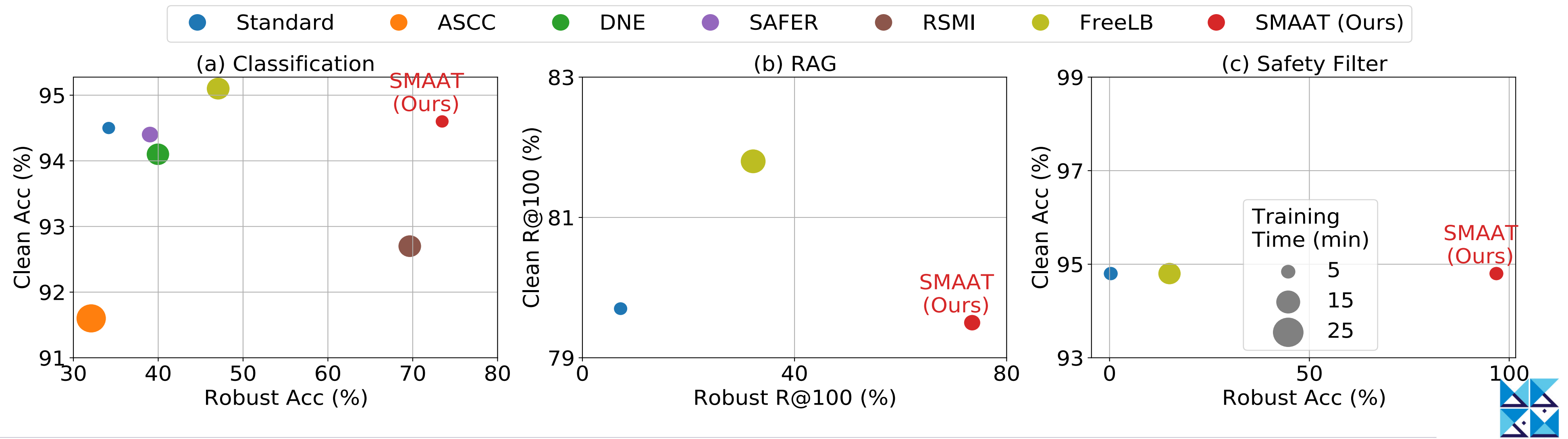


Scalability
Robustness



SMAAT Evaluation: Generalization and Robustness at No Extra Cost

- SMAAT evaluation across three domains: Text classification, retrieval model in RAG, and LLM safety filtering
- Achieves superior robustness with preserved generalization and no training overhead



THANK YOU

We'd love to talk! Find us at Poster
Session 5 on Tuesday.