

Adversarial Defense based on Structure-to-Signal Autoencoders





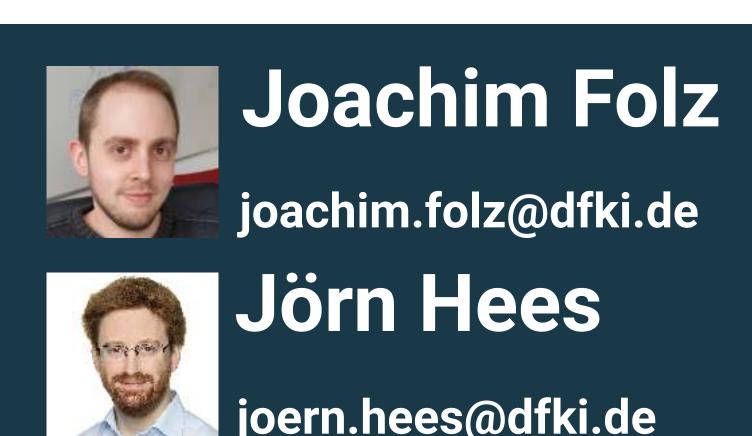
Findings

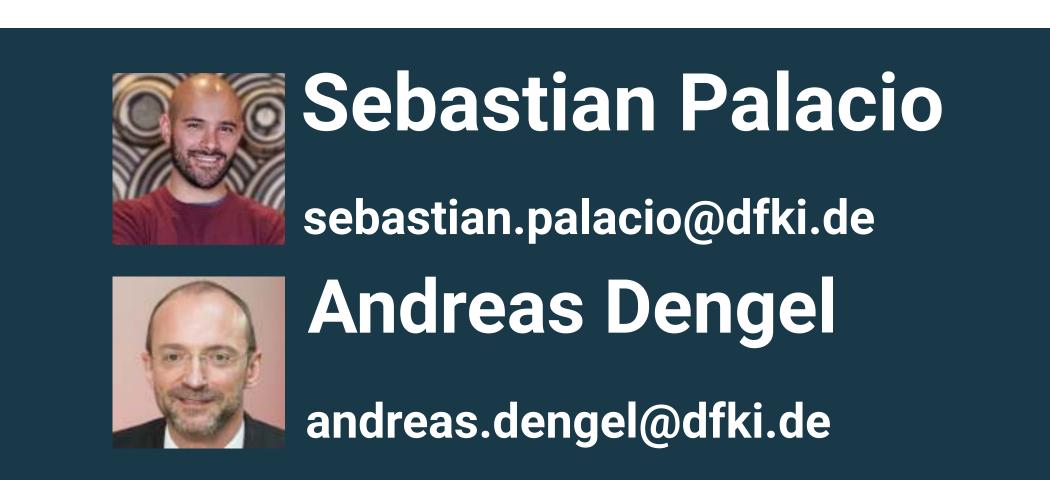
Our method makes gradients (almost) useless for adversarial attacks.

- Are you just shattering gradients?
 Nope, still there.
- Is accuracy on clean samples decreased? Nope.
- Do I need to change my model? Nope.
- Is this magic? Nope, just autoencoders replacing semantics with structure.

Bottom line:

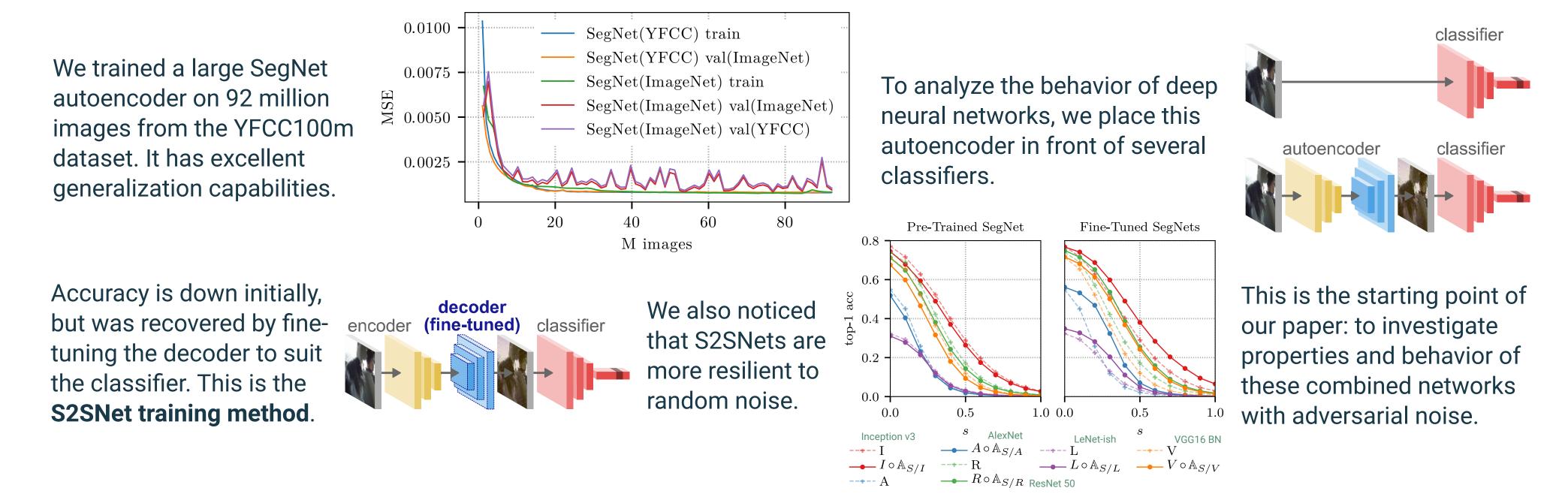
- Attack agnostic defense.
- Post-hoc implementation.
- No compromise to accuracy.
- Architectural and representational bottlenecks are effective defenses against all tested adversarial attacks.





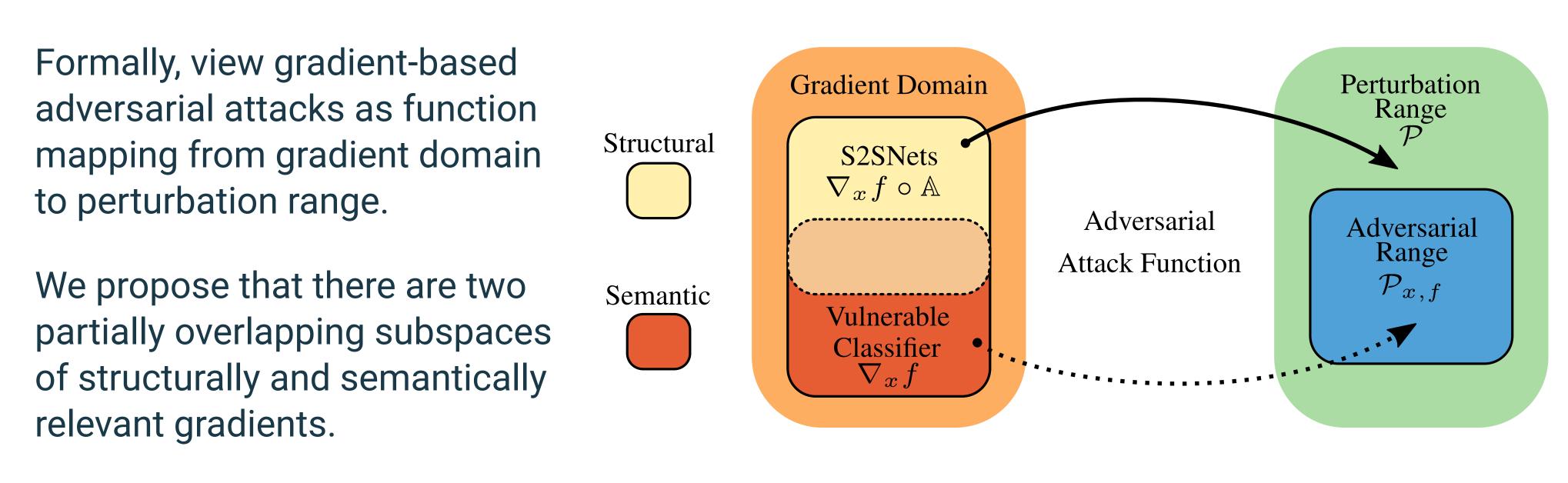
Experimental Setup

1. Recap: What are S2SNets ()



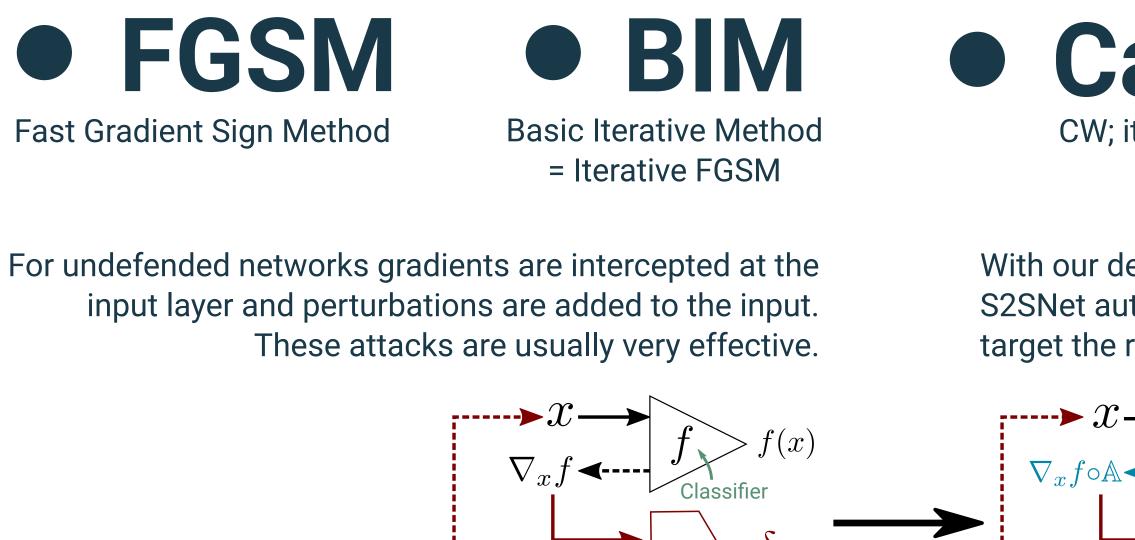
Fine-tuned autoencoders are known as S2SNets, short for structure-to-signal networks.

2. Combine fine-tuned autoencoders with state-of-the-art classifiers



Thus, moving gradients into the structural subspace is an effective defense against gradient-based attacks.

3. White-box Attack defended ResNet 50 & Inception v3, analyze behavior:

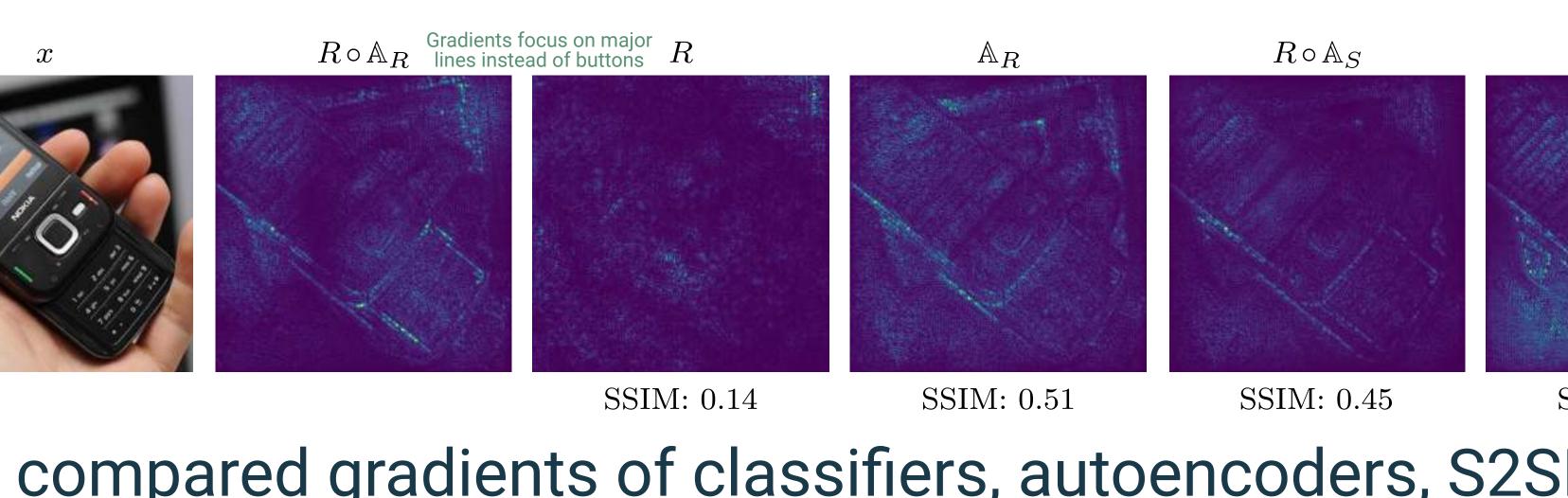


• Carlini-Wagner CW; iterative optimization-based method

With our defense gradients are intercepted at the input layer of the S2SNet autoencoder. These attacks are far less effective. They target the reconstruction, not the classification. $x \xrightarrow{\sum_{x \in \mathbb{Z}} \mathbb{N} \text{et autoencoder}} A(x) \xrightarrow{\sum_{x \in \mathbb{Z}} \mathbb{N} \text{et autoencoder}} f(\mathbb{A}(x)) \xrightarrow{\approx f(x)} \delta_x$

Results

1. Structural nature of gradients

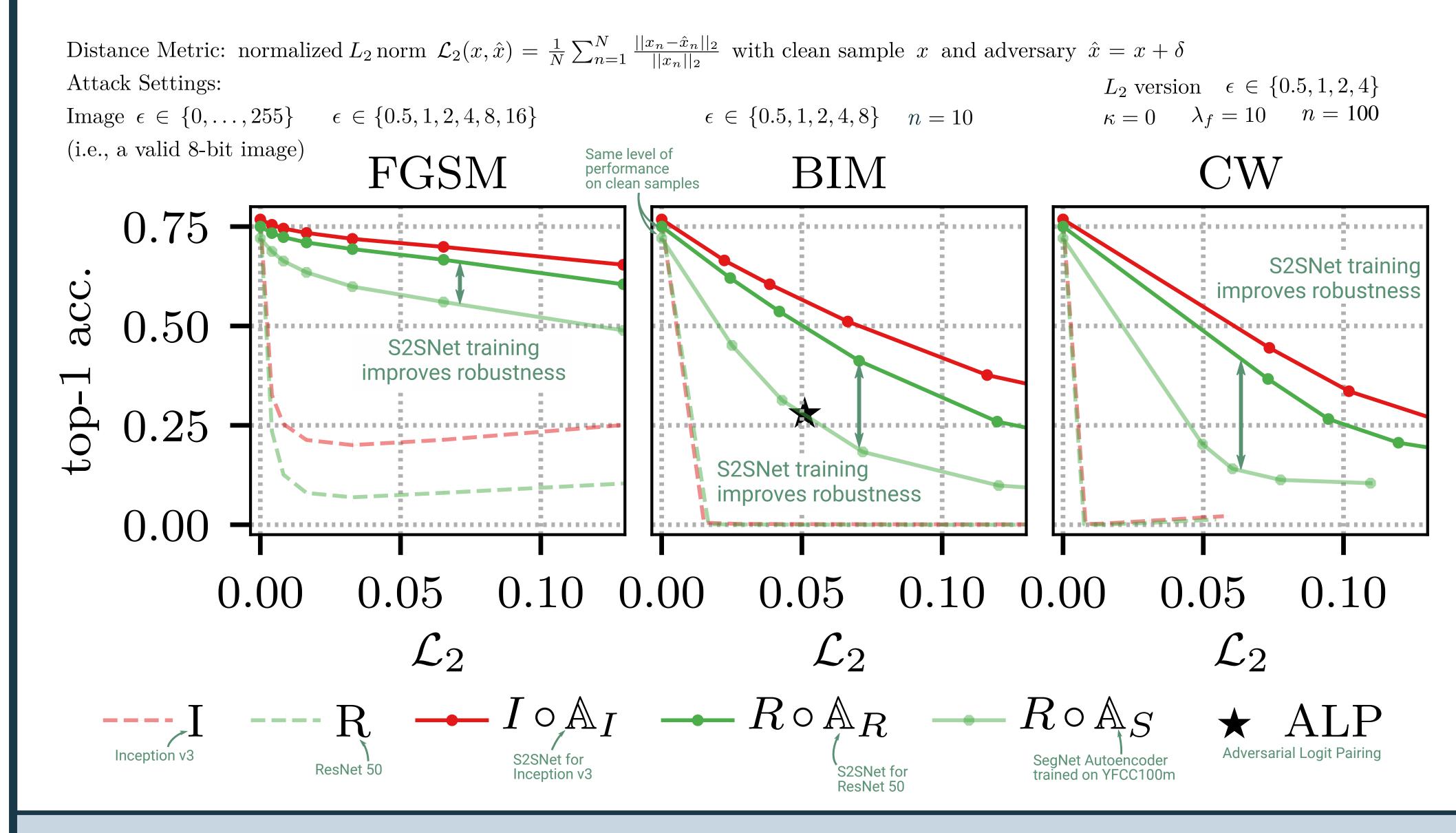


We compared gradients of classifiers, autoencoders, S2SNets, and combined networks using structural similarity (SSIM).

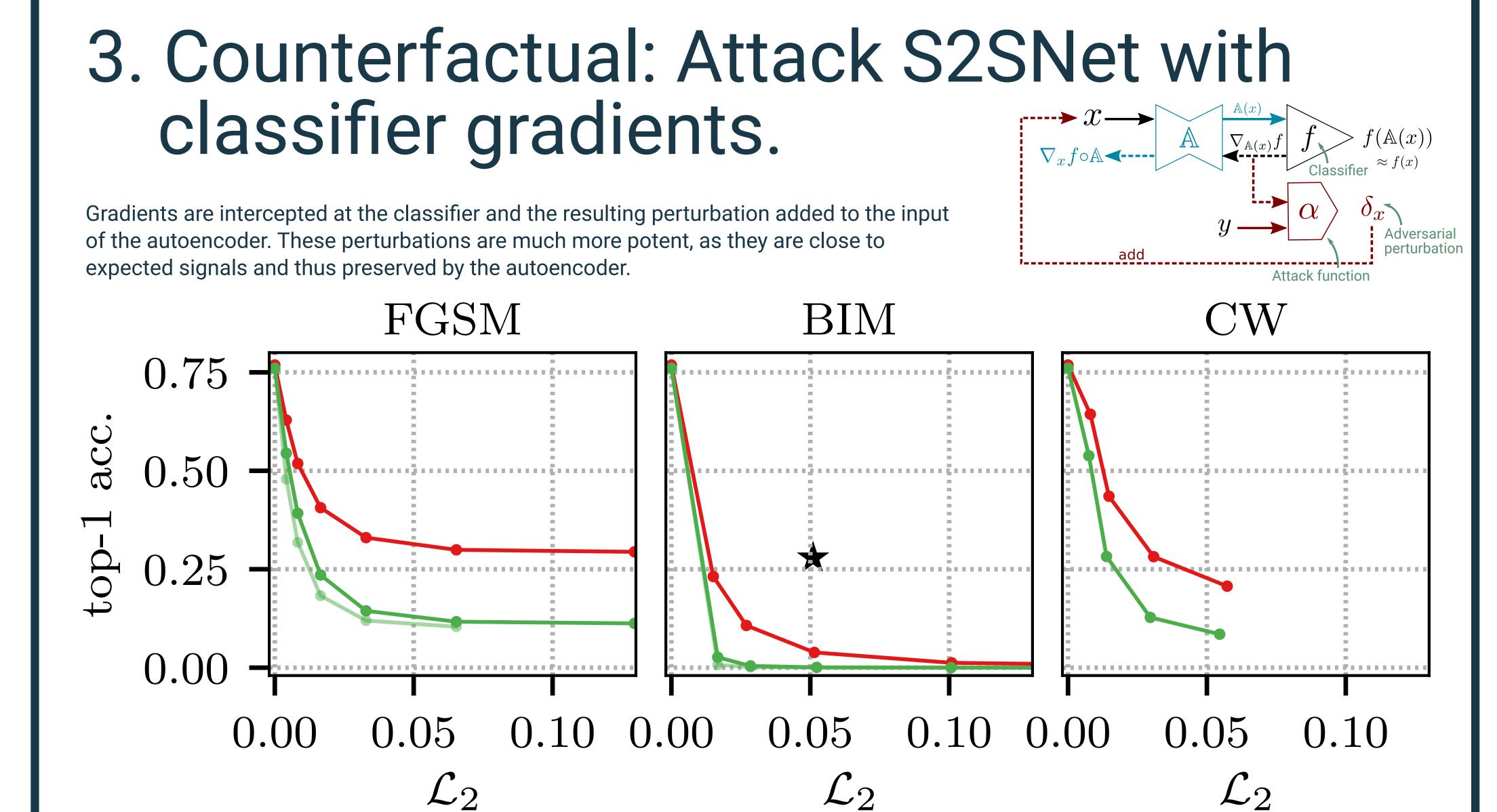
Gradients of autoencoders are more similar among themselves than compared to any classifier on its own.

2. S2SNets defend against well-known adversarial attacks

We performed a large-scale evaluation of our defense method on the full ILSVRC 2012 validation set.



The autoencoder on its own adds an **architectural bottleneck**. Fine-tuning into S2SNet adds a **representational bottleneck**, thus greatly improving robustness to tested attacks.



4. Counterfactual: Attack classifier with S2SNet gradients.

Gradients are intercepted at the autoencoder and the resulting perturbation added to the input of the classifier. These perturbations are not understood by the classifier and thus much less potent. This is further evidence that previously the autoencoder and thus the structure was attacked.

FGSM

BIM

CW 0.75 0.50 0.25

Future Work:

- Investigate whether interpretability techniques provide further insights on how adversarial attacks interact with S2SNets.
- Develop specific (lighter weight; production-ready)
 S2SNet autoencoder architecture.
- What about non-gradient-based attacks?





