# Adversarial Defense based on Structure-to-Signal Autoencoders

Joachim Folz*, Sebastian Palacio*, Joern Hees, Damian Borth, and Andreas Dengel

[1] German Research Center for Artificial Intelligence (DFKI)
[2] TU Kaiserslautern
`first.last@dfki.de`

**Abstract.** Adversarial attack methods have demonstrated the fragility of deep neural networks. Their imperceptible perturbations are frequently able fool classifiers into potentially dangerous misclassifications. We propose a novel way to interpret adversarial perturbations in terms of the effective input signal that classifiers actually use. Based on this, we apply specially trained autoencoders, referred to as S2SNets, as defense mechanism. They follow a two-stage training scheme: first unsupervised, followed by a fine-tuning of the decoder, using gradients from an existing classifier. S2SNets induce a shift in the distribution of gradients propagated through them, stripping them from class-dependent signal. We analyze their robustness against several white-box and gray-box scenarios on the large ImageNet dataset. Our approach reaches comparable resilience in white-box attack scenarios as other state-of-the-art defenses in gray-box scenarios. We further analyze the relationships of AlexNet, VGG 16, ResNet 50 and Inception v3 in adversarial space, and found that VGG 16 is the easiest to fool, while perturbations from ResNet 50 are the most transferable.

**Keywords:** adversarial defense, CNNs, deep learning, autoencoders

## 1 Introduction

We nowadays see an increasing adoption of deep learning techniques in production systems, partially even safety-relevant ones (as exemplified in [1]). Hence, it is not surprising that the discovery of adversarial spaces for neural networks [2] has sparked a lot of interest. A growing community has quickly focused on different ways to reach those spaces [3,4,5], understand their properties [1,6,7], and protect vulnerable models from their malicious nature [8,9,10].

Most prevalent methods exploiting adversarial spaces [3,4,5,10] use gradients as their main starting point to search for perturbations surrounding a clean sample. Due to the intractability of transformations modeled by neural networks and the limited amount of change that is allowed for perturbations to be considered adversarial, gradients from a classifier expose just enough information about

---

* Authors contributed equally

which parts of the input correlate highly with the label that a model associates to it. With this in mind, defenses against adversarial attacks have been devised upon those same gradients in two fundamental ways: 1) They complement traditional optimization schemes with a two-fold objective that minimizes the overall prediction cost while maximizing the perturbation space around clean images that classifiers can withstand [10,11,12]. 2) Gradients are blocked or obfuscated in such a way that attacking algorithms can no longer use them to find effective adversarial perturbations [13,14,15]. Type 1 methods enjoy mathematical rigor and hence, provide formal guarantees with respect to the kind of perturbations they are robust to. However, note that this can also be disadvantageous since networks become attack-dependent. Any other strategy for finding perturbations could circumvent such a defense mechanism [16]. While effective for small-scale problems such as MNIST and CIFAR, we found no empirical evidence that these methods scale to larger problems such as ImageNet. It has even been shown that defenses for adversarial attacks tested on small datasets do not scale well when applied to bigger problems [6].

Currently, large scale state-of-the-art defenses rely on the second use of gradients: suppression [17] and blockage [14,18]. As defined by Athalye *et al.* [16], gradient obfuscation is the result of instabilities from vanishing or exploding gradients, and the use of stochastic or non-differentiable preprocessing steps. All these alternatives can be modeled as lossy identities where the original signal contained in the input is preserved, while the adversarial perturbation is destroyed. The usefulness of this principle fits well with findings from a recent study showing that image classifiers only use a small fraction of the entire signal within original input. Therefore, a portion of its information can indeed be dropped without affecting performance [19].

In this paper, we propose an alternative defense that affects the information contained in gradients by reforming its class-related signal into a structural one. Intuitively, we learn an identity function that encodes structure and decodes only the structural parts of the input necessary for classification, dropping everything else. To this end, an autoencoder (AE) is trained to approximate the identity function that preserves only the part of the signal that is *useful* for a target classifier. The structural information is preserved by training both encoder and decoder unsupervised, and fine-tuning only the decoder with gradients coming from an existing classifier. By using a function that only looks at structure, gradients are devoid of any class-related information, therefore invalidating the fundamental assumptions about gradients that attackers rely on. We call this defense a *Structure-To-Signal Network* (S2SNet).

**Formal Definitions**  Let $f : \mathbb{R}^3 \to \{1, \dots, k\}$ be an image classifier, and $\tilde{x}_f$ the portion of the signal in the original input $x$ that is effectively being used by $f$. In other words, $f(x) = f(\tilde{x}_f)$ subject to $I(x, \tilde{x}_f) < I(x, x)$, where $I$ is a measure of information *e.g.*, the *normalized mutual information* [20]. Let $\mathcal{P} = \{\delta : ||\delta||_L < \epsilon\} \subseteq \mathbb{R}^3$ be the space of all adversarial and non-adversarial perturbations under a given norm $L \in \{1, 2, \infty\}$. We also define the sub-space
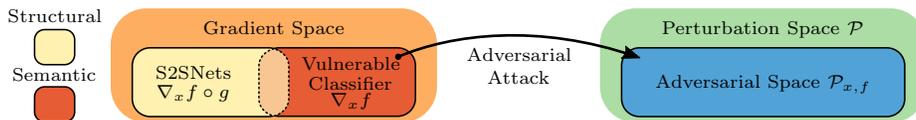
Fig. 1: Overview of gradient and perturbation spaces, and their relation to S2SNets and adversarial attacks.

$\mathcal{P}_{x,f} \subseteq \mathcal{P}$ as the set of adversarial perturbations that are reachable from the input gradients $\nabla_x C(f(x), y_{1...k}) \subseteq \mathbb{R}^3$, where $C$ represents the cost function, and $y_i$ is the ground truth used for training. Note that said spaces depend on $f$ and $x$, but not on the class $y$ or the cost function $C$. Hence, unless noted otherwise, we use the simplified notation $\nabla_x f = \nabla_x C(f(x), y_{1...k})$ from now on. Adversarial perturbations are elements $\delta_{adv} \in \mathcal{P}_{x,f}$ such that $f(x) \neq f(x + \delta_{adv})$ [3]. As attacks use gradients to compute adversarial perturbations, and such gradients can only correspond to parts of the signal that are used by the model, it follows that $\mathcal{P}_{x,f} \equiv \mathcal{P}_{\tilde{x},f}$. An S2SNet can hence be defined as a function $g : \mathbb{R}^3 \to \mathbb{R}^3$ such that $g(x) = \tilde{x}_f$. This way, $g$ defends $f$ by serving as a proxy for incoming, potentially malicious inputs, as well as exposing gradients via the function composition $f \circ g$. The relation between $x$ and $\tilde{x}$ concerning information implies that $g$ produces lossy reconstructions of the input space. We train $g$ in a way that gradients $\nabla_x f \circ g$ lie in a different space than those leading attackers to $\mathcal{P}_{\tilde{x},f}$ (Figure 1). In other words, we train $g$ such that $|\nabla_x f \circ g \cap \nabla_x f|$ is minimized. This in turn, will cause the intersection $|\mathcal{P}_{x,f \circ g} \cap \mathcal{P}_{\tilde{x},f}|$ to be smaller, resulting in perturbations that are non-adversarial (i.e., $\delta \in \mathcal{P} - \mathcal{P}_{\hat{x},f}$). Further details of the architecture and training of $f$ and $g$ are discussed in Section 3.1.

The architecture of this defense offers important advantages when dealing with adversarial attacks:

– **Zero compromise for safe use-cases**: there is no drop in performance when the defense is deployed, but the network is not under attack. As the transformation done by S2SNet is preserving all required signal $\tilde{x}_f$, using either $g(x)$ or $x$ has no impact to the classifier when clean images are used.
– **Removing S2SNet is a defense strategy**: a novel gray-box defense (i.e., when the attacker knows the classifier but not the defense mechanism) works by giving away gradient information from an S2SNet to the attacker and then removing it, applying adversarial attacks (based on gradients from S2SNet) to the original classifier instead.
– **Attack agnostic**: S2SNets rely on the same information used by adversarial attacks but not on the attacks themselves. Therefore, S2SNets do not require any assumptions with respect to the specific way an attack works.

---

[3] Strictly speaking, adversarial perturbations can be reached through other domains that do not depend on gradients (as shown by Ngueyen et al. [21]) but so far, all instances of strong adversarial attack methods, base their entire strategy on the information of gradients.

– ***Post-hoc* implementation**: our defense uses gradients from a trained network, and can be used to defend models that are already in production. Likewise, no special considerations need to be made when training a new classifier from scratch.
– **Compatibility with other defenses**: due to the compositional nature of this approach, any additional defense strategies that work with the original classifier can be implemented for the ensemble $f \circ g$.

We test S2SNets on two high-performing image classifiers (ResNet50 [22] and Inception-v3 [23]) against three attack methods (Fast Gradient-Sign Method (FGSM) [3], Basic Iterative Method (BIM) [24] and Carlini-Wager (CW) [5]) on the large scale ImageNet [25] dataset. Experiments are conducted on both classifiers under white-box and gray-box conditions. An evaluation of the effectiveness of S2SNets with respect to regular AEs (*e.g.*, as proposed in [26]) is presented in Section 3.1, empirically proving that S2SNets are a better approximation of the signal $\tilde{x}_f$ that is used by a classifier. Furthermore, an evaluation of the gradient space $\mathcal{P}_{x,f}$ and $\mathcal{P}_{x,f \circ g}$ is conducted, showing that their intersection does indeed approach the empty set.

The main contributions of this paper are threefold: First, we propose a novel way to interpret adversarial perturbations, namely in terms of the effective input signal that classifiers use $\tilde{x}_f$. Second, we introduce a robust and flexible defense against large-scale adversarial attacks based on S2SNets. And third, we provide a comprehensive baseline evaluation of adversarial attacks for several state-of-the-art models on a large dataset.

## 2    Related Work

The fast growing interest in the phenomenon of adversarial attacks has gained momentum since its discovery [2] and has had three main areas of focus. The first area is the one that seeks new and more effective ways of reaching adversarial spaces. In [3], a first comprehensive analysis of the extent of adversarial spaces was explored, proposing a fast method to compute perturbations based on the sign of gradients. An iterative version of this method was later introduced [24] and shown to work significantly better, even when applied to images that were physically printed and digitized again. A prominent exception to attacks based on gradients succeeded using evolutionary algorithms [21]. Nevertheless, this has not been a practical wide spread method, mostly due to how costly it is to compute. Papernot *et al.* [1] showed how effective adversarial attacks could get, even with very few assumptions about the attacked model. Finally, there are methods that go beyond a greedy iteration over the gradient space and perform different kinds of optimization that maximize misclassification while minimizing the norm of the perturbation [4,5].

The second area focuses on understanding the properties of adversarial perturbations. The work of Goodfellow *et al.* [3] was already pointing at the linear nature of neural networks as the main enabler of adversarial attacks. This went in opposition of what was initially theorized, where the claim was that

non-linearities were the main vulnerability. Not only was it possible to perturb natural looking images to look like something completely different, but it was also possible to get models issuing predictions with high confidence using either noise images, or highly artificial patterns [2,21]. The transferability of adversarial perturbations was shown to be possible by crafting attacks on one network and using them to fool a second classifier [27]. However, transferable attacks are limited to the simpler methods as iterative ones tend to exploit particularities of each model, and hence lose power when used in different architectures [6]. As it turns out, not only is adversarial noise transferable between models but it is also possible to transfer a single universal adversarial perturbation to all samples in a dataset to achieve high misclassification rates [7]. Said individual perturbations can even be applied to physical objects and bias a model towards a specific class [28].

The third and arguably the most popular area of research has focused on how networks can be protected against such attacks. Strategies include changing the optimization objective to account for possible adversarial spaces [10,11], detection [9], dataset augmentation that includes adversarial examples [3,12], suppressing perturbations [8,15,17,26,29] or obfuscating the gradients to prevent attackers from estimating an effective perturbation [13,14,18,30,31].

In this work, we build on the idea of using AEs as a compressed representation of the input [15,26], but tailored towards a specific characteristic of adversarial perturbations [17], using the notion of useful input signal (*i.e.*, the one effectively used by a classifier) [19]. Furthermore, we explore the nature of adversarial perturbations and its relationship with network capacity, in terms of used signal.

## 3 Methods

This section explains in detail the architecture of an S2SNet and its particular signal-preserving training scheme, followed by an empirical evaluation of the gradients it provides. We start by testing the robustness of S2SNets in a white-box setting, and compare it to a simple baseline using regular AEs. To recreate realistic attack conditions, we further test the ensemble network, simulating a re-parametrization technique similar to [16], aimed at circumventing the defense, and explore further strategies to cope with this attack. Next, we examine the performance of S2SNets in a gray-box scenario. Finally, we provide an evaluation of the transferability of single step attacks for different models, and correlate their overlap in terms of input signal from the perspective of adversarial perturbations.

### 3.1 Structure-to-Signal Networks

S2SNets start out as plain AEs that are trained on the large-scale YFCC100m data set [32]. Only a single pass is required, as the $\approx 100$ million images are more

than sufficient to train the underlying SegNet architecture [33] to convergence[4]. This network, referred to as $\mathbb{A}_S$, is able to reproduce input signals required by a diverse set of classifiers such that their top-1 accuracy is within $\approx 3$ percentage points of the original classification performance [19].

To model the effective input signal used by a trained classifier $f$, Palacio *et al.* propose to further fine-tune the decoder of $\mathbb{A}_S$ using gradients from $f$ itself. This allows $\mathbb{A}_S$ to learn a way to decode the input that retains the signal required by $f$. This fine-tuned variant, called $\mathbb{A}_f$, reconstructs images such that the original top-1 accuracy of $f$ is preserved, while amount of information in the reconstructed image (measured as normalized mutual information) decreases with respect to the original sample.

Note that, since the encoder of $\mathbb{A}_f$ was trained unsupervised, any intermediate representations produced by $\mathbb{A}_f$ are entirely class-agnostic. This means that, during backpropagation through an S2SNet, gradients that can be read at the shallowest layer correspond only to information about structure. Intuitively, gradients from $\mathbb{A}_f$ point to parts of the image that can be changed to influence the *reconstruction* error.

In the following, we measure the extent to which gradients shift when images are forwarded through these networks. Furthermore we explore and quantify their emerging resilience to adversarial attacks in Section 4.


### 3.2  Properties of Gradient Distributions

To verify that the Structure-to-Signal training scheme produces large shifts in the distribution of gradients, we forward images through a classifier, a pre-trained SegNet, and through the fine-tuned counterpart to compare their gradients. We use the magnitude of the gradients instead of raw values to stress the differences of their spatial distribution. A large change in the position where gradients originally occur within the image is a good indicator that the information conveyed by gradients has changed.

Concretely, we quantify the structural similarity [34] (SSIM; a locally normalized mean square error measured in a sliding window) of gradients obtained by the same image when passed through a ResNet50 ($R$), an S2SNet fine-tuned to defend the same ResNet50 ($R \circ \mathbb{A}_R$), a plain SegNet AE coupled with ResNet50 ($R \circ \mathbb{A}_S$), the fine-tuned S2SNet without the classifier ($\mathbb{A}_R$), and the plain SegNet AE alone ($\mathbb{A}_S$). While the first three models require gradients to be computed with respect to a class label, the last two are produced by measuring the reconstruction error. Also, note that for $\mathbb{A}_R$ the true reconstruction cannot be directly obtained as it is indirectly defined by the classifier it was fine-tuned on. In this case, reconstruction gradients are computed by comparing its output to the original input.

---

[4] We choose a large architecture as an upper bound to the ideal identity function $g(x) = \tilde{x}_f$, because it was proven capable of encoding the semantics of the deep image classifiers tested in this work.

Table 1 reports the mean SSIM of the gradients over ImageNet's validation set for all combinations of network pairs. Note that the dissimilarity between $\mathbb{A}_S$ and $\mathbb{A}_R$ indicates that the Structure-to-Signal training scheme has indeed changed the reconstruction process *i.e.*, the identity function being computed based on the input. Similarly, comparing the SSIM values of $R \circ \mathbb{A}_S$ and $R \circ \mathbb{A}_R$ reveals a difference in the information contained by their gradients. Most importantly, the similarity between the gradients coming from classifier $R$ and either of the ensembles $R \circ \mathbb{A}_S$ and $R \circ \mathbb{A}_R$ is considerably smaller (0.18 and 0.17) than any combination involving AEs exclusively ($> 0.32$). This disparity indicates how much the position of gradients change (and hence the information contained within), when passed through an S2SNet.

Further evidence for the class-agnostic nature of gradients propagated through AEs can be found when comparing the SSIM of gradient magnitudes when different target labels are used to compute gradients. Let $x$ be some input image and $y^*$ its true label. We then randomly select a different label $\hat{y} \neq y^*$ and compute

$$SSIM(||\nabla_x C(f(x), y^*)||, ||\nabla_x C(f(x), \hat{y})||)$$

for all $x$ in ImageNet's validation set, and $f \in \{R, R \circ \mathbb{A}_R, R \circ \mathbb{A}_S\}$. We measured mean SSIM values of 0.50 for $R \circ \mathbb{A}_R$, 0.47 for $R \circ \mathbb{A}_S$, and 0.34 for $R$. This means that the influence of the label is smaller when gradients are propagated through the AE, but also emphasizes how dissimilar gradients of just ResNet 50 are, compared to any AE at just 0.12 to 0.18 SSIM (Table 1).

Figure 2 visualizes this phenomenon. Here, gradient magnitudes observed for just AEs ($\mathbb{A}_S$, $\mathbb{A}_R$) predominantly highlight edges as source of error. This is expected, since it is more difficult to accurately reproduce the high frequencies required by sharp edges, compared to the lower frequencies of blobs. Extracting magnitudes based on classification from AEs show similar structures ($R \circ \mathbb{A}_R$, $R \circ \mathbb{A}_S$). Some coincidental overlap between ResNet 50 ($R$) and AE variants is unavoidable, since edges are also important for classification [35]. However, the classifier on its own differs considerably from all other patterns. Overall, SSIM is at least twice as high between AE variants, than between the classifier and any of the AE configurations.

Table 1: Pairwise mean SSIM of input gradient magnitudes for ResNet 50 ($R$) on the ImageNet validation set, with and without being passed through $\mathbb{A}_R$ or $\mathbb{A}_S$. SSIM values of $R$ w.r.t. any AE variant show the least similarity.

|  | $R$ | $R \circ \mathbb{A}_R$ | $R \circ \mathbb{A}_S$ | $\mathbb{A}_R$ | $\mathbb{A}_S$ |
|---|---|---|---|---|---|
| $R$ | 1.00 | 0.17 | 0.18 | 0.12 | 0.14 |
| $R \circ \mathbb{A}_R$ | 0.17 | 1.00 | 0.40 | 0.46 | 0.32 |
| $R \circ \mathbb{A}_S$ | 0.18 | 0.40 | 1.00 | 0.37 | 0.36 |
| $\mathbb{A}_R$ | 0.12 | 0.46 | 0.37 | 1.00 | 0.36 |
| $\mathbb{A}_S$ | 0.14 | 0.32 | 0.36 | 0.36 | 1.00 |

## 4   Experiments

This section presents experiments quantifying the robustness of S2SNets when used as a defense against adversarial attacks (Figure 3). The experimental setup closely follows the conditions from Guo *et al.* [14] in order to facilitate comparability:

- **Dataset**: We use ImageNet to test S2SNets in a challenging, large-scale scenario. Classifiers are trained on its training set and evaluations are carried out on the full validation set $X$ (50000 images).
- **Image classifier under attack**: we use ResNet 50 ($R$) and Inception v3 ($I$), both pre-trained on ImageNet as target models. The classifiers have been trained under clean conditions *i.e.*, no special considerations with respect to adversarial attacks were made during training.
- **Defense**: we train an S2SNet for each of the classifiers under attack, following the scheme described in Section 3.1. These defenses are denoted as $\mathbb{A}_R$ and $\mathbb{A}_I$ for $R$ and $I$ respectively.
- **Perturbation Magnitude**: we use the normalized $L_2$ norm ($\mathcal{L}_2$) between a clean sample $x$ and its adversary $\hat{x} = x + \epsilon \cdot \delta$, as defined by Guo *et al.* [14]. Epsilon values for each attack are listed below.
- **Defense Strength Metric**: vulnerability to adversarial attacks is measured in terms of the number of newly misclassified samples. More precisely, for any given attack to classifier $f$, we calculate $\frac{\sum_{x \in TP} \mathcal{L}_2(x,\hat{x})}{|TP|}$, where $TP = \{x \in X | f(x) = y^*\}$ is the set of true positives and $\hat{x}$ is the adversarial example generated by the attack, based on $x$.
- **Attack Methods**: protected models are tested against a single step method, an iterative variant, and an optimization-based alternative. Note that to replicate realistic threat conditions, all resulting adversarial samples are cast to the discrete RGB range [0, 255].
  - *Fast Gradient Sign Method (FGSM)* [3]: a simple, yet effective attack that works also when transferred to different models. Epsilon values used for this method are $\epsilon \in \{0.5, 1, 2, 4, 8, 16\}$.
  - *Basic Iterative Method (BIM)* [24]: an iterative version of FGSM that shows more attack effectiveness but less transferable properties. Epsilon
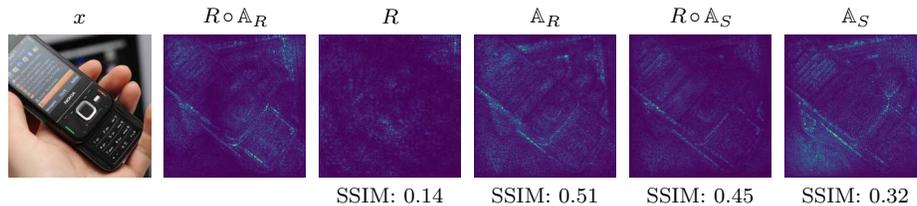


Fig. 2: Gradient magnitudes for ResNet 50 ($R$) given a single input $x$, propagated through $\mathbb{A}_R$ or $\mathbb{A}_S$. SSIM values are in comparison to $||\nabla_x R \circ \mathbb{A}_R(x)||$.

Fig. 3: Overview of analyzed attack/defense scenarios.

values used for BIM are $\epsilon \in \{0.5, 1, 2, 4, 8\}$. The number of iterations is fixed at 10.

- *Carlini-Wagner $L_2$ (CW)* [5]: an optimization-based method that has proven to be effective even against hardened models. Note that this attack issues perturbations that lay in a continuous domain. Epsilon values used for CW are $\epsilon \in \{0.5, 1, 2, 4\}$; the number of iterations is fixed at 100, $\kappa = 0$ and $\lambda_f = 10$.

– **Attack Conditions**: we test the proposed defenses under two conditions.
  - *White-box setting*: the attacker has knowledge about the classifier and the defense. This includes reading access to both the predictions of the classifier, intermediate activations and backpropagated gradients. The attacker is forced to forward valid images through the defended network *i.e.*, through the composed classifier $R \circ \mathbb{A}_R$ or $I \circ \mathbb{A}_I$.
  - *Gray-box setting*: the attacker has access to a classification model but is unaware of its defense strategy.

### 4.1  White-Box

For this setting, input images flow first through the S2SNet before reaching the original classifier. Similarly, gradients are read from the shallowest layer of the S2SNet. For comparison, attacks are also run on the unprotected versions of ResNet 50 and Inception v3. Results are summarized in Figure 4.

Under these conditions, plain models are only able to moderately resist FGSM attacks, and fail completely for the more capable BIM and CW attacks.
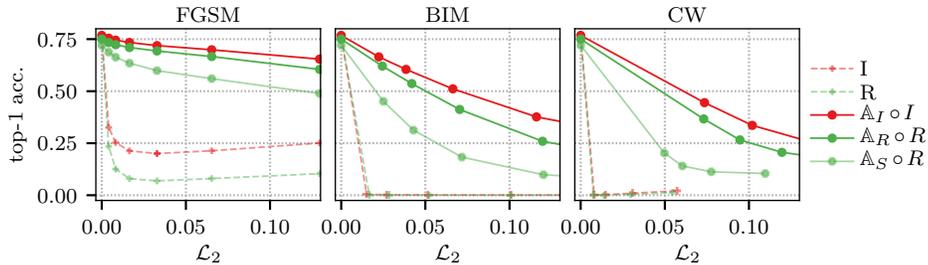


Fig. 4: White-box attacks on ResNet 50 ($R$) and Inception v3 ($I$), with (solid) and without (dashed) S2SNet as a defense.

In contrast, S2SNets provide high levels of protection. CW, being an optimization attack and generally the most capable in our comparison, is expected to be more effective at breaking through the S2SNet defense. However, it does so by introducing large perturbations even for small values of $\epsilon$. In fact, none of the attacking configurations was able to entirely fool any of our defended classifiers for even the highest levels of $\mathcal{L}_2$ in our tests. Note that these white-box results are, to our surprise, already comparable with some state-of-the-art gray-box defenses [14,18] (*i.e.*, known classifier but unknown defense). Despite the less favorable conditions for a white-box defense, S2SNets already match alternative state-of-the-art protections that were tested under the more permissive assumptions allowed in gray-box settings.

For completeness, we also evaluate the defense using a pre-trained SegNet($\mathbb{A}_S$) instead of an S2SNet for ResNet 50 (shown in Figure 4 with light solid green). This is most similar to the defense proposed by Meng *et al.* [26]. We can confirm that having just an AE does not suffice to guard the classifier against adversarial attacks and its performance is consistently lower than S2SNets. For a visual analysis of the attacks and their reconstructions by S2SNets, we refer the reader to the supplementary material.

**Bypassing S2SNets through Reparametrization** In order to push the limits of S2SNets, we now simulate a more hostile scenario where an attacker tries to actively circumvent the defense mechanism. Based on the work of Athalye *et al.* [16], a reparametrization of the input space can be implemented for defenses that operate under the composition of functions $f \circ g$. This works by defining the original input $x$ as a function of a hidden state $x = h(z)$ such that $g(h(z)) = h(z)$. Note that the main motivation behind reparametrization is to alleviate instabilities of gradients caused by defenses that *rely* on said instability. Although S2SNets are very deep architectures, the resilience of this defense lays in the directed change that has been induced within the information contained in gradients. In addition, there is no trivial way to come up with a suitable $h(z)$ that does not end up inducing the same transformation in the gradients, as done by S2SNets.

Despite all these concerns, we assume for this experiment that such $h(z)$ can be found, and that attacking $h(z)$ does indeed circumvent S2SNets altogether. We simulate the potential strength of this attack by using gradients of the original (unprotected) classifier, and applying them to the input directly. The resulting adversarial attack is then passed through the hardened classifier. We refer to this hypothetical scenario as White-Box$^+$. Results are shown in Figure 5.

Under these conditions, we observe that hardened models revert back to the behavior shown by their corresponding unprotected versions for FGSM and BIM. In general, this condition is expected, and confirms once more that S2SNets are being trained to preserve the information that is useful to the classifier. Gradients collected directly from a vulnerable model, have by definition, only information that is useful for classification and hence, perturbations based on those gradients will be preserved by S2SNets. Interestingly, while CW is most successful in the

Fig. 5: White-Box$^+$ attacks on ResNet 50 ($R$) and Inception v3 ($I$) classifiers.

previous white-box experiment, its highly optimized perturbations are less effective here. We believe that CW is "overfitting" more strongly on the adversarial signal of the original model than what S2SNets find useful to preserve.

At this point, we can further exploit the benefits that S2SNets enjoy as a defense, by adding another layer of protection and showing how it affects the effectiveness of an attack. As shown by Palacio *et al.* [19], AEs following a structure-to-signal training scheme exhibit strong resilience to random noise, as opposed to traditionally trained AEs. We demonstrate that it is straightforward to add a layer of protection based on random noise. Said stochastic strategy is added after the adversarial image has been computed but before it passes through an S2SNet defense. We experiment with three sources of noise:

- *Gaussian Noise*: $x \sim \frac{1}{\sqrt{2\pi\sigma^2}}e^{-(x^2/(2\sigma^2))}$, where $\sigma = \eta$.
- *Uniform Noise*: $\eta x$, where $x \sim \mathcal{U}(-1, 1)$.
- *Sign Noise*: $\eta \mathtt{sgn}(x)$, where $x \sim \mathcal{U}(-0.5, 0.5)$

Figure 6 shows the results of a BIM attack for ResNet 50, with noise levels $\eta \in \{0.5, 4, 8, 12, 15, 20, 30\}$ under White-Box$^+$ conditions. Overall, as noise strength increases, the resilience to adversarial attacks improves. The initial degradation under zero adversarial attacks is dependent on the amount of noise $\eta$ which can be tuned as the trade-off between maximum accuracy and adversarial robustness.



Fig. 6: White-Box$^+$ attacks on $R \circ \mathbb{A}_R$ with BIM, mitigated with additive sign, random and Gaussian noise of varying strength $\eta$.

Fig. 7: Gray-Box$^-$ attacks on ResNet 50 ($R$) and Inception v3 ($I$).

We observe that both Gaussian and sign noise perform almost identically, while uniform noise offers the least effective protection.

### 4.2   Gray-Box

In contrast to white-box scenarios, gray-box attacks assume that there is limited access to information about the attack target. More specifically, the conditions for gray-box define that the attacker has knowledge about the network but not about the defense strategy.

Given the compositional nature of S2SNets, a novel way to defend a classification network under gray-box conditions consists in giving access to the gradients of S2SNets. Once an attack is ready, the S2SNet is *removed* and the perturbed image is processed by the classifier only. In other words, for a classifier $f$ and corresponding defense $g$, the attack is crafted based on the combined network $f \circ g$ but forwarded to just the classifier $f$. We refer to this type of defense as Gray-Box$^-$.

We run the attacks with the same experimental settings described in Section 4.1, but enforce a defense policy following Gray-Box$^-$ conditions. Results are presented in Figure 7. Overall, we see that this scenario is consistently robust to any adversarial attack. Fooling perturbations are clearly visible and cannot be considered adversarial anymore. With BIM-based attacks, classifiers gain back roughly half as much accuracy as in the white-box case. Again, CW is more comparable to FGSM than BIM, i.e., mostly ineffective, and consistent with observations made in the White-Box$^+$ case.

By combining the results of Gray-Box$^-$ and White-Box$^+$, we can reason about the relationship of $\mathcal{P}_{x,f}$ and $\mathcal{P}_{x,f\circ g}$. First, the White-Box$^+$ experiment tells us that attacks crafted for $f$ are valid attacks for $f \circ g$. Following from the definition of these spaces, it holds that $\mathcal{P}_{x,f} \subsetneq \mathcal{P}_{x,f\circ g}$. Furthermore, looking at the Gray-Box$^-$ experiments, taking elements from $\nabla_x f \circ g$ to create an attack, produces elements in $\mathcal{P}_{x,f\circ g}$. As these do not attack $f$, we conclude that those same elements do not lie in $\mathcal{P}_{x,f}$. We finally conclude that with the tested attacks $\mathcal{P}_{x,f}$ cannot be reachable from $\nabla_x f \circ g$.

Fig. 8: Relative accuracy for black-box attacks on AlexNet ($A$), VGG 16 ($V$), ResNet 50 ($R$), and Inception v3 ($I$).

### 4.3   Classifier Relationships in Adversarial Space

Given the high robustness of S2SNets in the Gray-Box$^-$ scenario, we do not expect further insights from traditional black-box attacks. Instead, we explore the relationships between the signal that different image classifiers use, as reported in [19]. Intuitively, if a signal used by a classifier $f_1$ also encompasses the signal that another classifier $f_2$ uses, then an adversarial attack crafted for $f_1$ should also fool $f_2$. Note that the relation is directional, and it may not hold in the opposite way.

To test this, we construct adversarial samples using FGSM for four reference classifiers: AlexNet [36], VGG 16 [37], ResNet 50, and Inception v3. For each classifier we run black-box attacks with perturbations computed on the other three models. The resulting accuracies (in terms of their relative drop in accuracy) are shown in Figure 8.

As expected, adversarial examples from higher compatible architectures result in higher fooling rations (lower accuracy). Overall, adversarial samples created for AlexNet are the least compatible among the four, mainly due to its comparably lower accuracy and limited amount of input signal used. On the other hand, ResNet 50 produces the most compatible perturbations, being the model generating attacks that were most effective when tried on other architectures.

Furthermore, we can also analyze how easy to fool a network is, based on the following criteria:

– Fastest drop in accuracy per classifier.
– Mean accuracy of adversarial samples coming from other classifiers that have the same $\mathcal{L}_2$.

With either criterion, results clearly show that VGG 16 is the easiest to fool, coming up on first place, followed by ResNet 50, AlexNet, and finally Inception v3. This relationship in terms of useful signal, aligns with results in [19]. Furthermore, similar experiments in [7] based on universal perturbations, indicate that a selection of similar architectures (CaffeNet, VGG 19, ResNet 101 and GoogLeNet) maintain the same relationship.

## 5   Conclusions & Future Work

We have proposed S2SNets as a new method to defend neural networks against adversarial examples. We model the defense strategy as a transformation of the domain used by attackers, namely gradients coming from the attacked classifier. Instead of focusing on gradient obfuscation via non-differentiable methods or any other instabilities, we purposely induce a transformation on the gradients that strip them from any semantic information. S2SNets work by masking the classifier via the function composition. That way, the information of inputs is preserved for classification, but gradients point to structural changes for reconstructing the original sample. Such a defense is possible by using a novel two-stage *Structure-to-Signal* training scheme for deep AEs. On the first stage, the AE is trained unsupervised in the traditional way. For the second part, only the decoder gets fine-tuned with gradients from the model that is to be defended.

We evaluate the proposed defense under white-box and gray-box settings using a large scale dataset, against three different attack methods, on two highly performing deep image classifiers. A baseline comparison shows that the two-staged training scheme performs better than using regular AEs. Most interestingly, we show that resiliency to adversarial noise under white-box conditions, exhibit comparable performance to state-of-the-art under more favorable gray-box settings. Furthermore, we show how the properties of S2SNets can be exploited to add more defense mechanisms to maintain robustness even under the harshest, albeit currently hypothetical conditions, where the protection of S2SNets is circumvented. A gray-box scenario was also tested where the defense consists on the removal of S2SNets, showing high levels of robustness for all attacks. Finally, a comparison between the resiliency of four well-known deep CNNs is presented, providing further evidence that a relation of order exists between these classifiers, in terms of the amount of signal they use; this time, in terms of the effectiveness of adversarial noise.

S2SNets are only one way in which the transformation can occur in gradient space. We would like to explore other ways in which such transformation can occur, and even if the intersection between gradients yielding successful adversarial perturbations can be effectively zero. The signal-preserving nature of S2SNets make this defense a potential mechanism to explore and understand the nature of attacks. Comparing classification consistency of a clean sample, before and after being passed through a S2SNet, has potential implications for detection of adversarial attacks by learning abnormal distribution fluctuations.

## Acknowledgments

# References

1. Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z.B., Swami, A.: Practical black-box attacks against deep learning systems using adversarial examples. arXiv preprint (2016)
2. Szegedy, C., Zaremba, W., Sutskever, I., Bruna, J., Erhan, D., Goodfellow, I., Fergus, R.: Intriguing properties of neural networks. arXiv preprint arXiv:1312.6199 (2013)
3. Goodfellow, I.J., Shlens, J., Szegedy, C.: Explaining and harnessing adversarial examples. arXiv preprint arXiv:1412.6572 (2014)
4. Moosavi Dezfooli, S.M., Fawzi, A., Frossard, P.: Deepfool: a simple and accurate method to fool deep neural networks. In: Proceedings of 2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR). Number EPFL-CONF-218057 (2016)
5. Carlini, N., Wagner, D.: Towards evaluating the robustness of neural networks. In: Security and Privacy (SP), 2017 IEEE Symposium on, IEEE (2017) 39–57
6. Kurakin, A., Goodfellow, I., Bengio, S.: Adversarial machine learning at scale. In: International Conference on Learning Representations. (2017)
7. Moosavi-Dezfooli, S.M., Fawzi, A., Fawzi, O., Frossard, P.: Universal adversarial perturbations. arXiv preprint (2017)
8. Papernot, N., McDaniel, P., Wu, X., Jha, S., Swami, A.: Distillation as a defense to adversarial perturbations against deep neural networks. In: Security and Privacy (SP), 2016 IEEE Symposium on, IEEE (2016) 582–597
9. Feinman, R., Curtin, R.R., Shintre, S., Gardner, A.B.: Detecting adversarial samples from artifacts. arXiv preprint arXiv:1703.00410 (2017)
10. Madry, A., Makelov, A., Schmidt, L., Tsipras, D., Vladu, A.: Towards deep learning models resistant to adversarial attacks. In: International Conference on Learning Representations. (2018)
11. Sinha, A., Namkoong, H., Duchi, J.: Certifiable distributional robustness with principled adversarial training. In: International Conference on Learning Representations. (2018)
12. Tramr, F., Kurakin, A., Papernot, N., Goodfellow, I., Boneh, D., McDaniel, P.: Ensemble adversarial training: Attacks and defenses. International Conference on Learning Representations (2018)
13. Buckman, J., Roy, A., Raffel, C., Goodfellow, I.: Thermometer encoding: One hot way to resist adversarial examples. In: International Conference on Learning Representations. (2018)
14. Guo, C., Rana, M., Cisse, M., van der Maaten, L.: Countering adversarial images using input transformations. International Conference on Learning Representations (2018) accepted as poster.
15. Song, Y., Kim, T., Nowozin, S., Ermon, S., Kushman, N.: Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. In: International Conference on Learning Representations. (2018)
16. Athalye, A., Carlini, N., Wagner, D.: Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. arXiv preprint arXiv:1802.00420 (2018)
17. Liao, F., Liang, M., Dong, Y., Pang, T., Zhu, J., Hu, X.: Defense against adversarial attacks using high-level representation guided denoiser. arXiv preprint arXiv:1712.02976 (2017)

18. Xie, C., Wang, J., Zhang, Z., Ren, Z., Yuille, A.: Mitigating adversarial effects through randomization. In: International Conference on Learning Representations. (2018)
19. Palacio, S., Folz, J., Hees, J., Raue, F., Borth, D., Dengel, A.: What do deep networks like to see. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR). (June 2018)
20. Strehl, A., Ghosh, J.: Cluster ensembles—a knowledge reuse framework for combining multiple partitions. Journal of machine learning research **3**(Dec) (2002) 583–617
21. Nguyen, A., Dosovitskiy, A., Yosinski, J., Brox, T., Clune, J.: Synthesizing the preferred inputs for neurons in neural networks via deep generator networks. In: Advances in Neural Information Processing Systems. (2016) 3387–3395
22. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. (2016) 770–778
23. Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the inception architecture for computer vision. CoRR **abs/1512.00567** (2015)
24. Kurakin, A., Goodfellow, I.J., Bengio, S.: Adversarial examples in the physical world. CoRR **abs/1607.02533** (2016)
25. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., Berg, A.C., Fei-Fei, L.: ImageNet Large Scale Visual Recognition Challenge. International Journal of Computer Vision (IJCV) **115**(3) (2015) 211–252
26. Meng, D., Chen, H.: Magnet: a two-pronged defense against adversarial examples. In: Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security, ACM (2017) 135–147
27. Liu, Y., Chen, X., Liu, C., Song, D.: Delving into transferable adversarial examples and black-box attacks. arXiv preprint arXiv:1611.02770 (2016)
28. Brown, T.B., Mané, D., Roy, A., Abadi, M., Gilmer, J.: Adversarial patch. arXiv preprint arXiv:1712.09665 (2017)
29. Dhillon, G.S., Azizzadenesheli, K., Bernstein, J.D., Kossaifi, J., Khanna, A., Lipton, Z.C., Anandkumar, A.: Stochastic activation pruning for robust adversarial defense. In: International Conference on Learning Representations. (2018)
30. Wang, Q., Guo, W., Zhang, K., Ororbia, I., Alexander, G., Xing, X., Giles, C.L., Liu, X.: Learning adversary-resistant deep neural networks. arXiv preprint arXiv:1612.01401 (2016)
31. Pouya Samangouei, Maya Kabkab, R.C.: Defense-gan: Protecting classifiers against adversarial attacks using generative models. International Conference on Learning Representations (2018) accepted as poster.
32. Thomee, B., Shamma, D.A., Friedland, G., Elizalde, B., Ni, K., Poland, D., Borth, D., Li, L.: The new data and new challenges in multimedia research. CoRR **abs/1503.01817** (2015)
33. Badrinarayanan, V., Handa, A., Cipolla, R.: Segnet: A deep convolutional encoder-decoder architecture for robust semantic pixel-wise labelling. arXiv preprint arXiv:1505.07293 (2015)
34. Wang, Z., Bovik, A.C., Sheikh, H.R., Simoncelli, E.P.: Image quality assessment: from error visibility to structural similarity. IEEE Transactions on Image Processing **13**(4) (April 2004) 600–612
35. Zeiler, M.D., Fergus, R.: Visualizing and understanding convolutional networks. In: European conference on computer vision, Springer (2014) 818–833

36. Krizhevsky, A.: One weird trick for parallelizing convolutional neural networks. CoRR **abs/1404.5997** (2014)
37. Simonyan, K., Zisserman, A.: Very deep convolutional networks for large-scale image recognition. CoRR **abs/1409.1556** (2014)

## A   On the Relationship between $\mathcal{P}_{x,f}$ and $\mathcal{P}_{x,f\circ g}$: Formal Proof

Experiments of Section 3.2 show that the domain of attacks is being shifted such that their range does not lie on the vulnerable space $\mathcal{P}_{x,f}$. We conduct a formal analysis of the relationship between perturbation spaces $\mathcal{P}_{x,f}$ and $\mathcal{P}_{x,f\circ g}$ and show that one is a subset of the other. There are a few simplifications required for the proof, which will be accounted for at the end of this section. First, we use the results of BIM for an $\mathcal{L}_2 = 0.05$ tested on ResNet 50 as a reference. Values for the baseline (no defense), White-Box$^+$ and, Gray-Box$^-$ are 0.0, 0.0 and 0.672041 respectively. Computing the relative drop in performance (*i.e.*, normalizing by 0.75004: the accuracy under no attack) yields 0.0, 0.0 and 0.8960. To simplify the handling of fuzzy sets, we assign the membership functions $\mu(\alpha(x))$ for each perturbation set to be either 0 or 1 if the fooling ratio is below random chance or above 0.8 respectively. With this in mind, we can define four axioms that come from the domain of the adversarial attack $\alpha : \nabla_f \cup \nabla_{f\circ g} \rightarrow \mathcal{P}$, and the aforementioned experiments with simplified membership functions. The proof is a simple proof by contradiction with case analysis on the first disjunction.

*Proof.*

| | |
|---|---|
| 1. $\forall_x (x \in \nabla_f \lor x \in \nabla_{f\circ g})$ | (Def. domain of $\alpha$) |
| 2. $\forall_x (x \in \nabla_f \rightarrow \alpha(x) \in \mathcal{P}_{x,f})$ | (Baseline Exp.) |
| 3. $\forall_x (x \in \nabla_f \rightarrow \alpha(x) \in \mathcal{P}_{x,f\circ g})$ | (White-Box$^+$) |
| 4. $\forall_x (x \in \nabla_{f\circ g} \rightarrow \alpha(x) \notin \mathcal{P}_{x,f})$ | (Gray-Box$^-$) |

| | |
|---|---|
| 5. $\exists_x (\alpha(x) \in \mathcal{P}_{x,f} \land \alpha(x) \notin \mathcal{P}_{x,f\circ g})$ | (Assumption, $\not\subseteq$) |
| 6. $\alpha(a) \in \mathcal{P}_{a,f} \land \alpha(a) \notin \mathcal{P}_{a,f\circ g}$ | (Skolemization $x \rightarrow a, 5$) |
| 7. $a \in \nabla_f \lor a \in \nabla_{f\circ g}$ | (U.I. $x \rightarrow a$, 1) |
| 8. $a \in \nabla_f \rightarrow \alpha(a) \in \mathcal{P}_{a,f\circ g}$ | (U.I. $x \rightarrow a$, 3) |
| 9. $a \in \nabla_{f\circ g} \rightarrow \alpha(a) \notin \mathcal{P}_{a,f}$ | (U.I. $x \rightarrow a$, 4) |

| | |
|---|---|
| 10. $a \in \nabla_{f\circ g}$ | (Assumption, 7) |
| 11. $\alpha(a) \notin \mathcal{P}_{a,f}$ | ($\rightarrow$ 10, 9) |
| 12. $\alpha(a) \notin \mathcal{P}_{a,f} \land \alpha(a) \in \mathcal{P}_{a,f}$ | ($\land$, 11, 6) |
| 13. Contradiction! | ⚡ |

| | |
|---|---|
| 14. $\neg(a \in \nabla_{f\circ g})$ | (Q.E.A. 10) |
| 15. $a \in \nabla_f$ | (D.Syllogism, 14, 7) |
| 16. $\alpha(a) \in \mathcal{P}_{a,f\circ g}$ | ($\rightarrow$ 15, 8) |
| 17. $\alpha(a) \notin \mathcal{P}_{a,f\circ g} \land \alpha(a) \in \mathcal{P}_{a,f\circ g}$ | ($\land$, 6, 16) |
| 18. Contradiction! | ⚡ |

| | |
|---|---|
| 19. $\neg\exists_x (\alpha(x) \in \mathcal{P}_{x,f} \land \alpha(x) \notin \mathcal{P}_{x,f\circ g})$ | (Q.E.A., 5) |
| 20. $\forall_x \neg(\alpha(x) \in \mathcal{P}_{x,f} \land \alpha(x) \notin \mathcal{P}_{x,f\circ g})$ | ($\neg\exists$, 19) |
| 21. $\forall_x (\alpha(x) \notin \mathcal{P}_{x,f} \lor \alpha(x) \in \mathcal{P}_{x,f\circ g})$ | (Distr. $\neg$, 20) |

22. $\forall_x(\alpha(x) \in \mathcal{P}_{x,f} \to \alpha(x) \in \mathcal{P}_{x,f \circ g})$      (MI $\to$., 21)

23. $\mathcal{P}_{x,f} \subseteq \mathcal{P}_{x,f \circ g}$      (Def. $\subseteq$, 22)

Going back to the simplified membership function, it follows that different reference experiments (network, attack and $\mathcal{L}_2$) will naturally yield different results. This is especially true if we take the raw accuracy as the membership function, instead of the simplified one. However, one can argue that the subset relationship is, in general terms, valid since overall, most experiments under the simplified membership function yield the same axioms). This is why we denote such a relationship by the approximate subset relationship $\subseteq$ to refer to this result.



Fig. 9: Analysis of the relationship between gradient space and perturbation space. By doing case analysis (middle) and a proof by contradiction, we infer that the perturbation space approximates $\mathcal{P}_{x,f} \subsetneq \mathcal{P}_{x,f \circ g}$

A similar analysis can be done graphically as shown in Figure 9 (middle). Case A is covered by Gray-Box$^-$ experiments which show that $\mathcal{P}_{x,f}$ can only be reached 10.4% of the time, which we know now, actually corresponds to case B due to the inclusion $\mathcal{P}_{x,f} \subsetneq \mathcal{P}_{x,f \circ g}$. This leaves cases C, D which are covered by White-Box experiments. Here, using the BIM on ResNet 50 and $\mathcal{L}_2$ = 0.05 as reference, yields that $\alpha(x) \in \mathcal{P}_{x,f \circ g}$ with a membership $\mu(\alpha(x)) = 1 - (0.5/0.75) = 0.333$. That leaves case D with perturbations in the remaining 0.667 in $\mathcal{P} - \{\mathcal{P}_{x,f} \cup \mathcal{P}_{x,f \circ g}\}$.

Likewise, for cases where the domain is $\nabla_f$, we see that they all fall into $\mathcal{P}_{x,f}$ which we know lies in $\mathcal{P}_{x,f \circ g}$ hence, they all fall into case F, eliminating samples falling into the remaining ones E, G, H.

# B    Perturbed Images and their Reconstructions

The following figures (10-16) show attack attempts on ResNet 50 ($R$) with correctly classified images from our random shuffle of the ImageNet validation set. $x$ (column 1) denotes the clean image and $\hat{x}$ its perturbed variants. Columns 2-7 show a combination of attack method (FGSM, BIM, CW) and gradient sources ($\nabla_x R(x)$, $\nabla_x R \circ \mathbb{A}_A(x)$). Below are perturbations $\delta$ (row 2), reconstructions $\hat{x}' = \mathbb{A}_R(\hat{x})$ (row 3), and remaining perturbation in reconstructions $\delta' = \mathbb{A}_R(\hat{x}) - \mathbb{A}_R(x)$ (row 4). Perturbation images are subject to histogram equalization for increases visibility. These examples further illustrate the structural nature of attacks through S2SNets.



Fig. 10: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.032, 0.032, 0.098, 0.115, 0.052, 0.064.

Fig. 11: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.036, 0.036, 0.111, 0.104, 0.059, 0.172.



Fig. 12: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.028, 0.028, 0.084, 0.105, 0.045, 0.095.

Fig. 13: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.028, 0.028, 0.088, 0.119, 0.046, 0.054.



Fig. 14: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.025, 0.025, 0.077, 0.124, 0.041, 0.212.

Fig. 15: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.030, 0.030, 0.095, 0.107, 0.050, 0.150.



Fig. 16: Example of different attacks on ResNet 50 ($R$) with and without $\mathbb{A}_R$, along with perturbations and reconstructions. $\mathcal{L}_2$ (left to right): 0.018, 0.018, 0.053, 0.067, 0.029, 0.045.

## C    Raw Experiment Data

This section contains raw values for all experiments we conducted. A subset was used to create the plots seen in the paper, but we also provide results we ended up not using for our argumentation. In the following tables, $\mathcal{L}_2$ denotes the normalized $L2$ dissimilarity (Equation 1) wrt. reference image $x_1$ and another image $x_2$. We also provide $L_\infty$ values, the usual infinity norm divided by 255 (Equation 2). Settings for attacks besides $\epsilon$ are defined in Section 4.

$$\mathcal{L}_2(x_1, x_2) = \frac{||x_1 - x_2||_2}{||x_1||_2} \tag{1}$$

$$L_\infty(x_1, x_2) = \frac{max(|x_1 - x_2|)}{255} \tag{2}$$

Tables 2, 3, and 4 contain White-Box results for the respective attack method. They cover experiments described in Section 4.1 Tables 5, 6, and 7 contain results attacks where the network that supplied the gradients (source network) that the attack is based on is different from the network that is attacked (target network). They cover White-Box$^+$ and Gray-Box$^-$ experiments described in Section 4.1 and 4.2. Finally, Table 8 contains results on BIM attacks mitigated by different types and strengths ($\eta$) of noise.

Table 2: Accuracy and distance statistics for FGSM attacks.

| Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) | $\mathcal{L}_2$ | $L_\infty$ |
|---|---|---|---|---|---|
| A | 0.5 | 56.48 | 20.30 | 0.00413 | 0.00196 |
| A | 1.0 | 56.48 | 8.60 | 0.00825 | 0.00392 |
| A | 2.0 | 56.48 | 2.80 | 0.01647 | 0.00784 |
| A | 4.0 | 56.48 | 1.14 | 0.03283 | 0.01569 |
| A | 8.0 | 56.48 | 0.83 | 0.06533 | 0.03137 |
| A | 16.0 | 56.48 | 0.69 | 0.12945 | 0.06275 |
| I | 0.5 | 76.86 | 32.58 | 0.00413 | 0.00196 |
| I | 1.0 | 76.86 | 25.32 | 0.00825 | 0.00392 |
| I | 2.0 | 76.86 | 21.32 | 0.01647 | 0.00784 |
| I | 4.0 | 76.86 | 20.04 | 0.03284 | 0.01569 |
| I | 4.0 | 76.86 | 20.04 | 0.03284 | 0.01569 |
| I | 8.0 | 76.86 | 21.37 | 0.06535 | 0.03137 |
| I | 16.0 | 76.86 | 25.09 | 0.12948 | 0.06275 |
| R | 0.5 | 76.03 | 23.49 | 0.00413 | 0.00196 |
| R | 1.0 | 76.03 | 12.58 | 0.00826 | 0.00392 |
| R | 2.0 | 76.03 | 8.00 | 0.01648 | 0.00784 |
| R | 4.0 | 76.03 | 6.90 | 0.03284 | 0.01569 |
| R | 8.0 | 76.03 | 8.02 | 0.06535 | 0.03137 |
| R | 16.0 | 76.03 | 10.37 | 0.12947 | 0.06275 |
| V | 0.5 | 73.28 | 8.82 | 0.00413 | 0.00196 |
| V | 1.0 | 73.28 | 3.75 | 0.00826 | 0.00392 |

Table 2: Accuracy and distance statistics for FGSM attacks.

| Network | $\epsilon$ | top-1 acc. (%) | ...(adversarial) | $\mathcal{L}_2$ | $L_\infty$ |
|---|---|---|---|---|---|
| $V$ | 2.0 | 73.28 | 2.74 | 0.01648 | 0.00784 |
| $V$ | 4.0 | 73.28 | 2.85 | 0.03285 | 0.01569 |
| $V$ | 8.0 | 73.28 | 3.79 | 0.06537 | 0.03137 |
| $V$ | 16.0 | 73.28 | 6.64 | 0.12951 | 0.06275 |
| $A \circ \mathbb{A}_A$ | 0.5 | 56.38 | 54.90 | 0.00413 | 0.00196 |
| $A \circ \mathbb{A}_A$ | 1.0 | 56.38 | 53.79 | 0.00827 | 0.00392 |
| $A \circ \mathbb{A}_A$ | 2.0 | 56.38 | 52.37 | 0.01650 | 0.00784 |
| $A \circ \mathbb{A}_A$ | 4.0 | 56.38 | 50.68 | 0.03289 | 0.01569 |
| $A \circ \mathbb{A}_A$ | 8.0 | 56.38 | 48.56 | 0.06545 | 0.03137 |
| $A \circ \mathbb{A}_A$ | 16.0 | 56.38 | 42.86 | 0.12970 | 0.06275 |
| $I \circ \mathbb{A}_I$ | 0.5 | 76.79 | 75.46 | 0.00413 | 0.00196 |
| $I \circ \mathbb{A}_I$ | 1.0 | 76.79 | 74.51 | 0.00825 | 0.00392 |
| $I \circ \mathbb{A}_I$ | 2.0 | 76.79 | 73.41 | 0.01647 | 0.00784 |
| $I \circ \mathbb{A}_I$ | 4.0 | 76.79 | 71.91 | 0.03283 | 0.01569 |
| $I \circ \mathbb{A}_I$ | 8.0 | 76.79 | 69.89 | 0.06532 | 0.03137 |
| $I \circ \mathbb{A}_I$ | 16.0 | 76.79 | 65.40 | 0.12942 | 0.06275 |
| $R \circ \mathbb{A}_R$ | 0.5 | 75.00 | 73.35 | 0.00413 | 0.00196 |
| $R \circ \mathbb{A}_R$ | 1.0 | 75.00 | 72.33 | 0.00826 | 0.00392 |
| $R \circ \mathbb{A}_R$ | 2.0 | 75.00 | 70.99 | 0.01648 | 0.00784 |
| $R \circ \mathbb{A}_R$ | 4.0 | 75.00 | 69.33 | 0.03285 | 0.01569 |
| $R \circ \mathbb{A}_R$ | 8.0 | 75.00 | 66.65 | 0.06535 | 0.03137 |
| $R \circ \mathbb{A}_R$ | 16.0 | 75.00 | 60.47 | 0.12948 | 0.06275 |
| $V \circ \mathbb{A}_V$ | 0.5 | 68.14 | 65.20 | 0.00413 | 0.00196 |
| $V \circ \mathbb{A}_V$ | 1.0 | 68.14 | 63.38 | 0.00827 | 0.00392 |
| $V \circ \mathbb{A}_V$ | 2.0 | 68.14 | 61.01 | 0.01650 | 0.00784 |
| $V \circ \mathbb{A}_V$ | 4.0 | 68.14 | 58.16 | 0.03289 | 0.01569 |
| $V \circ \mathbb{A}_V$ | 8.0 | 68.14 | 54.85 | 0.06544 | 0.03137 |
| $V \circ \mathbb{A}_V$ | 16.0 | 68.14 | 47.38 | 0.12965 | 0.06275 |
| $I \circ \mathbb{A}_S$ | 0.5 | 73.80 | 71.37 | 0.00413 | 0.00196 |
| $I \circ \mathbb{A}_S$ | 1.0 | 73.80 | 69.71 | 0.00825 | 0.00392 |
| $I \circ \mathbb{A}_S$ | 2.0 | 73.80 | 67.63 | 0.01647 | 0.00784 |
| $I \circ \mathbb{A}_S$ | 4.0 | 73.80 | 65.41 | 0.03283 | 0.01569 |
| $I \circ \mathbb{A}_S$ | 8.0 | 73.80 | 63.10 | 0.06532 | 0.03137 |
| $I \circ \mathbb{A}_S$ | 16.0 | 73.80 | 59.38 | 0.12943 | 0.06275 |
| $R \circ \mathbb{A}_S$ | 0.5 | 72.07 | 68.75 | 0.00413 | 0.00196 |
| $R \circ \mathbb{A}_S$ | 1.0 | 72.07 | 66.33 | 0.00825 | 0.00392 |
| $R \circ \mathbb{A}_S$ | 2.0 | 72.07 | 63.47 | 0.01647 | 0.00784 |
| $R \circ \mathbb{A}_S$ | 4.0 | 72.07 | 59.89 | 0.03283 | 0.01569 |
| $R \circ \mathbb{A}_S$ | 8.0 | 72.07 | 56.03 | 0.06533 | 0.03137 |
| $R \circ \mathbb{A}_S$ | 16.0 | 72.07 | 48.87 | 0.12943 | 0.06275 |

Table 3: Accuracy and distance statistics for BIM attacks.

| Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) | $\mathcal{L}_2$ | $L_\infty$ |
|---|---|---|---|---|---|
| $A$ | 0.5 | 56.48 | 0.11 | 0.02278 | 0.01961 |
| $A$ | 1.0 | 56.48 | 0.05 | 0.03842 | 0.03922 |
| $A$ | 2.0 | 56.48 | 0.03 | 0.06572 | 0.07843 |
| $A$ | 4.0 | 56.48 | 0.02 | 0.11770 | 0.15686 |
| $A$ | 8.0 | 56.48 | 0.01 | 0.22130 | 0.31373 |
| $I$ | 0.5 | 76.86 | 0.50 | 0.01508 | 0.01961 |
| $I$ | 1.0 | 76.86 | 0.20 | 0.02698 | 0.03921 |
| $I$ | 2.0 | 76.86 | 0.15 | 0.05147 | 0.07841 |
| $I$ | 4.0 | 76.86 | 0.10 | 0.10090 | 0.15682 |
| $I$ | 8.0 | 76.86 | 0.08 | 0.19724 | 0.31366 |
| $R$ | 0.5 | 76.03 | 0.10 | 0.01669 | 0.01961 |
| $R$ | 1.0 | 76.03 | 0.05 | 0.02847 | 0.03922 |
| $R$ | 2.0 | 76.03 | 0.02 | 0.05232 | 0.07843 |
| $R$ | 4.0 | 76.03 | 0.01 | 0.10067 | 0.15686 |
| $R$ | 8.0 | 76.03 | 0.00 | 0.19510 | 0.31373 |
| $V$ | 0.5 | 73.28 | 0.34 | 0.01724 | 0.01961 |
| $V$ | 1.0 | 73.28 | 0.27 | 0.02920 | 0.03922 |
| $V$ | 2.0 | 73.28 | 0.23 | 0.05305 | 0.07843 |
| $V$ | 4.0 | 73.28 | 0.18 | 0.10117 | 0.15686 |
| $V$ | 8.0 | 73.28 | 0.12 | 0.19571 | 0.31372 |
| $A \circ \mathbb{A}_A$ | 0.5 | 56.38 | 44.11 | 0.02637 | 0.01961 |
| $A \circ \mathbb{A}_A$ | 1.0 | 56.38 | 35.57 | 0.04537 | 0.03922 |
| $A \circ \mathbb{A}_A$ | 2.0 | 56.38 | 24.47 | 0.07559 | 0.07843 |
| $A \circ \mathbb{A}_A$ | 4.0 | 56.38 | 12.69 | 0.12549 | 0.15686 |
| $A \circ \mathbb{A}_A$ | 8.0 | 56.38 | 4.00 | 0.22330 | 0.31373 |
| $I \circ \mathbb{A}_I$ | 0.5 | 76.79 | 66.48 | 0.02235 | 0.01961 |
| $I \circ \mathbb{A}_I$ | 1.0 | 76.79 | 60.47 | 0.03852 | 0.03922 |
| $I \circ \mathbb{A}_I$ | 2.0 | 76.79 | 51.12 | 0.06634 | 0.07843 |
| $I \circ \mathbb{A}_I$ | 4.0 | 76.79 | 37.65 | 0.11604 | 0.15686 |
| $I \circ \mathbb{A}_I$ | 8.0 | 76.79 | 23.03 | 0.21053 | 0.31373 |
| $R \circ \mathbb{A}_R$ | 0.5 | 75.00 | 62.05 | 0.02445 | 0.01961 |
| $R \circ \mathbb{A}_R$ | 1.0 | 75.00 | 53.64 | 0.04197 | 0.03922 |
| $R \circ \mathbb{A}_R$ | 2.0 | 75.00 | 41.21 | 0.07040 | 0.07843 |
| $R \circ \mathbb{A}_R$ | 4.0 | 75.00 | 25.93 | 0.11964 | 0.15686 |
| $R \circ \mathbb{A}_R$ | 8.0 | 75.00 | 11.78 | 0.21596 | 0.31373 |
| $V \circ \mathbb{A}_V$ | 0.5 | 68.14 | 46.16 | 0.02467 | 0.01961 |
| $V \circ \mathbb{A}_V$ | 1.0 | 68.14 | 34.45 | 0.04230 | 0.03922 |
| $V \circ \mathbb{A}_V$ | 2.0 | 68.14 | 22.14 | 0.07086 | 0.07843 |
| $V \circ \mathbb{A}_V$ | 4.0 | 68.14 | 12.65 | 0.11932 | 0.15686 |
| $V \circ \mathbb{A}_V$ | 8.0 | 68.14 | 6.33 | 0.21275 | 0.31373 |
| $R \circ \mathbb{A}_S$ | 0.5 | 72.07 | 45.11 | 0.02507 | 0.01961 |
| $R \circ \mathbb{A}_S$ | 1.0 | 72.07 | 31.29 | 0.04294 | 0.03922 |

Table 3: Accuracy and distance statistics for BIM attacks.

| Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) | $\mathcal{L}_2$ | $L_\infty$ |
|---|---|---|---|---|---|
| $R \circ \mathbb{A}_S$ | 2.0 | 72.07 | 18.32 | 0.07170 | 0.07843 |
| $R \circ \mathbb{A}_S$ | 4.0 | 72.07 | 9.86 | 0.12014 | 0.15686 |
| $R \circ \mathbb{A}_S$ | 8.0 | 72.07 | 5.42 | 0.21249 | 0.31373 |
| $R \circ \mathbb{A}_S$ | 16.0 | 72.07 | 1.09 | 0.40976 | 0.62744 |

Table 4: Accuracy and distance statistics for CW attacks.

| Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) | $\mathcal{L}_2$ | $L_\infty$ |
|---|---|---|---|---|---|
| $I$ | 0.5 | 76.86 | 0.16 | 0.00805 | 0.00587 |
| $I$ | 1.0 | 76.86 | 0.29 | 0.01475 | 0.00966 |
| $I$ | 2.0 | 76.86 | 0.99 | 0.03084 | 0.01848 |
| $I$ | 4.0 | 76.86 | 2.13 | 0.05731 | 0.03350 |
| $R$ | 0.5 | 76.03 | 0.10 | 0.00752 | 0.00511 |
| $R$ | 1.0 | 76.03 | 0.02 | 0.01387 | 0.00850 |
| $R$ | 2.0 | 76.03 | 0.42 | 0.02964 | 0.01644 |
| $R$ | 4.0 | 76.03 | 1.32 | 0.05459 | 0.02917 |
| $I \circ \mathbb{A}_I$ | 0.5 | 76.80 | 44.49 | 0.07357 | 0.11869 |
| $I \circ \mathbb{A}_I$ | 1.0 | 76.80 | 33.58 | 0.10195 | 0.17756 |
| $I \circ \mathbb{A}_I$ | 2.0 | 76.80 | 25.19 | 0.13851 | 0.24872 |
| $I \circ \mathbb{A}_I$ | 4.0 | 76.80 | 19.72 | 0.18622 | 0.31932 |
| $R \circ \mathbb{A}_R$ | 0.5 | 75.01 | 36.66 | 0.07310 | 0.11166 |
| $R \circ \mathbb{A}_R$ | 1.0 | 75.01 | 26.59 | 0.09461 | 0.15410 |
| $R \circ \mathbb{A}_R$ | 2.0 | 75.01 | 20.63 | 0.11958 | 0.19562 |
| $R \circ \mathbb{A}_R$ | 4.0 | 75.01 | 16.68 | 0.15458 | 0.23814 |
| $R \circ \mathbb{A}_S$ | 0.5 | 72.07 | 20.26 | 0.04971 | 0.06407 |
| $R \circ \mathbb{A}_S$ | 1.0 | 72.07 | 14.08 | 0.06044 | 0.07838 |
| $R \circ \mathbb{A}_S$ | 2.0 | 72.07 | 11.27 | 0.07758 | 0.09839 |
| $R \circ \mathbb{A}_S$ | 4.0 | 72.07 | 10.45 | 0.10965 | 0.13377 |

Table 5: Accuracy for FGSM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| $A$ | $I$ | 0.5 | 76.86 | 56.11 |
| $A$ | $R$ | 0.5 | 76.03 | 55.93 |
| $A$ | $V$ | 0.5 | 73.28 | 56.01 |
| $A$ | $A \circ \mathbb{A}_A$ | 0.5 | 56.38 | 56.26 |
| $A$ | $I$ | 1.0 | 76.86 | 55.86 |
| $A$ | $R$ | 1.0 | 76.03 | 55.43 |
| $A$ | $V$ | 1.0 | 73.28 | 55.58 |

Table 5: Accuracy for FGSM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|:---:|:---:|:---:|:---:|:---:|
| $A$ | $A \circ \mathbb{A}_A$ | 1.0 | 56.38 | 56.19 |
| $A$ | $I$ | 2.0 | 76.86 | 55.02 |
| $A$ | $R$ | 2.0 | 76.03 | 54.16 |
| $A$ | $V$ | 2.0 | 73.28 | 54.38 |
| $A$ | $A \circ \mathbb{A}_A$ | 2.0 | 56.38 | 55.89 |
| $A$ | $I$ | 4.0 | 76.86 | 52.48 |
| $A$ | $R$ | 4.0 | 76.03 | 50.82 |
| $A$ | $V$ | 4.0 | 73.28 | 51.57 |
| $A$ | $A \circ \mathbb{A}_A$ | 4.0 | 56.38 | 55.25 |
| $A$ | $I$ | 8.0 | 76.86 | 44.76 |
| $A$ | $R$ | 8.0 | 76.03 | 41.18 |
| $A$ | $V$ | 8.0 | 73.28 | 43.48 |
| $A$ | $A \circ \mathbb{A}_A$ | 8.0 | 56.38 | 52.33 |
| $A$ | $A \circ \mathbb{A}_A$ | 16.0 | 56.38 | 41.49 |
| $I$ | $A$ | 0.5 | 56.48 | 76.31 |
| $I$ | $R$ | 0.5 | 76.03 | 74.63 |
| $I$ | $V$ | 0.5 | 73.28 | 75.11 |
| $I$ | $I \circ \mathbb{A}_I$ | 0.5 | 76.79 | 76.71 |
| $I$ | $A$ | 1.0 | 56.48 | 75.70 |
| $I$ | $R$ | 1.0 | 76.03 | 72.26 |
| $I$ | $V$ | 1.0 | 73.28 | 73.29 |
| $I$ | $I \circ \mathbb{A}_I$ | 1.0 | 76.79 | 76.45 |
| $I$ | $A$ | 2.0 | 56.48 | 74.27 |
| $I$ | $R$ | 2.0 | 76.03 | 68.06 |
| $I$ | $V$ | 2.0 | 73.28 | 69.89 |
| $I$ | $I \circ \mathbb{A}_I$ | 2.0 | 76.79 | 75.90 |
| $I$ | $A$ | 4.0 | 56.48 | 71.20 |
| $I$ | $R$ | 4.0 | 76.03 | 61.82 |
| $I$ | $V$ | 4.0 | 73.28 | 64.71 |
| $I$ | $I \circ \mathbb{A}_I$ | 4.0 | 76.79 | 74.65 |
| $I$ | $A$ | 8.0 | 56.48 | 64.46 |
| $I$ | $R$ | 8.0 | 76.03 | 54.70 |
| $I$ | $V$ | 8.0 | 73.28 | 58.73 |
| $I$ | $I \circ \mathbb{A}_I$ | 8.0 | 76.79 | 72.45 |
| $I$ | $I \circ \mathbb{A}_I$ | 16.0 | 76.79 | 68.18 |
| $R$ | $A$ | 0.5 | 56.48 | 75.29 |
| $R$ | $I$ | 0.5 | 76.86 | 74.01 |
| $R$ | $V$ | 0.5 | 73.28 | 73.22 |
| $R$ | $R \circ \mathbb{A}_R$ | 0.5 | 75.00 | 75.71 |
| $R$ | $R \circ \mathbb{A}_S$ | 0.5 | 72.07 | 75.03 |
| $R$ | $A$ | 1.0 | 56.48 | 74.39 |
| $R$ | $I$ | 1.0 | 76.86 | 71.87 |

Table 5: Accuracy for FGSM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| $R$ | $V$ | 1.0 | 73.28 | 70.38 |
| $R$ | $R \circ \mathbb{A}_R$ | 1.0 | 75.00 | 75.42 |
| $R$ | $R \circ \mathbb{A}_S$ | 1.0 | 72.07 | 73.89 |
| $R$ | $A$ | 2.0 | 56.48 | 72.36 |
| $R$ | $I$ | 2.0 | 76.86 | 68.23 |
| $R$ | $V$ | 2.0 | 73.28 | 65.24 |
| $R$ | $R \circ \mathbb{A}_R$ | 2.0 | 75.00 | 74.68 |
| $R$ | $R \circ \mathbb{A}_S$ | 2.0 | 72.07 | 71.70 |
| $R$ | $A$ | 4.0 | 56.48 | 68.07 |
| $R$ | $I$ | 4.0 | 76.86 | 63.06 |
| $R$ | $V$ | 4.0 | 73.28 | 58.30 |
| $R$ | $R \circ \mathbb{A}_R$ | 4.0 | 75.00 | 73.01 |
| $R$ | $R \circ \mathbb{A}_S$ | 4.0 | 72.07 | 67.98 |
| $R$ | $A$ | 8.0 | 56.48 | 58.65 |
| $R$ | $I$ | 8.0 | 76.86 | 56.61 |
| $R$ | $V$ | 8.0 | 73.28 | 51.20 |
| $R$ | $R \circ \mathbb{A}_R$ | 8.0 | 75.00 | 69.59 |
| $R$ | $R \circ \mathbb{A}_S$ | 8.0 | 72.07 | 62.50 |
| $R$ | $R \circ \mathbb{A}_R$ | 16.0 | 75.00 | 62.54 |
| $V$ | $A$ | 0.5 | 56.48 | 72.30 |
| $V$ | $I$ | 0.5 | 76.86 | 70.97 |
| $V$ | $R$ | 0.5 | 76.03 | 69.25 |
| $V$ | $V \circ \mathbb{A}_V$ | 0.5 | 68.14 | 72.52 |
| $V$ | $A$ | 1.0 | 56.48 | 71.19 |
| $V$ | $I$ | 1.0 | 76.86 | 68.69 |
| $V$ | $R$ | 1.0 | 76.03 | 65.32 |
| $V$ | $V \circ \mathbb{A}_V$ | 1.0 | 68.14 | 71.79 |
| $V$ | $A$ | 2.0 | 56.48 | 68.84 |
| $V$ | $I$ | 2.0 | 76.86 | 64.78 |
| $V$ | $R$ | 2.0 | 76.03 | 59.18 |
| $V$ | $V \circ \mathbb{A}_V$ | 2.0 | 68.14 | 70.12 |
| $V$ | $A$ | 4.0 | 56.48 | 63.51 |
| $V$ | $I$ | 4.0 | 76.86 | 59.33 |
| $V$ | $R$ | 4.0 | 76.03 | 51.14 |
| $V$ | $V \circ \mathbb{A}_V$ | 4.0 | 68.14 | 67.02 |
| $V$ | $A$ | 8.0 | 56.48 | 52.28 |
| $V$ | $I$ | 8.0 | 76.86 | 51.99 |
| $V$ | $R$ | 8.0 | 76.03 | 42.35 |
| $V$ | $V \circ \mathbb{A}_V$ | 8.0 | 68.14 | 62.11 |
| $V$ | $V \circ \mathbb{A}_V$ | 16.0 | 68.14 | 53.49 |
| $A \circ \mathbb{A}_A$ | $A$ | 0.5 | 56.48 | 41.39 |
| $A \circ \mathbb{A}_A$ | $A$ | 1.0 | 56.48 | 29.07 |

Table 5: Accuracy for FGSM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| $A \circ \mathbb{A}_A$ | $A$ | 2.0 | 56.48 | 14.42 |
| $A \circ \mathbb{A}_A$ | $A$ | 4.0 | 56.48 | 5.32 |
| $A \circ \mathbb{A}_A$ | $A$ | 8.0 | 56.48 | 2.23 |
| $A \circ \mathbb{A}_A$ | $A$ | 16.0 | 56.48 | 1.34 |
| $I \circ \mathbb{A}_I$ | $I$ | 0.5 | 76.86 | 62.88 |
| $I \circ \mathbb{A}_I$ | $I$ | 1.0 | 76.86 | 51.84 |
| $I \circ \mathbb{A}_I$ | $I$ | 2.0 | 76.86 | 40.63 |
| $I \circ \mathbb{A}_I$ | $I$ | 4.0 | 76.86 | 33.02 |
| $I \circ \mathbb{A}_I$ | $I$ | 8.0 | 76.86 | 29.92 |
| $I \circ \mathbb{A}_I$ | $I$ | 16.0 | 76.86 | 29.42 |
| $R \circ \mathbb{A}_R$ | $R$ | 0.5 | 76.03 | 54.46 |
| $R \circ \mathbb{A}_R$ | $R$ | 1.0 | 76.03 | 39.20 |
| $R \circ \mathbb{A}_R$ | $R$ | 2.0 | 76.03 | 23.53 |
| $R \circ \mathbb{A}_R$ | $R$ | 4.0 | 76.03 | 14.43 |
| $R \circ \mathbb{A}_R$ | $R$ | 8.0 | 76.03 | 11.67 |
| $R \circ \mathbb{A}_R$ | $R$ | 16.0 | 76.03 | 11.25 |
| $V \circ \mathbb{A}_V$ | $V$ | 0.5 | 73.28 | 35.55 |
| $V \circ \mathbb{A}_V$ | $V$ | 1.0 | 73.28 | 18.50 |
| $V \circ \mathbb{A}_V$ | $V$ | 2.0 | 73.28 | 8.27 |
| $V \circ \mathbb{A}_V$ | $V$ | 4.0 | 73.28 | 5.41 |
| $V \circ \mathbb{A}_V$ | $V$ | 8.0 | 73.28 | 5.39 |
| $V \circ \mathbb{A}_V$ | $V$ | 16.0 | 73.28 | 7.72 |
| $R \circ \mathbb{A}_S$ | $R$ | 0.5 | 76.03 | 47.90 |
| $R \circ \mathbb{A}_S$ | $R$ | 1.0 | 76.03 | 31.81 |
| $R \circ \mathbb{A}_S$ | $R$ | 2.0 | 76.03 | 18.28 |
| $R \circ \mathbb{A}_S$ | $R$ | 4.0 | 76.03 | 11.96 |
| $R \circ \mathbb{A}_S$ | $R$ | 8.0 | 76.03 | 10.40 |

Table 6: Accuracy for BIM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| $A$ | $I$ | 0.5 | 76.86 | 55.04 |
| $A$ | $R$ | 0.5 | 76.03 | 53.95 |
| $A$ | $V$ | 0.5 | 73.28 | 54.42 |
| $A$ | $A \circ \mathbb{A}_A$ | 0.5 | 56.38 | 54.61 |
| $A$ | $I$ | 1.0 | 76.86 | 53.51 |
| $A$ | $R$ | 1.0 | 76.03 | 51.69 |
| $A$ | $V$ | 1.0 | 73.28 | 52.67 |
| $A$ | $A \circ \mathbb{A}_A$ | 1.0 | 56.38 | 51.96 |
| $A$ | $I$ | 2.0 | 76.86 | 48.11 |
| $A$ | $R$ | 2.0 | 76.03 | 45.05 |

Table 6: Accuracy for BIM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ...(adversarial) |
|---|---|---|---|---|
| $A$ | $V$ | 2.0 | 73.28 | 47.17 |
| $A$ | $A \circ \mathbb{A}_A$ | 2.0 | 56.38 | 46.29 |
| $A$ | $I$ | 4.0 | 76.86 | 33.50 |
| $A$ | $R$ | 4.0 | 76.03 | 29.14 |
| $A$ | $V$ | 4.0 | 73.28 | 34.14 |
| $A$ | $A \circ \mathbb{A}_A$ | 4.0 | 56.38 | 34.28 |
| $A$ | $I$ | 8.0 | 76.86 | 11.08 |
| $A$ | $R$ | 8.0 | 76.03 | 8.70 |
| $A$ | $V$ | 8.0 | 73.28 | 12.73 |
| $A$ | $A \circ \mathbb{A}_A$ | 8.0 | 56.38 | 14.31 |
| $I$ | $A$ | 0.5 | 56.48 | 71.70 |
| $I$ | $R$ | 0.5 | 76.03 | 67.32 |
| $I$ | $V$ | 0.5 | 73.28 | 70.39 |
| $I$ | $I \circ \mathbb{A}_I$ | 0.5 | 76.79 | 74.27 |
| $I$ | $A$ | 1.0 | 56.48 | 67.01 |
| $I$ | $R$ | 1.0 | 76.03 | 60.49 |
| $I$ | $V$ | 1.0 | 73.28 | 65.72 |
| $I$ | $I \circ \mathbb{A}_I$ | 1.0 | 76.79 | 71.67 |
| $I$ | $A$ | 2.0 | 56.48 | 58.45 |
| $I$ | $R$ | 2.0 | 76.03 | 49.84 |
| $I$ | $V$ | 2.0 | 73.28 | 57.29 |
| $I$ | $I \circ \mathbb{A}_I$ | 2.0 | 76.79 | 67.21 |
| $I$ | $A$ | 4.0 | 56.48 | 41.90 |
| $I$ | $R$ | 4.0 | 76.03 | 34.22 |
| $I$ | $V$ | 4.0 | 73.28 | 44.44 |
| $I$ | $I \circ \mathbb{A}_I$ | 4.0 | 76.79 | 58.85 |
| $I$ | $A$ | 8.0 | 56.48 | 18.61 |
| $I$ | $R$ | 8.0 | 76.03 | 17.36 |
| $I$ | $V$ | 8.0 | 73.28 | 28.68 |
| $I$ | $I \circ \mathbb{A}_I$ | 8.0 | 76.79 | 46.60 |
| $R$ | $A$ | 0.5 | 56.48 | 68.98 |
| $R$ | $I$ | 0.5 | 76.86 | 68.41 |
| $R$ | $V$ | 0.5 | 73.28 | 64.93 |
| $R$ | $R \circ \mathbb{A}_R$ | 0.5 | 75.00 | 72.31 |
| $R$ | $R \circ \mathbb{A}_S$ | 0.5 | 72.07 | 63.43 |
| $R$ | $A$ | 1.0 | 56.48 | 62.48 |
| $R$ | $I$ | 1.0 | 76.86 | 62.54 |
| $R$ | $V$ | 1.0 | 73.28 | 57.52 |
| $R$ | $R \circ \mathbb{A}_R$ | 1.0 | 75.00 | 68.83 |
| $R$ | $R \circ \mathbb{A}_S$ | 1.0 | 72.07 | 51.71 |
| $R$ | $A$ | 2.0 | 56.48 | 50.59 |
| $R$ | $I$ | 2.0 | 76.86 | 52.16 |

Table 6: Accuracy for BIM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| $R$ | $V$ | 2.0 | 73.28 | 46.09 |
| $R$ | $R \circ \mathbb{A}_R$ | 2.0 | 75.00 | 63.07 |
| $R$ | $R \circ \mathbb{A}_S$ | 2.0 | 72.07 | 35.85 |
| $R$ | $A$ | 4.0 | 56.48 | 30.60 |
| $R$ | $I$ | 4.0 | 76.86 | 37.46 |
| $R$ | $V$ | 4.0 | 73.28 | 32.11 |
| $R$ | $R \circ \mathbb{A}_R$ | 4.0 | 75.00 | 53.10 |
| $R$ | $R \circ \mathbb{A}_S$ | 4.0 | 72.07 | 20.81 |
| $R$ | $A$ | 8.0 | 56.48 | 10.25 |
| $R$ | $I$ | 8.0 | 76.86 | 21.29 |
| $R$ | $V$ | 8.0 | 73.28 | 18.94 |
| $R$ | $R \circ \mathbb{A}_R$ | 8.0 | 75.00 | 36.36 |
| $R$ | $R \circ \mathbb{A}_S$ | 8.0 | 72.07 | 11.40 |
| $V$ | $A$ | 0.5 | 56.48 | 63.81 |
| $V$ | $I$ | 0.5 | 76.86 | 64.25 |
| $V$ | $R$ | 0.5 | 76.03 | 53.11 |
| $V$ | $V \circ \mathbb{A}_V$ | 0.5 | 68.14 | 64.68 |
| $V$ | $A$ | 1.0 | 56.48 | 55.92 |
| $V$ | $I$ | 1.0 | 76.86 | 57.50 |
| $V$ | $R$ | 1.0 | 76.03 | 41.59 |
| $V$ | $V \circ \mathbb{A}_V$ | 1.0 | 68.14 | 56.26 |
| $V$ | $A$ | 2.0 | 56.48 | 42.54 |
| $V$ | $I$ | 2.0 | 76.86 | 45.69 |
| $V$ | $R$ | 2.0 | 76.03 | 26.86 |
| $V$ | $V \circ \mathbb{A}_V$ | 2.0 | 68.14 | 43.38 |
| $V$ | $A$ | 4.0 | 56.48 | 23.27 |
| $V$ | $I$ | 4.0 | 76.86 | 30.33 |
| $V$ | $R$ | 4.0 | 76.03 | 13.02 |
| $V$ | $V \circ \mathbb{A}_V$ | 4.0 | 68.14 | 27.12 |
| $V$ | $A$ | 8.0 | 56.48 | 6.32 |
| $V$ | $I$ | 8.0 | 76.86 | 15.10 |
| $V$ | $R$ | 8.0 | 76.03 | 4.42 |
| $V$ | $V \circ \mathbb{A}_V$ | 8.0 | 68.14 | 13.25 |
| $A \circ \mathbb{A}_A$ | $A$ | 0.5 | 56.48 | 0.95 |
| $A \circ \mathbb{A}_A$ | $A$ | 1.0 | 56.48 | 0.13 |
| $A \circ \mathbb{A}_A$ | $A$ | 2.0 | 56.48 | 0.04 |
| $A \circ \mathbb{A}_A$ | $A$ | 4.0 | 56.48 | 0.02 |
| $A \circ \mathbb{A}_A$ | $A$ | 8.0 | 56.48 | 0.01 |
| $I \circ \mathbb{A}_I$ | $I$ | 0.5 | 76.86 | 23.14 |
| $I \circ \mathbb{A}_I$ | $I$ | 1.0 | 76.86 | 10.72 |
| $I \circ \mathbb{A}_I$ | $I$ | 2.0 | 76.86 | 3.87 |
| $I \circ \mathbb{A}_I$ | $I$ | 4.0 | 76.86 | 1.23 |

Table 6: Accuracy for BIM attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ...(adversarial) |
|---|---|---|---|---|
| $I \circ \mathbb{A}_I$ | $I$ | 8.0 | 76.86 | 0.40 |
| $R \circ \mathbb{A}_R$ | $R$ | 0.5 | 76.03 | 2.65 |
| $R \circ \mathbb{A}_R$ | $R$ | 1.0 | 76.03 | 0.44 |
| $R \circ \mathbb{A}_R$ | $R$ | 2.0 | 76.03 | 0.08 |
| $R \circ \mathbb{A}_R$ | $R$ | 4.0 | 76.03 | 0.01 |
| $R \circ \mathbb{A}_R$ | $R$ | 8.0 | 76.03 | 0.00 |
| $V \circ \mathbb{A}_V$ | $V$ | 0.5 | 73.28 | 0.59 |
| $V \circ \mathbb{A}_V$ | $V$ | 1.0 | 73.28 | 0.33 |
| $V \circ \mathbb{A}_V$ | $V$ | 2.0 | 73.28 | 0.24 |
| $V \circ \mathbb{A}_V$ | $V$ | 4.0 | 73.28 | 0.18 |
| $V \circ \mathbb{A}_V$ | $V$ | 8.0 | 73.28 | 0.12 |
| $R \circ \mathbb{A}_S$ | $R$ | 0.5 | 76.03 | 0.73 |
| $R \circ \mathbb{A}_S$ | $R$ | 1.0 | 76.03 | 0.11 |
| $R \circ \mathbb{A}_S$ | $R$ | 2.0 | 76.03 | 0.03 |
| $R \circ \mathbb{A}_S$ | $R$ | 4.0 | 76.03 | 0.01 |
| $R \circ \mathbb{A}_S$ | $R$ | 8.0 | 76.03 | 0.00 |

Table 7: Accuracy for CW attacks when source and target network differ.

| Source Network | Target Network | $\epsilon$ | top-1 acc. (%) | ...(adversarial) |
|---|---|---|---|---|
| $I$ | $I \circ \mathbb{A}_I$ | 0.5 | 76.80 | 72.90 |
| $I$ | $I \circ \mathbb{A}_I$ | 1.0 | 76.80 | 70.08 |
| $I$ | $I \circ \mathbb{A}_I$ | 2.0 | 76.80 | 65.54 |
| $I$ | $I \circ \mathbb{A}_I$ | 4.0 | 76.80 | 59.46 |
| $R$ | $R \circ \mathbb{A}_R$ | 0.5 | 75.01 | 71.15 |
| $R$ | $R \circ \mathbb{A}_R$ | 1.0 | 75.01 | 67.98 |
| $R$ | $R \circ \mathbb{A}_R$ | 2.0 | 75.01 | 63.42 |
| $R$ | $R \circ \mathbb{A}_R$ | 4.0 | 75.01 | 57.05 |
| $I \circ \mathbb{A}_I$ | $I$ | 0.5 | 76.86 | 64.36 |
| $I \circ \mathbb{A}_I$ | $I$ | 1.0 | 76.86 | 43.53 |
| $I \circ \mathbb{A}_I$ | $I$ | 2.0 | 76.86 | 28.18 |
| $I \circ \mathbb{A}_I$ | $I$ | 4.0 | 76.86 | 20.69 |
| $R \circ \mathbb{A}_R$ | $R$ | 0.5 | 76.03 | 53.83 |
| $R \circ \mathbb{A}_R$ | $R$ | 1.0 | 76.03 | 28.23 |
| $R \circ \mathbb{A}_R$ | $R$ | 2.0 | 76.03 | 12.76 |
| $R \circ \mathbb{A}_R$ | $R$ | 4.0 | 76.03 | 8.50 |

Table 8: BIM attacks on ResNet 50, mitigated with different types and strength ($\eta$) of noise.

| Noise Type | $\epsilon$ | $\eta$ | top-1 acc. (%) | ...(adversarial) |
|---|---|---|---|---|
| Gaussian | 0.5 | 0.5 | 76.03 | 2.77 |
| Gaussian | 0.5 | 1.0 | 76.03 | 3.15 |
| Gaussian | 0.5 | 2.0 | 76.03 | 4.73 |
| Gaussian | 0.5 | 4.0 | 76.03 | 11.15 |
| Gaussian | 0.5 | 8.0 | 76.03 | 29.17 |
| Gaussian | 0.5 | 12.0 | 76.03 | 41.43 |
| Gaussian | 0.5 | 15.0 | 76.03 | 46.46 |
| Gaussian | 0.5 | 20.0 | 76.03 | 50.48 |
| Gaussian | 0.5 | 30.0 | 76.03 | 48.21 |
| Gaussian | 1.0 | 0.5 | 76.03 | 0.45 |
| Gaussian | 1.0 | 1.0 | 76.03 | 0.50 |
| Gaussian | 1.0 | 2.0 | 76.03 | 0.65 |
| Gaussian | 1.0 | 4.0 | 76.03 | 1.88 |
| Gaussian | 1.0 | 8.0 | 76.03 | 11.37 |
| Gaussian | 1.0 | 12.0 | 76.03 | 24.19 |
| Gaussian | 1.0 | 15.0 | 76.03 | 31.81 |
| Gaussian | 1.0 | 20.0 | 76.03 | 40.08 |
| Gaussian | 1.0 | 30.0 | 76.03 | 43.35 |
| Gaussian | 2.0 | 0.5 | 76.03 | 0.08 |
| Gaussian | 2.0 | 1.0 | 76.03 | 0.08 |
| Gaussian | 2.0 | 2.0 | 76.03 | 0.10 |
| Gaussian | 2.0 | 4.0 | 76.03 | 0.16 |
| Gaussian | 2.0 | 8.0 | 76.03 | 1.05 |
| Gaussian | 2.0 | 12.0 | 76.03 | 5.00 |
| Gaussian | 2.0 | 15.0 | 76.03 | 10.38 |
| Gaussian | 2.0 | 20.0 | 76.03 | 20.76 |
| Gaussian | 2.0 | 30.0 | 76.03 | 32.58 |
| Gaussian | 4.0 | 0.5 | 76.03 | 0.01 |
| Gaussian | 4.0 | 1.0 | 76.03 | 0.01 |
| Gaussian | 4.0 | 2.0 | 76.03 | 0.01 |
| Gaussian | 4.0 | 4.0 | 76.03 | 0.02 |
| Gaussian | 4.0 | 8.0 | 76.03 | 0.04 |
| Gaussian | 4.0 | 12.0 | 76.03 | 0.10 |
| Gaussian | 4.0 | 15.0 | 76.03 | 0.30 |
| Gaussian | 4.0 | 20.0 | 76.03 | 1.85 |
| Gaussian | 4.0 | 30.0 | 76.03 | 11.53 |
| Gaussian | 8.0 | 0.5 | 76.03 | 0.00 |
| Gaussian | 8.0 | 1.0 | 76.03 | 0.00 |
| Gaussian | 8.0 | 2.0 | 76.03 | 0.00 |
| Gaussian | 8.0 | 4.0 | 76.03 | 0.00 |
| Gaussian | 8.0 | 8.0 | 76.03 | 0.00 |

Table 8: BIM attacks on ResNet 50, mitigated with different types and strength ($\eta$) of noise.

| Noise Type | $\epsilon$ | $\eta$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| Gaussian | 8.0 | 12.0 | 76.03 | 0.00 |
| Gaussian | 8.0 | 15.0 | 76.03 | 0.01 |
| Gaussian | 8.0 | 20.0 | 76.03 | 0.01 |
| Gaussian | 8.0 | 30.0 | 76.03 | 0.26 |
| Sign | 0.5 | 0.5 | 76.03 | 2.77 |
| Sign | 0.5 | 1.0 | 76.03 | 3.14 |
| Sign | 0.5 | 2.0 | 76.03 | 4.81 |
| Sign | 0.5 | 4.0 | 76.03 | 11.25 |
| Sign | 0.5 | 8.0 | 76.03 | 29.59 |
| Sign | 0.5 | 12.0 | 76.03 | 41.88 |
| Sign | 0.5 | 15.0 | 76.03 | 47.01 |
| Sign | 0.5 | 20.0 | 76.03 | 50.73 |
| Sign | 0.5 | 30.0 | 76.03 | 47.74 |
| Sign | 1.0 | 0.5 | 76.03 | 0.45 |
| Sign | 1.0 | 1.0 | 76.03 | 0.47 |
| Sign | 1.0 | 2.0 | 76.03 | 0.65 |
| Sign | 1.0 | 4.0 | 76.03 | 1.94 |
| Sign | 1.0 | 8.0 | 76.03 | 11.61 |
| Sign | 1.0 | 12.0 | 76.03 | 24.76 |
| Sign | 1.0 | 15.0 | 76.03 | 32.74 |
| Sign | 1.0 | 20.0 | 76.03 | 40.71 |
| Sign | 1.0 | 30.0 | 76.03 | 43.30 |
| Sign | 2.0 | 0.5 | 76.03 | 0.08 |
| Sign | 2.0 | 1.0 | 76.03 | 0.08 |
| Sign | 2.0 | 2.0 | 76.03 | 0.09 |
| Sign | 2.0 | 4.0 | 76.03 | 0.16 |
| Sign | 2.0 | 8.0 | 76.03 | 1.13 |
| Sign | 2.0 | 12.0 | 76.03 | 5.38 |
| Sign | 2.0 | 15.0 | 76.03 | 11.19 |
| Sign | 2.0 | 20.0 | 76.03 | 21.74 |
| Sign | 2.0 | 30.0 | 76.03 | 33.49 |
| Sign | 4.0 | 0.5 | 76.03 | 0.01 |
| Sign | 4.0 | 1.0 | 76.03 | 0.01 |
| Sign | 4.0 | 2.0 | 76.03 | 0.01 |
| Sign | 4.0 | 4.0 | 76.03 | 0.02 |
| Sign | 4.0 | 8.0 | 76.03 | 0.04 |
| Sign | 4.0 | 12.0 | 76.03 | 0.12 |
| Sign | 4.0 | 15.0 | 76.03 | 0.38 |
| Sign | 4.0 | 20.0 | 76.03 | 2.25 |
| Sign | 4.0 | 30.0 | 76.03 | 13.05 |
| Sign | 8.0 | 0.5 | 76.03 | 0.00 |

Table 8: BIM attacks on ResNet 50, mitigated with different types and strength ($\eta$) of noise.

| Noise Type | $\epsilon$ | $\eta$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| Sign | 8.0 | 1.0 | 76.03 | 0.00 |
| Sign | 8.0 | 2.0 | 76.03 | 0.00 |
| Sign | 8.0 | 4.0 | 76.03 | 0.00 |
| Sign | 8.0 | 8.0 | 76.03 | 0.00 |
| Sign | 8.0 | 12.0 | 76.03 | 0.00 |
| Sign | 8.0 | 15.0 | 76.03 | 0.01 |
| Sign | 8.0 | 20.0 | 76.03 | 0.01 |
| Sign | 8.0 | 30.0 | 76.03 | 0.44 |
| Uniform | 0.5 | 0.5 | 76.03 | 2.72 |
| Uniform | 0.5 | 1.0 | 76.03 | 2.81 |
| Uniform | 0.5 | 2.0 | 76.03 | 3.37 |
| Uniform | 0.5 | 4.0 | 76.03 | 5.57 |
| Uniform | 0.5 | 8.0 | 76.03 | 13.82 |
| Uniform | 0.5 | 12.0 | 76.03 | 24.69 |
| Uniform | 0.5 | 15.0 | 76.03 | 31.86 |
| Uniform | 0.5 | 20.0 | 76.03 | 40.42 |
| Uniform | 0.5 | 30.0 | 76.03 | 49.01 |
| Uniform | 1.0 | 0.5 | 76.03 | 0.44 |
| Uniform | 1.0 | 1.0 | 76.03 | 0.44 |
| Uniform | 1.0 | 2.0 | 76.03 | 0.49 |
| Uniform | 1.0 | 4.0 | 76.03 | 0.76 |
| Uniform | 1.0 | 8.0 | 76.03 | 2.69 |
| Uniform | 1.0 | 12.0 | 76.03 | 8.04 |
| Uniform | 1.0 | 15.0 | 76.03 | 13.55 |
| Uniform | 1.0 | 20.0 | 76.03 | 23.04 |
| Uniform | 1.0 | 30.0 | 76.03 | 37.01 |
| Uniform | 2.0 | 0.5 | 76.03 | 0.08 |
| Uniform | 2.0 | 1.0 | 76.03 | 0.08 |
| Uniform | 2.0 | 2.0 | 76.03 | 0.08 |
| Uniform | 2.0 | 4.0 | 76.03 | 0.10 |
| Uniform | 2.0 | 8.0 | 76.03 | 0.20 |
| Uniform | 2.0 | 12.0 | 76.03 | 0.63 |
| Uniform | 2.0 | 15.0 | 76.03 | 1.43 |
| Uniform | 2.0 | 20.0 | 76.03 | 4.51 |
| Uniform | 2.0 | 30.0 | 76.03 | 15.73 |
| Uniform | 4.0 | 0.5 | 76.03 | 0.01 |
| Uniform | 4.0 | 1.0 | 76.03 | 0.01 |
| Uniform | 4.0 | 2.0 | 76.03 | 0.01 |
| Uniform | 4.0 | 4.0 | 76.03 | 0.01 |
| Uniform | 4.0 | 8.0 | 76.03 | 0.02 |
| Uniform | 4.0 | 12.0 | 76.03 | 0.03 |

Table 8: BIM attacks on ResNet 50, mitigated with different types and strength ($\eta$) of noise.

| Noise Type | $\epsilon$ | $\eta$ | top-1 acc. (%) | ... (adversarial) |
|---|---|---|---|---|
| Uniform | 4.0 | 15.0 | 76.03 | 0.04 |
| Uniform | 4.0 | 20.0 | 76.03 | 0.09 |
| Uniform | 4.0 | 30.0 | 76.03 | 0.83 |
| Uniform | 8.0 | 0.5 | 76.03 | 0.00 |
| Uniform | 8.0 | 1.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 2.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 4.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 8.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 12.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 15.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 20.0 | 76.03 | 0.00 |
| Uniform | 8.0 | 30.0 | 76.03 | 0.01 |