
Detecting and Diagnosing Adversarial Images with Class-Conditional Capsule Reconstructions

Yao Qin*[†]
UC San Diego

Nicholas Frosst*
Google Brain

Sara Sabour
Google Brain

Colin Raffel
Google Brain

Garrison Cottrell
UC San Diego

Geoffrey Hinton
Google Brain

Abstract

Adversarial examples raise questions about whether neural network models are sensitive to the same visual features as humans. Most of the proposed methods for mitigating adversarial examples have subsequently been defeated by stronger attacks. Motivated by these issues, we take a different approach and propose to instead *detect* adversarial examples based on class-conditional reconstructions of the input. Our method uses the reconstruction network proposed as part of Capsule Networks (CapsNets), but is general enough to be applied to standard convolutional networks. We find that adversarial or otherwise corrupted images result in much larger reconstruction errors than normal inputs, prompting a simple detection method by thresholding the reconstruction error. Based on these findings, we propose the Reconstructive Attack which seeks both to cause a misclassification and a low reconstruction error. While this attack produces undetected adversarial examples, we find that for CapsNets the resulting perturbations can cause the images to appear visually more like the target class. This suggests that CapsNets utilize features that are more aligned with human perception and address the central issue raised by adversarial examples.

1 Introduction

Adversarial examples [Szegedy et al., 2014] are inputs that are designed by an adversary to cause a machine learning system to make a misclassification. A series of studies on adversarial attacks has shown that it is easy to cause misclassifications using visually imperceptible changes to an image under ℓ_p -norm based similarity metrics [Goodfellow et al., 2015, Kurakin et al., 2017, Madry et al., 2018, Carlini and Wagner, 2017b]. Since the discovery of adversarial examples, there has been a constant “arms race” between better attacks and better defenses. Many new defenses have been proposed [Song et al., 2018, Gong et al., 2017, Grosse et al., 2017, Metzen et al., 2017], only to be broken shortly thereafter [Carlini and Wagner, 2017a, Athalye et al., 2018]. Currently, the most effective approach to reduce network’s vulnerability to adversarial examples is “adversarial training”, in which a network is trained on both clean images and adversarially perturbed ones [Goodfellow et al., 2015, Madry et al., 2018]. However, adversarial training is very time-consuming because it requires generating adversarial examples during training. It also typically only helps improve a network’s robustness to adversarial examples that are generated in a similar way to those on which the network was trained. Hinton et al. [2018] showed that capsule models are more robust to simple adversarial attacks than CNNs but Michels et al. [2019] showed that this is not the case for all attacks.

*Equal Contributions.

[†]This work was done while the author interned at Google Brain.

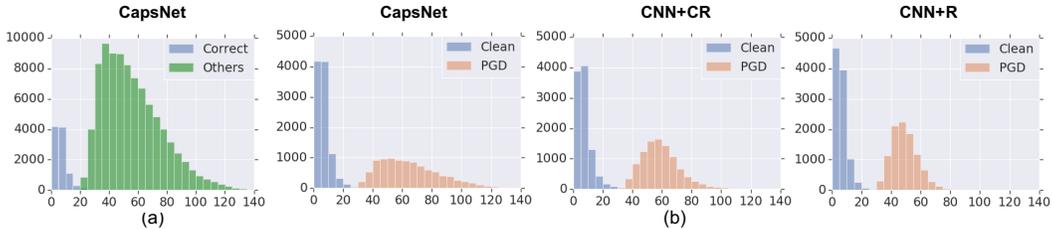


Figure 1: (a) The histogram of ℓ_2 distances between the input and the reconstruction using the winning capsule or other capsules in CapsNet on the real MNIST images. (b) The histograms of ℓ_2 distances between the reconstruction and the input for real and adversarial images for the three models explored in this paper on the MNIST dataset. We use PGD with the ℓ_∞ bound $\epsilon = 0.3$ to create the attacks. Notice the stark difference between the distributions of reconstructions of the winning capsule and the other capsules.

The cycle of attacks and defenses motivates us to rethink both how we can improve the general robustness of neural networks as well as the high-level motivation for this pursuit. One potential path forward is to detect adversarial inputs, instead of attempting to accurately classify them [Schott et al., 2018]. Recent work [Jetley et al., 2018, Gilmer et al., 2018b] argues that adversarial examples can exist within the data distribution, which implies that detecting adversarial examples based on an estimate of the data distribution alone might be insufficient. Instead, in this paper we develop methods for detecting adversarial examples by making use of *class-conditional reconstruction networks*. These sub-networks, first proposed by Sabour et al. [2017] as part of a Capsule Network (CapsNet), allow a model to produce a reconstruction of its input based on the identity and instantiation parameters of the winning capsule. Interestingly, we find that reconstructing an input from the capsule corresponding to the correct class results in a much lower reconstruction error than reconstructing the input from capsules corresponding to incorrect classes, as shown in Figure 1(a). Motivated by this, we propose using the reconstruction sub-network in a CapsNet as an attack-independent detection mechanism. Specifically, we reconstruct a given input from the resulting capsule pose parameters of the winning capsule and then detect adversarial examples by comparing the difference between the reconstruction distributions for natural and adversarial (or otherwise corrupted) images. The histograms of ℓ_2 reconstruction distance between the reconstruction and the input for real and adversarial images are shown in Figure 1(b).

We extend this detection mechanism to standard convolutional neural networks (CNNs) and show its effectiveness against black box and white box attacks on three image datasets; MNIST, Fashion-MNIST and SVHN. We show that capsule models achieve the strongest attack detection rates and accuracy on these attacks. We then test our method against a stronger attack, the Reconstructive Attack, specifically designed to attack our detection mechanism by generating adversarial examples with a small reconstruction error. With this attack we are able to create adversarial examples that change the model’s prediction and are not detected, but we show that this attack is less successful than a non-reconstructive attack and that CapsNets achieve the highest detection rates on these attacks as well.

We then explore two interesting observations specific to CapsNets. First, we illustrate how the success of the targeted reconstructive attack is highly dependent on the visual similarity between the source image and the target class. Second, we show that many of the resultant attacks resemble members of the target class and so cease to be “adversarial” – i.e., they may also be misclassified by humans. These findings suggest that CapsNets with class conditional reconstructions have the potential to address the real issue with adversarial examples – networks should make predictions based on the same properties of the image as people use rather than using features that can be manipulated by an imperceptible adversarial attack.

2 Background

Adversarial examples were first introduced in [Biggio et al., 2013, Szegedy et al., 2014], where a given image was modified by following the gradient of a classifier’s output with respect to the image’s pixels. Importantly, only an extremely small (and thus imperceptible) perturbation was required to cause a misclassification. Goodfellow et al. [2015] then developed the more efficient Fast Gradient

Sign method (FGSM), which can change the label of the input image X with a similarly imperceptible perturbation which is constructed by taking an ϵ step in the direction of the gradient. Later, the Basic Iterative Method (BIM) [Kurakin et al., 2017] and Project Gradient Descent [Madry et al., 2018] can generate stronger attacks improved on FGSM by taking multiple steps in the direction of the gradient and clipping the overall change to ϵ after each step. In addition, Carlini and Wagner [2017b] proposed another iterative optimization-based method to construct strong adversarial examples with small perturbations. An early approach to reducing vulnerability to adversarial examples was proposed by [Goodfellow et al., 2015], where a network was trained on both clean images and adversarially perturbed ones. Since then, there has been a constant “arms race” between better attacks and better defenses; Kurakin et al. [2018] provide an overview of this field.

A recent thread of research focuses on the generation of (and defense against) adversarial examples which are not simply slightly-perturbed versions of clean images. For example, several approaches were proposed which use generative models to create novel images which appear realistic but which result in a misclassification [Samangouei et al., 2018, Ilyas et al., 2017, Meng and Chen, 2017]. These adversarial images are not imperceptibly close to some existing image, but nevertheless resemble members of the data distribution to humans and are strongly misclassified by neural networks. [Sabour et al., 2016] also consider adversarial examples which are not the result of pixel-space perturbations by manipulating the hidden representation of a neural network in order to generate an adversarial example.

Another line of work, surveyed by [Carlini and Wagner, 2017a], attempts to circumvent adversarial examples by detecting them with a separately-trained classifier [Gong et al., 2017, Grosse et al., 2017, Metzen et al., 2017] or using statistical properties [Hendrycks and Gimpel, 2017, Li and Li, 2017, Feinman et al., 2017, Grosse et al., 2017]. However, many of these approaches were subsequently shown to be flawed [Carlini and Wagner, 2017a, Athalye et al., 2018]. Most recently, [Schott et al., 2018] investigated the effectiveness of a class-conditional generative model as a defense mechanism for MNIST digits. In comparison, our method does not increase the computational overhead of the classification and tries to detect adversarial examples by attempting to reconstruct them.

3 Preliminaries

Adversarial Examples Given a clean test image x , its corresponding label y , and a classifier $f(\cdot)$ which predicts a class label given an input, we refer to $x' = x + \delta$ as an adversarial example if it is able to fool the classifier into making a wrong prediction $f(x') \neq f(x) = y$. The small adversarial perturbation δ (where “small” is measured under some norm) causes the adversarial example x' to appear visually similar to the clean image x but to be classified differently. In the unrestricted case where we only require that $f(x') \neq y$, we refer to x' as an “untargeted adversarial example”. A more powerful attack is to generate a “targeted adversarial example”: instead of simply fooling the classifier to make a wrong prediction, we force the classifier to predict some targeted label $f(x') = t \neq y$. In this paper, the target label t is selected uniformly at random as any label which is not the ground-truth correct label. As is standard practice in the literature, in this paper we use ℓ_∞ and ℓ_2 as norms to measure the size of δ . When generating an adversarial example, we use these norms to constrain the constructed adversarial examples.

We test our detection mechanism on three ℓ_∞ norm based attacks (fast gradient sign method (FGSM) [Goodfellow et al., 2015], the basic iterative method (BIM) [Kurakin et al., 2017], projected gradient descent (PGD) [Madry et al., 2018]) and one ℓ_2 norm based attack (Carlini-Wagner (CW) [Carlini and Wagner, 2017b]).

Capsule Networks Capsule Networks (CapsNets) are an alternative architecture for neural networks [Sabour et al., 2017, Hinton et al., 2018]. In this work we make use of the CapsNet architecture detailed by [Sabour et al., 2017]. Unlike a standard neural network which is made up of layers of scalar-valued units, CapsNets are made up of layers of capsules, which output a vector or matrix. Intuitively, just as one can think of the activation of a unit in a normal neural network as the presence of a feature in the input, the activation of a capsule can be thought of as both the presence of a feature and the pose parameters that represent attributes of that feature. A top-level capsule in a classification network therefore outputs both a classification and pose parameters that represent the instance of that class in the input. This high level representation allows us to train a reconstruction network.

Threat Model Prior work considers adversarial examples mainly in two categories: white-box and black-box. In this paper, we test our detection mechanism against both of these attacks. For white-box attacks, the adversary has full access to the model as well as its parameters. In particular, the adversary is allowed to compute the gradient through the model to generate adversarial examples. To perform black-box attacks, the adversary is allowed to know the network architecture but not parameters. Therefore, we retrain a substitute model that has the same architecture as the target model and generate adversarial examples by attacking the substitute model. Then we transfer these attacks to the target model. For ℓ_∞ based attacks, we always control the ℓ_∞ norm of the adversarial examples to be within a relatively small bound ϵ , which is specific to each dataset.

4 Detecting Adversarial Images by Reconstruction

To detect adversarial images, we make use of the reconstruction network proposed in [Sabour et al., 2017], which takes pose parameters v as input and outputs the reconstructed image $g(v)$. The reconstruction network is simply a fully connected neural network with two ReLU hidden layers with 512 and 1024 units respectively, with a sigmoid output with the same dimensionality as the dataset. The reconstruction network is trained to minimize the ℓ_2 distance between the input image and the reconstructed image. This same network architecture is used for all the models and datasets we explore. The only difference is what is given to the reconstruction network as input.

4.1 Models

CapsNet The reconstruction network of the CapsNet is *class-conditional*: It takes in the pose parameters of all the class capsules and masks all values to 0 except for the pose parameters of the predicted class. We use this reconstruction network for detecting adversarial attacks by measuring the Euclidean distance between the input and a class conditional reconstruction. Specifically, for any given input x , the CapsNet can output a prediction $f(x)$ as well as the the pose parameters v for all classes. The reconstruction network takes in the pose parameters and then selects the pose parameter corresponding to the predicted class, denoted as $v_{f(x)}$, to generate a reconstruction $g(v_{f(x)})$. Then we compute the ℓ_2 reconstruction distance $d(x) = \|g(v_{f(x)}), x\|_2$ between the reconstructed image and the input image, and compare it with a pre-defined detection threshold p (described below in Section 4.2). If the reconstruction distance $d(x)$ is higher than the detection threshold p , we flag the input as an adversarial example. Figure 1 (b) shows example histograms of reconstruction distances for natural images and typical adversarial examples.

CNN+CR Although our strategy is inspired by the reconstruction networks used in CapsNets, the strategy can be extended to standard convolutional neural networks (CNNs). We create a similar architecture, CNN with conditional reconstruction (CNN+CR), by dividing the penultimate hidden layer of a CNN into groups corresponding to each class. The sum of each neuron group serves as the logit for that particular class and the group itself serves the same purpose as the pose parameters in the CapsNet. We use the same masking mechanism as Sabour et al. [2017] to select the pose parameter corresponding to the predicted label $v_{f(x)}$ and generate the reconstruction based on the selected pose parameters. In this way we extend the class-conditional reconstruction network to standard CNNs.

CNN+R We can also create a more naïve implementation of our strategy by simply computing the reconstruction from the activations in the entire penultimate layer without any masking mechanism. We call this model the "CNN+R" model. In this way we are able to study the effect of conditioning on the predicted class.

4.2 Detection Threshold

We find the threshold p for detecting adversarial inputs by measuring the reconstruction error between a validation input image and its reconstruction. If the distance between the input and the reconstruction is above the chosen threshold p , we classify the data as adversarial. Choosing the detection threshold p involves a trade-off between false positive and false negative detection rates. The optimal threshold depends on the probability of the system being attacked. Such a trade-off is discussed by Gilmer et al. [2018a]. In our experiments we don't tune this parameter and simply set it as the 95th percentile of validation distances. This means our false positive rate on the validation data is 5%.

Networks	Targeted (%)				Untargeted (%)			
	FGSM	BIM	PGD	CW	FGSM	BIM	PGD	CW
CapsNet	3/0	82/0	86/0	99/2	11/0	99/0	99/0	100/19
CNN+CR	16/0	93/0	95/0	89/8	85/0	100/0	100/0	100/28
CNN+R	37/0	100/0	100/0	100/47	64/0	100/0	100/0	100/63

Table 1: **Success Rate / Undetected Rate** of white-box targeted and untargeted attacks on the MNIST dataset. In the table, S_t/R_t is shown for targeted attacks and S_u/R_u is presented for untargeted attacks. Full results for FashionMNIST and SVHN can be seen in the Table 4 in the Appendix.

4.3 Evaluation Metrics

We can measure the ability of different networks to detect adversarial examples by computing the proportion of adversarial examples that both successfully fool the network and go undetected. For targeted attacks, the success rate S_t is defined as the proportion of inputs which are classified as the target class, $S_t = \frac{1}{N} \sum_i^N (f(x'_i) = t_i)$, while the success rate for untargeted attacks is defined as the proportion of inputs which are misclassified, $S_u = \frac{1}{N} \sum_i^N (f(x'_i) \neq y_i)$. The undetected rate for targeted attacks R_t is defined as the proportion of attacks that are successful and undetected, $R_t = \frac{1}{N} \sum_i^N (f(x'_i) = t_i) \cap (d(x'_i) \leq p)$, where $d(\cdot)$ computes the reconstruction distance of the input and p denotes the detection threshold introduced in section 4.2. Similarly, the undetected rate for untargeted attacks R_u can be defined as $R_u = \frac{1}{N} \sum_i^N (f(x'_i) \neq y_i) \cap (d(x'_i) \leq p)$. The smaller the undetected rate R_t or R_u is, the stronger the model is in detecting adversarial examples.

4.4 Test Models and Datasets

In all experiments, all three models (CapsNet, CNN+R, and CNN+CR) have the same number of parameters and were trained with Adam [Kingma and Ba, 2015] for the same number of epochs. In general, all models achieved similar test accuracy. We did not do an exhaustive hyperparameter search on these models, instead we chose hyper-parameters that allowed each model to perform roughly equivalently on the test sets. We ran experiments on three datasets: MNIST [LeCun et al., 1998], FashionMNIST [Xiao et al., 2017], and SVHN [Netzer et al., 2011]. The test error rate for each model on these three datasets, as well as details of the model architectures, can be seen in Section A and Section B in the Appendix.

5 Experiments

We first demonstrate how reconstruction networks can detect standard white and black-box attacks in addition to naturally corrupted images. Then, we introduce the “reconstructive attack”, which is specifically designed to circumvent our defense and show that it is a more powerful attack in this setting. Based on this finding, we qualitatively study the kind of misclassifications caused by the reconstructive attack and argue that they suggest that CapsNets learn features that are better aligned with human perception.

5.1 Standard Attacks

White Box We present the success and undetected rates for several targeted and untargeted attacks on MNIST (Table 1), FashionMNIST, and SVHN (Table 4 presented in the Appendix). Our method is able to accurately detect many attacks with very low undetected rates. Capsule models almost always have the lowest undetected rates out of our three models. It is worth noting that this method performs best with the simplest dataset, MNIST, and that the highest undetected rates are found with the Carlini-Wagner attack on the SVHN dataset. This illustrates both the strength of this attack and a shortcoming of our defense, namely that our detection mechanism relies on ℓ_2 image distance as a proxy for visual similarity, and in the case of higher dimensional color datasets such as SVHN, this proxy is less meaningful.

Black Box We also tested our detection mechanism results on black box attacks. Given the low undetected rates in the white-box settings, it is not surprising that our detection method is able to detect black box attacks as well. In fact, on the MNIST dataset the capsule model is able to detect all targeted and untargeted PGD attacks. Both the CNN-R and the CNN-CR models are able to detect the black box attacks as well, but with a relatively higher success rate. A table of these results can be seen in the Table 7 in the Appendix.

5.2 Corruption Attack

Recent work has argued that improving the robustness of neural networks to epsilon-bounded adversarial attacks should not come at the expense of increasing error rates under distributional shifts that do not affect human classification rates and are likely to be encountered in the “real-world” [Gilmer et al., 2018a]. For example, if an image is corrupted due to adverse weather, lighting, or occlusion, we might hope that our model can continue to provide reliable predictions or detect the distributional shift. We can test our detection method on its ability to detect these distributional shifts by making use of the Corrupted MNIST dataset [Mu and Gilmer, 2019]. This data set contains many visual transformations of MNIST that do not seem to affect human performance, but nevertheless are strongly misclassified by state-of-the-art MNIST models. Our three models can almost always detect these distributional shifts (in all corruptions CapsNets have either a small undetected rate or an undetected rate of 0). Please refer to Table 5 in the Appendix for detailed test results and Figure 5 and Figure 6 for visualization of Corrupted MNIST.

5.3 Reconstructive Attacks

Thus far we have only evaluated previously-defined attacks. To evaluate our defense more rigorously, we introduce an attack specifically designed to take into account our defense mechanism. In order to construct adversarial examples that cannot be detected by the network, we propose a two-stage optimization method we call a “reconstructive attack”. Specifically, in each step, we first attempt to fool the network by following a standard attack which computes the gradient of the cross-entropy loss function with respect to the input. Then, in the second stage, we take the reconstruction error into account by updating the adversarial perturbation based on the ℓ_2 reconstruction loss. In this way, we endeavor to construct adversarial examples that can fool the network and also have a small ℓ_2 reconstruction error. The untargeted and targeted reconstructive attacks based are described in more detail below.

Untargeted Reconstructive Attacks To construct untargeted reconstructive attacks, we first update the perturbation based on the gradient of the cross-entropy loss function following a standard FGSM attack [Goodfellow et al., 2015], that is:

$$\delta \leftarrow \text{clip}_\epsilon(\delta + c \cdot \beta \cdot \text{sign}(\nabla_\delta \ell_{net}(f(x + \delta), y))), \quad (1)$$

where $\ell_{net}(f(\cdot), y)$ is the cross-entropy loss function, ϵ is the ℓ_∞ bound for our attacks, c is a hyperparameter controlling the step size in each iteration and β is a hyperparameter which balances the importance of the cross-entropy loss and the reconstruction loss (explained further below). In the second stage, we focus on constraining the reconstructed image from the newly predicted label to have a small reconstruction distance by updating δ according to

$$\delta \leftarrow \text{clip}_\epsilon(\delta - c \cdot (1 - \beta) \cdot \text{sign}(\nabla_\delta (\|g(v_{f(x+\delta)}) - (x + \delta)\|_2))), \quad (2)$$

where $g(v_{f(x+\delta)})$ is the class-conditional reconstruction based on the predicted label $f(x + \delta)$ in a CapsNet or CNN+CR network. The δ used here is the optimized δ from the first stage. $\|g(v_{f(x+\delta)}) - (x + \delta)\|_2$ is the ℓ_2 reconstruction distance between the reconstructed image and the input image. Since the CNN+R network does not use the class conditional reconstruction, we simply use the reconstructed image without the masking mechanism. According to Eqn 1 and Eqn 2, we can see that β balances the importance between the success rate of attacks and the reconstruction distance. This hyperparameter was tuned for each model and each dataset in order to create the strongest attacks with the worst undetected rate (**worst case** undetected rate). The plots showing that success rate and undetected rate changes as this parameter varied between 0 and 1 can be seen in the Figure 7 of the Appendix.

Targeted Reconstructive Attacks We perform a similar two-stage optimization to construct targeted reconstructive attacks, by defining a target label and attempting to maximize the classification probability of this label, and minimize the reconstruction error from the corresponding capsule. Because the targeted label is given, another way to construct targeted reconstructive attacks is to combine these two stages into one stage via minimizing the loss function $\ell = \beta \cdot \ell_{net}(f(x + \delta), y) + (1 - \beta) \cdot \|g(v_{f(x+\delta)}) - (x + \delta)\|_2$. We implemented both of these targeted reconstructive attacks and found that the two-stage version is a stronger attack. Therefore, all the Reconstructive Attack experiments performed in this paper are based on two-stage optimization.

	MNIST		FASHION		SVHN	
	Targeted R-PGD	Untargeted R-PGD	Targeted R-PGD	Untargeted R-PGD	Targeted R-PGD	Untargeted R-PGD
CapsNet	50.7/33.7	88.1/37.9	53.7/29.8	84.9/75.5	82.0/79.2	98.9/97.5
CNN+CR	98.6/68.1	99.4/87.7	89.8/84.4	91.5/86.0	99.0/97.9	99.9/99.5
CNN+R	95.5/71.2	95.1/70.5	94.6/88.4	98.9/90.0	99.5/99.3	100.0/99.9

Table 2: Success rate and the **worst case** undetected rate of white-box targeted and untargeted reconstructive attacks. S_t/R_t is shown for targeted attacks and S_u/R_u is presented for untargeted attacks. The worst case undetected rate is reported via tuning the hyperparameter β in Eqn 1 and Eqn 2. The best models are shown in **bold**. All the numbers are shown in %. A full table with more attacks can be seen in the Table 6 in the Appendix.

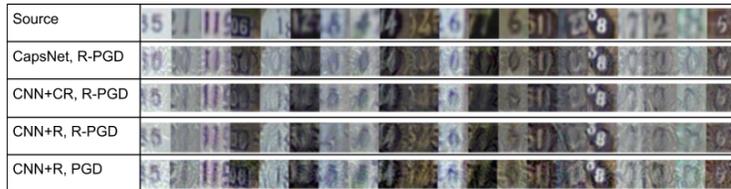


Figure 2: These are randomly sampled (not cherry picked) successful and undetected adversarial attacks created by R-PGD with a target class of 0 using SVHN for each model. They were created with an epsilon bound of 25/255 based on the pixel range, as is standard in the adversarial literature. We can see that for the capsule model, many of the attacks are not “adversarial” as they resemble members of the target class.

We build our reconstructive attack based on the standard PGD attack, denoted as R-PGD, and test the performance of our detection models against this reconstructive attack in a white-box setting (white-box Reconstructive FGSM and BIM are reported in Table 6 in the Appendix). Comparing Table 1 and Table 2, we can see that the Reconstructive Attack is significantly less successful at changing the model’s prediction (lower success rates than the standard attack). However, this attack is more successful at fooling our detection method. For all reconstructive attacks on the three datasets, the capsule model has the lowest attack success rate and the lowest undetected rate. We report results for black-box R-PGD attacks in Table 7 in the Appendix, which suggest similar conclusions. Note that if our detection method was perfect, then the undetected rate would be 0 as is the case for the standard (non-reconstructive) MNIST attacks. This shortcoming warrants further investigation (see below).

5.4 Visual Coherence of the Reconstructive Attack

Thus far we have treated all attacks as equal. However, this is not the case – a key component of an adversarial example is that it is visually similar to the source image, and that it does not resemble the adversarial target class. The adversarial research community makes use of a small epsilon bound as a mechanism for ensuring that the resultant adversarial attacks are visually unchanged from the source image. In a standard neural network, and with standard attacks this heuristic is sufficient, because taking gradient steps in the image space in order to have a network misclassify an image normally results in something visually similar to the source image. However, this is not the case for adversarial attacks which take the reconstruction error into account: As shown in Figure 2, when we use R-PGD to attack the CapsNet, many of the resultant attacks resemble members of the target class. In this way, they stop being “adversarial”. As such, an attack detection method which does not detect them as adversarial is arguably behaving correctly. This puts the previously undetected rates presented earlier in a new light, and illustrates a difficulty in the evaluation of adversarial attacks and defenses. In other words, the research community has assumed that small epsilon bounded implies visual similarity to the source image, but this is not always a valid assumption.

It is worth noting however that not all the adversarial examples created with R-PGD have the property of visual similarity to the target class, i.e. some of them are typically adversarial. Thus far we have also treated all misclassifications as equal, but this too is a simplification. If our true aim in

adversarial robustness research is to create models that make predictions based on reasonable and human-observable features, then we would prefer models that are more likely to misclassify a “shirt” as a “t-shirt” (in the case of FashionMNIST) than to misclassify a “bag” as a “sweater”. For a model to behave ideally, the success of an adversarial perturbation would be related to the visual similarity between the source and the target class. By visualizing a matrix of adversarial success rates between each pair of classes (shown in Figure 3), we can see that for the capsule model there is great variance between source and target class pairs and that the success rate is related to the visual similarity of the classes. This is not the case for either of the other two CNN models (additional confusion matrices presented in the Appendix as Figure 8).

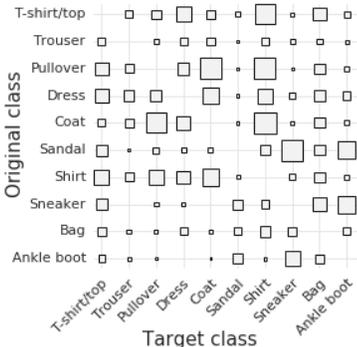


Figure 3: Success rate of targeted reconstructive PGD attacks for FashionMNIST with $\ell_\infty = 25/255$ for CapsNet. The size of the box at position x, y represents the success rate of adversarially perturbing inputs of class x to be classified at class y . For example, most attempts to cause an “Ankle boot” image to be misclassified were unsuccessful, with successful attacks only occurring when the target class was “Sandal” or “Sneaker”. On the other hand, more successful attacks occurred when trying to cause “Shirt” images to be misclassified as “T-shirt/top”, “Pullover”, “Dress”, or “Coat”.

5.5 Discussion

Our detection mechanism relies on a similarity metric (i.e. a measure of reconstruction error) between the reconstruction and the input. This metric is required both during training in order to train the reconstruction network and during test time in order to flag adversarial examples. In the three data sets we have evaluated, the distance between examples roughly correlates with semantic similarity. This is not the case, however, for images in more complex dataset such as CIFAR-10 [Krizhevsky, 2009] or ImageNet [Deng et al., 2009], in which two images may be similar in terms of semantic content but nevertheless have significant ℓ_2 distance. This issue will need to be resolved for this method to scale up to more complex problems, and offers a promising avenue for future research. Furthermore, our reconstruction network is trained on a hidden representation of one class but is trained to reconstruct the entire input. In datasets without distractors or backgrounds, this is not a problem. But in the case of ImageNet, in which the object responsible for the classification is not the only object in the image, attempting to reconstruct the entire input from a class encoding seems misguided.

6 Conclusion

We have presented a simple architectural extension that enables adversarial attack detection. Notably, this method does not rely on a specific predefined adversarial attack. We have shown that by reconstructing the input from the internal class-conditioned representation, our system is able to accurately detect black box and white box FGSM, BIM, PGD, and CW attacks. Of the three models we explored, we showed that the CapsNet performed best at this task, and was able to detect adversarial examples with greater accuracy on all the data sets we explored. We then proposed a new attack to beat our defense - the Reconstructive Attack - in which the adversary optimizes not only the classification loss but also minimizes the reconstruction loss. We showed that this attack was less successful than a standard attack in changing the model’s prediction, though it was able to fool our detection mechanism to some extent. We showed qualitatively that for the CapsNet, the success of this attack was proportional to the visual similarity between the target class and the source class. Finally we showed that many images generated by this attack when generated with the CapsNet are not typically adversarial, i.e. many of the resultant attacks resemble members of the target class even with a small epsilon bound. This is not the case for standard (non-reconstructive) attacks or for the CNN model. This implies that the gradient of the reconstructive attack may be better aligned with the true data manifold, and implies that the capsule model may be relying on human-detectable visually coherent features to make predictions. We believe this is a step towards solving the true problem posed by adversarial examples.

References

- Anish Athalye, Nicholas Carlini, and David Wagner. Obfuscated gradients give a false sense of security: Circumventing defenses to adversarial examples. *International Conference on Machine Learning*, 2018.
- Battista Biggio, Igino Corona, Davide Maiorca, Blaine Nelson, Nedim Šrđić, Pavel Laskov, Giorgio Giacinto, and Fabio Roli. Evasion attacks against machine learning at test time. In *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*, 2013.
- Nicholas Carlini and David Wagner. Adversarial examples are not easily detected: Bypassing ten detection methods. In *ACM Workshop on Artificial Intelligence and Security*, 2017a.
- Nicholas Carlini and David Wagner. Towards evaluating the robustness of neural networks. In *2017 IEEE Symposium on Security and Privacy (SP)*, 2017b.
- Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2009.
- Reuben Feinman, Ryan R. Curtin, Saurabh Shintre, and Andrew B. Gardner. Detecting adversarial samples from artifacts. *arXiv preprint arXiv:1703.00410*, 2017.
- Justin Gilmer, Ryan P. Adams, Ian Goodfellow, David Andersen, and George E. Dahl. Motivating the rules of the game for adversarial example research. *arXiv preprint arXiv:1807.06732*, 2018a.
- Justin Gilmer, Luke Metz, Fartash Faghri, Samuel S. Schoenholz, Maithra Raghu, Martin Wattenberg, and Ian Goodfellow. Adversarial spheres. *arXiv preprint arXiv:1801.02774*, 2018b.
- Zhitao Gong, Wenlu Wang, and Wei-Shinn Ku. Adversarial and clean data are not twins. *arXiv preprint arXiv:1704.04960*, 2017.
- Ian J. Goodfellow, Jonathon Shlens, and Christian Szegedy. Explaining and harnessing adversarial examples. *International Conference on Learning Representations*, 2015.
- Kathrin Grosse, Praveen Manoharan, Nicolas Papernot, Michael Backes, and Patrick McDaniel. On the (statistical) detection of adversarial examples. *arXiv preprint arXiv:1702.06280*, 2017.
- Dan Hendrycks and Kevin Gimpel. Early methods for detecting adversarial images. *International Conference on Learning Representations*, 2017.
- Geoffrey E. Hinton, Sara Sabour, and Nicholas Frosst. Matrix capsules with EM routing. *International Conference on Learning Representations*, 2018.
- Andrew Ilyas, Ajil Jalal, Eirini Asteri, Constantinos Daskalakis, and Alexandros G. Dimakis. The robust manifold defense: Adversarial training using generative models. *arXiv preprint arXiv:1712.09196*, 2017.
- Saumya Jetley, Nicholas Lord, and Philip Torr. With friends like these, who needs adversaries? In *Advances in Neural Information Processing Systems*, 2018.
- Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *International Conference on Learning Representations*, 2015.
- Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, University of Toronto, 2009.
- Alexey Kurakin, Ian Goodfellow, and Samy Bengio. Adversarial examples in the physical world. *International Conference on Learning Representations*, 2017.
- Alexey Kurakin, Ian Goodfellow, Samy Bengio, Yinpeng Dong, Fangzhou Liao, Ming Liang, Tianyu Pang, Jun Zhu, Xiaolin Hu, Cihang Xie, et al. Adversarial attacks and defences competition. *arXiv preprint arXiv:1804.00097*, 2018.

- Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
- Xin Li and Fuxin Li. Adversarial examples detection in deep networks with convolutional filter statistics. In *Proceedings of the IEEE International Conference on Computer Vision*, 2017.
- Aleksander Madry, Aleksandar Makelov, Ludwig Schmidt, Dimitris Tsipras, and Adrian Vladu. Towards deep learning models resistant to adversarial attacks. *International Conference on Learning Representations*, 2018.
- Dongyu Meng and Hao Chen. Magnet: a two-pronged defense against adversarial examples. In *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*, 2017.
- Jan Hendrik Metzen, Tim Genewein, Volker Fischer, and Bastian Bischoff. On detecting adversarial perturbations. *International Conference on Learning Representations*, 2017.
- Felix Michels, Tobias Uelwer, Eric Upschulte, and Stefan Harmeling. On the vulnerability of capsule networks to adversarial attacks. *arXiv preprint arXiv:1906.03612*, 2019.
- Norman Mu and Justin Gilmer. Mnist-c: A robustness benchmark for computer vision. In *ICML 2019 Workshop on Uncertainty and Robustness in Deep Learning*, 2019.
- Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. Reading digits in natural images with unsupervised feature learning. In *NIPS workshop on deep learning and unsupervised feature learning*, 2011.
- Sara Sabour, Yanshuai Cao, Fartash Faghri, and David J. Fleet. Adversarial manipulation of deep representations. *International Conference on Learning Representations*, 2016.
- Sara Sabour, Nicholas Frosst, and Geoffrey E. Hinton. Dynamic routing between capsules. In *Advances in Neural Information Processing Systems*, 2017.
- Pouya Samangouei, Maya Kabkab, and Rama Chellappa. Defense-gan: Protecting classifiers against adversarial attacks using generative models. *International Conference on Learning Representations*, 2018.
- Lukas Schott, Jonas Rauber, Wieland Brendel, and Matthias Bethge. Robust perception through analysis by synthesis. *arXiv preprint arXiv:1805.09190*, 2018.
- Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, and Nate Kushman. Pixeldefend: Leveraging generative models to understand and defend against adversarial examples. *International Conference on Learning Representations*, 2018.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *International Conference on Learning Representations*, 2014.
- Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017.

Appendix

A Network Architectures

Figure 4 shows the architecture of the capsule network and the CNN based reconstruction models used for experiments on MNIST, Fashion MNIST and SVHN dataset. MNIST and Fashion MNIST have exactly the same architectures while for we use larger models for SVHN. Note that the only difference between the CNN reconstruction (CNN+R) and the CNN conditional reconstruction (CNN+CR) is the masking procedure on the input to the reconstruction network based on the predicted class. In CNN+CR, only the hidden for the winning category are used to drive the reconstruction network. All three models have the same number of parameters.

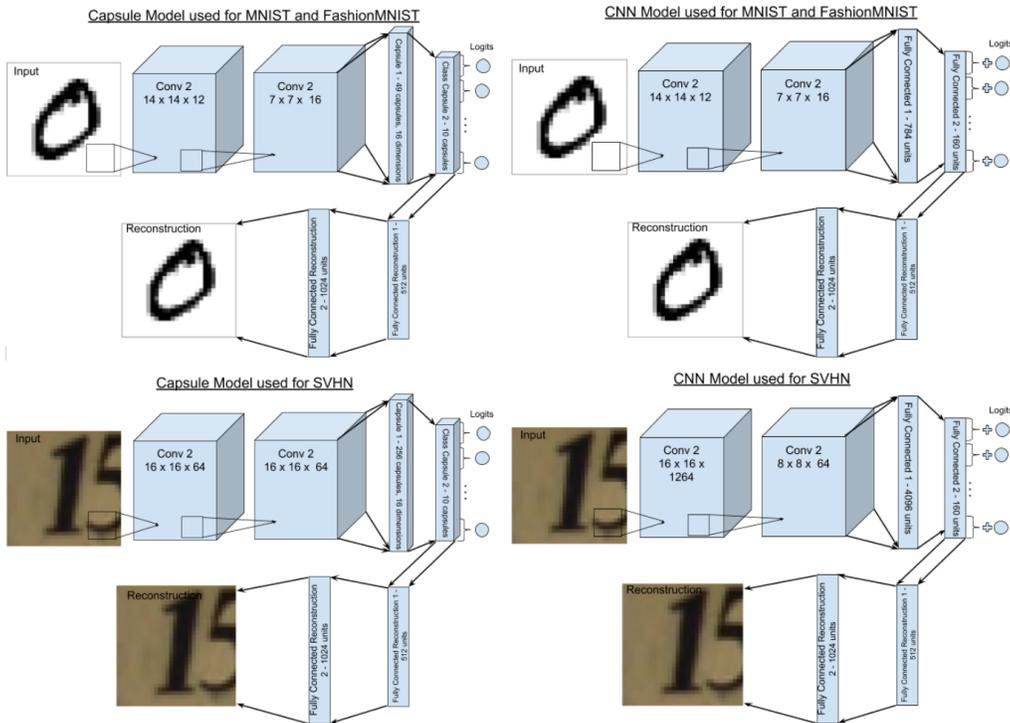


Figure 4: The architecture for the CapsNet and CNN based models used for our experiments on MNIST [LeCun et al., 1998], FashionMNIST [Xiao et al., 2017], and SVHN [Netzer et al., 2011].

B Test models

The error rate of each test model used in the paper are presented in Table 3, where the error rate measures the ratio of misclassified images, defined as $E = \frac{1}{N} \sum_i (f(x_i) \neq y_i)$. We ensure that they have similar performance.

Dataset	CapsNet	CNN+CR	CNN+R
MNIST	0.6%	0.7%	0.6%
FashionMNIST	9.6%	9.5%	9.3%
SVHN	10.7%	9.3%	9.5%

Table 3: Error rate of each model when the input are clean test images in each dataset.

C Implementation Details

For all the ℓ_∞ based adversarial examples, the ℓ_∞ norm of the perturbations is bound by ϵ , which is set to 0.3, 0.1, 0.1 for MNIST, Fashion MNIST and SVHN dataset, respectively, following previous work [Madry et al., 2018, Song et al., 2018]. In FGSM based attacks, the step size c is 0.05. In BIM-based [Kurakin et al., 2017] and PGD-based [Madry et al., 2018] attacks, the step size c is 0.01 for all the datasets and the number of iterations are 1000, 500 and 200 for MNIST, Fashion MNIST and SVHN dataset, respectively. We choose a sufficiently large number of iterations to ensure the attacks have converged.

We use the publicly released code from [Carlini and Wagner, 2017b] to perform the CW attack for our models. The number of iterations are set to 1000 for all three datasets.

D White Box Standard Attacks

The results of four white box standard attacks on the three datasets are shown in Table 4.

Networks	Targeted (%)				Untargeted (%)			
	FGSM	BIM	PGD	CW	FGSM	BIM	PGD	CW
MNIST Dataset								
CapsNet	3/0	82/0	86/0	99/2	11/0	99/0	99/0	100/19
CNN+CR	16/0	93/0	95/0	89/8	85/0	100/0	100/0	100/28
CNN+R	37/0	100/0	100/0	100/47	64/0	100/0	100/0	100/63
FASHION MNIST Dataset								
CapsNet	7/5	54/9	55/10	100/26	35/29	86/50	87/51	100/68
CNN+CR	19/13	89/28	89/28	87/37	74/33	100/25	100/24	100/72
CNN+R	23/16	98/19	98/19	99/81	62/48	100/35	100/34	100/87
SVHN Dataset								
CapsNet	22/20	83/45	84/46	100/90	74/67	99/70	99/68	100/94
CNN+CR	24/23	99/90	99/90	99/93	87/82	100/90	100/89	100/90
CNN+R	26/24	100/86	100/86	100/94	88/82	100/92	100/92	100/95

Table 4: Success rate and undetected rate of white-box targeted and untargeted attacks. In the table, S_t/R_t is shown for targeted attacks and S_u/R_u is presented for untargeted attacks.

E Corruption Attack

The error rate and undetected rate of the three test models on the Corrupted MNIST dataset is shown in Table 5, where the error rate measures the ratio of misclassified images.

Corruption	Clean	Gaussian Noise	Gaussian Blur	Line	Dotted Line	Elastic Transform
CapsNet	0.6/0.2	12.1/0.0	10.3/4.1	19.6/0.1	4.3/0.0	11.3/0.8
CNN+CR	0.7/0.3	9.8/0.0	6.7/4.2	17.6/0.1	4.2/0.0	11.1/1.1
CNN+R	0.6/0.4	6.7/0.0	8.9/6.4	18.9/0.1	3.1/0.0	12.2/2.1
Corruption	Saturate	JPEG	Quantize	Sheer	Spatter	Rotate
CapsNet	3.5/0.0	0.8/0.4	0.7/0.1	1.6/0.4	1.9/0.2	6.5/2.2
CNN+CR	1.5/0.0	0.8/0.5	0.9/0.1	2.1/0.4	1.8/0.4	6.1/1.6
CNN+R	1.2/0.0	0.7/0.5	0.7/0.2	2.2/0.7	1.8/0.4	6.5/3.4
Corruption	Contrast	Inverse	Canny Edge	Fog	Frost	Zigzag
CapsNet	92.0/0.0	91.0/0.0	21.5/0.0	83.7/0.0	70.6/0.0	16.9/0.0
CNN+CR	72.0/32.6	78.1/0.0	34.6/0.0	66.0/0.5	37.6/0.0	18.4/0.0
CNN+R	73.4/49.4	88.1/0.0	23.4/0.0	65.6/0.1	36.2/0.0	17.5/0.0

Table 5: Error rate/undetected rate on the Corrupted MNIST dataset.

E.1 Visualization of Corrupted MNIST Dataset

Visualization of examples from Corrupted MNIST dataset [Mu and Gilmer, 2019] and the corresponding reconstructed images for each model are shown in Figure 5 and Figure 6.

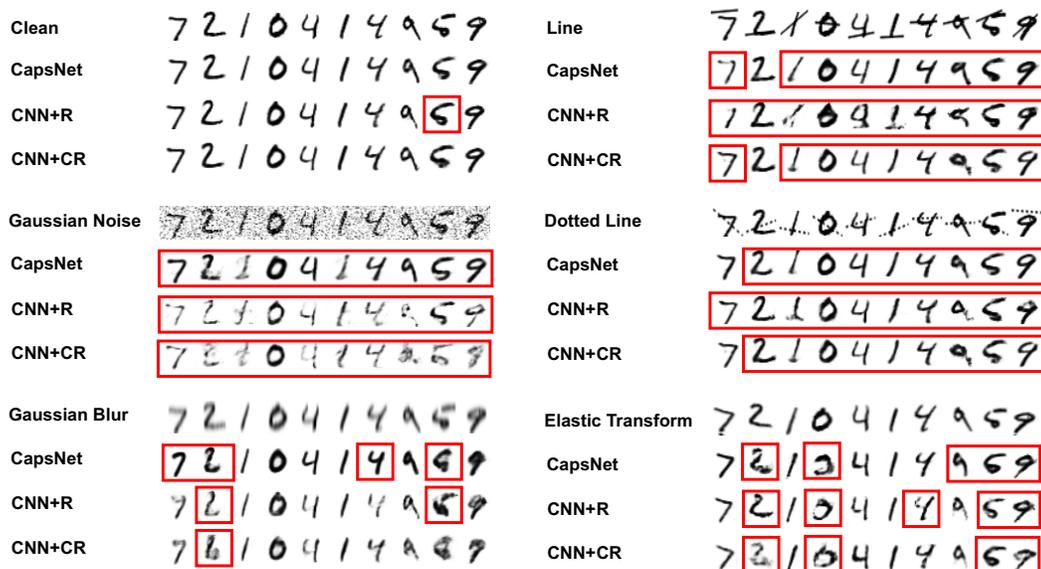


Figure 5: Examples of Corrupted MNIST and the reconstructed image for each model. Red boxes indicate that this input is flagged as an adversarial example while no boxes indicate that this input is correctly classified and is not flagged as an adversarial example.



Figure 6: Examples of Corrupted MNIST and the reconstructed image for each model. Red boxes indicate that this input is flagged as an adversarial example while green boxes indicate that this input has been misclassified and not been detected. No boxes indicate that this input is correctly classified and is not flagged as an adversarial example.

F Reconstructive Attacks

The results of Reconstructive FGSM, BIM and PGD on the three datasets are reported in Table 6.

Networks	Targeted (%)			Untargeted (%)		
	R-FGSM	R-BIM	R-PGD	R-FGSM	R-BIM	R-PGD
MNIST Dataset						
CapsNet	1.8/0.3	51.0/33.8	50.7/33.7	6.1/1.0	84.5/35.1	88.1/37.9
CNN+CR	7.6/0.5	98.0/68.1	98.6/68.1	41.7/3.2	96.5/86.8	99.4/87.7
CNN+R	16.9/3.3	86.3/65.9	95.5/71.2	25.9/8.1	82.9/67.8	95.1/70.5
FASHION MNIST Dataset						
CapsNet	6.5/5.8	53.3/28.4	53.7/29.8	33.3/29.9	85.3/75.9	84.9/75.5
CNN+CR	17.7/14.0	80.3/72.4	78.1/72.0	68.0/57.3	89.8/84.4	91.5/86.0
CNN+R	19.4/17.6	95.2/88.8	94.6/88.4	58.6/53.5	98.8/90.1	98.9/90.0
SVHN Dataset						
CapsNet	21.6/21.2	81.1/78.3	82.0/79.2	71.6/68.3	98.9/97.5	98.9/97.5
CNN+CR	24.2/22.6	98.5/97.6	99.0/97.9	86.0/82.3	99.9/99.5	99.9/99.5
CNN+R	26.6/25.8	99.6/99.4	99.5/99.3	87.1/84.5	100.0/99.9	100.0/99.9

Table 6: Success rate and the **worst case** undetected rate of white-box targeted and untargeted reconstructive attacks. S_t/R_t is shown for targeted attacks and S_u/R_u is presented for untargeted attacks.

F.1 Hyperparameter β

Figure 7 shows the plot of the success rate and the undetected rate versus the hyperparameter β which balances the importance between success rate and undetected rate in the targeted reconstructive PGD attacks on the MNIST dataset.

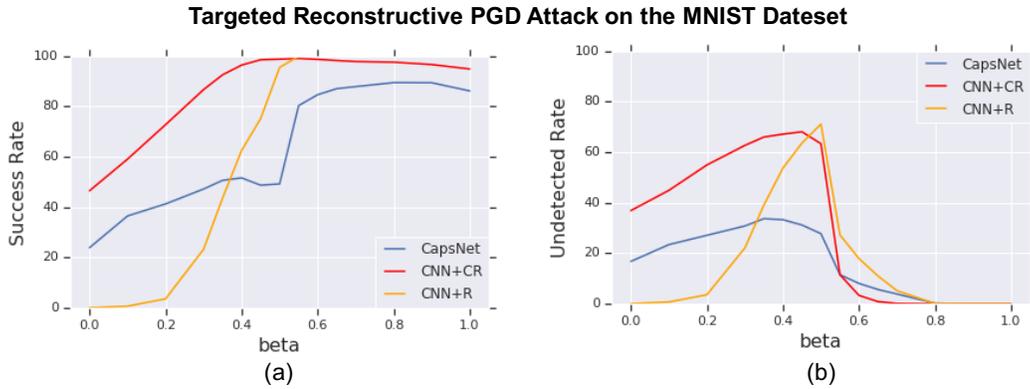


Figure 7: A plot of the success rate in (a) and undetected rate in (b) of targeted reconstructive PGD attack versus the hyperparameter β for each model on the MNIST test set. We set the max ℓ_∞ norm $\epsilon = 0.3$ to create the attacks.

E.2 Visual Coherence of the Reconstructive Attacks

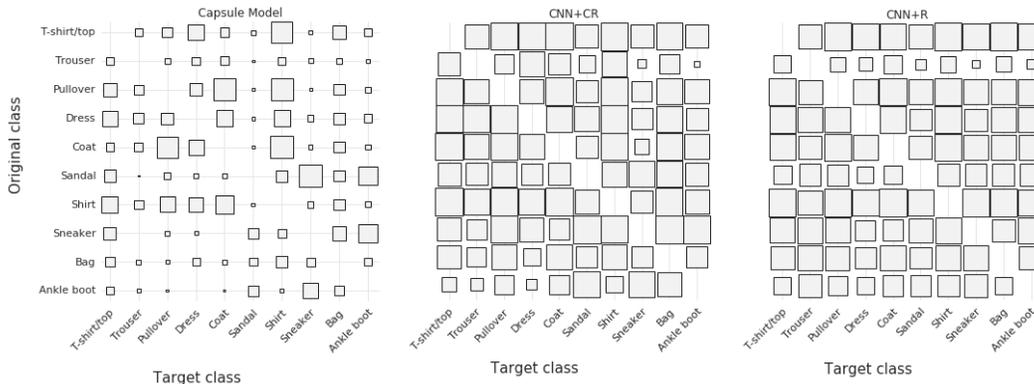


Figure 8: This diagram visualizes the adversarial success rates for each source/target pair for targeted reconstructive PGD attacks on FashionMNIST with $\ell_\infty = 25/255$. The size of the box at position x, y represents the success rate of adversarially perturbing inputs of class x to be classified as class y . We can see that there is significantly higher variance for the CapsNet model than for the two CNN models.

G Black Box Attacks

MNIST Dataset								
Targeted	CapsNet	CNN-CR	CNN-R		Untargeted	CapsNet	CNN-CR	CNN-R
PGD	1.5/0.0	7.8/0.0	7.4/0.0		PGD	8.5/0.0	32.6/0.0	27.6/0.0
R-PGD	4.2/1.0	18.3/11.0	11.3/4.8		R-PGD	10.4/2.4	42.7/24.9	25.2/8.9

Table 7: Success rate and undetected rate of **black-box** targeted and untargeted attacks. In the table, S_t/R_t is shown for targeted attacks and S_u/R_u is presented for untargeted attacks. All the numbers are shown in %.

H Visualization of Adversarial Examples and Reconstructions

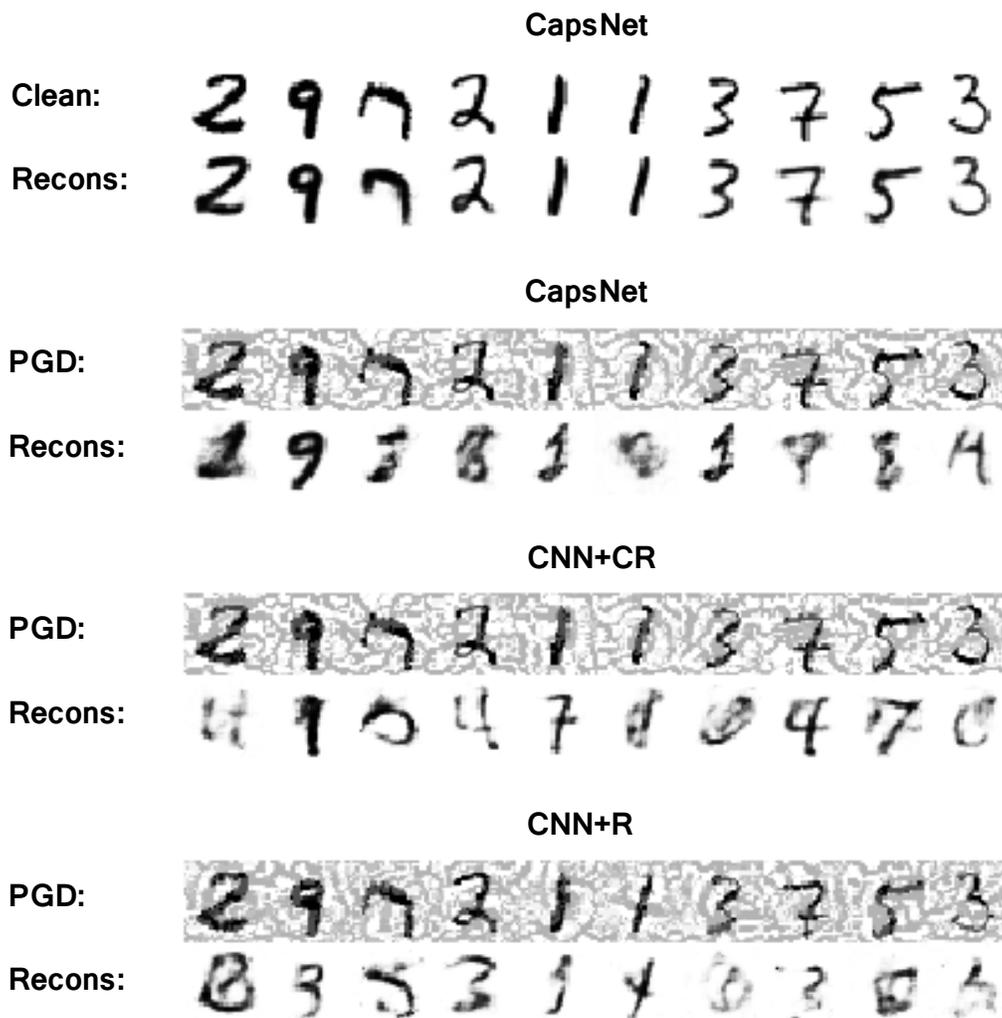


Figure 9: In the first set of images, the source clean image is presented in the first row with its reconstruction by Capsule Network presented in the second row. In the following three sets of images, the top row is the targeted PGD adversarial examples and the bottom is the corresponding reconstructed image by each model on the MNIST dataset.

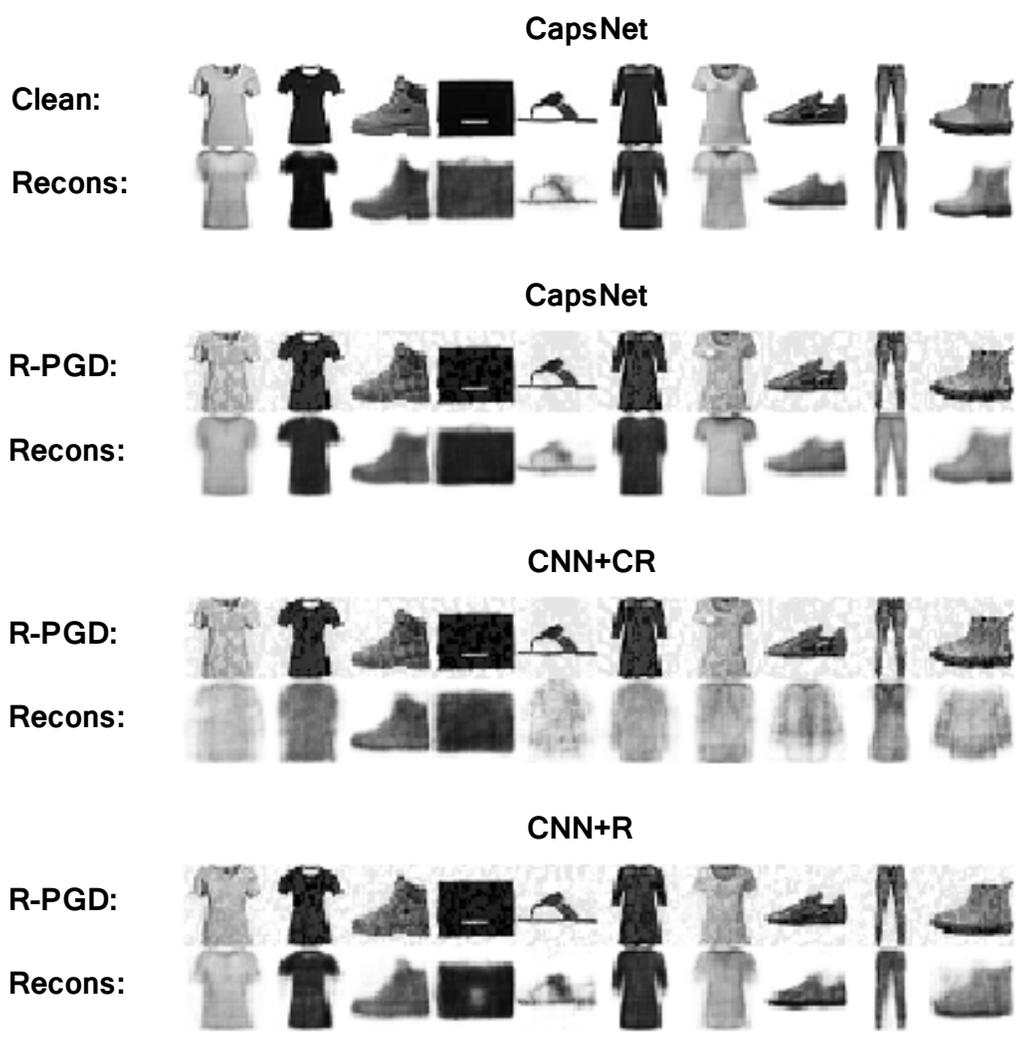


Figure 10: In the first set of images, the source clean image is presented in the first row with its reconstruction by Capsule Network presented in the second row. In the following three sets of images, the top row is the targeted Reconstructive R-PGD adversarial examples and the bottom is the corresponding reconstructed image by each model on the Fashion MNIST dataset.

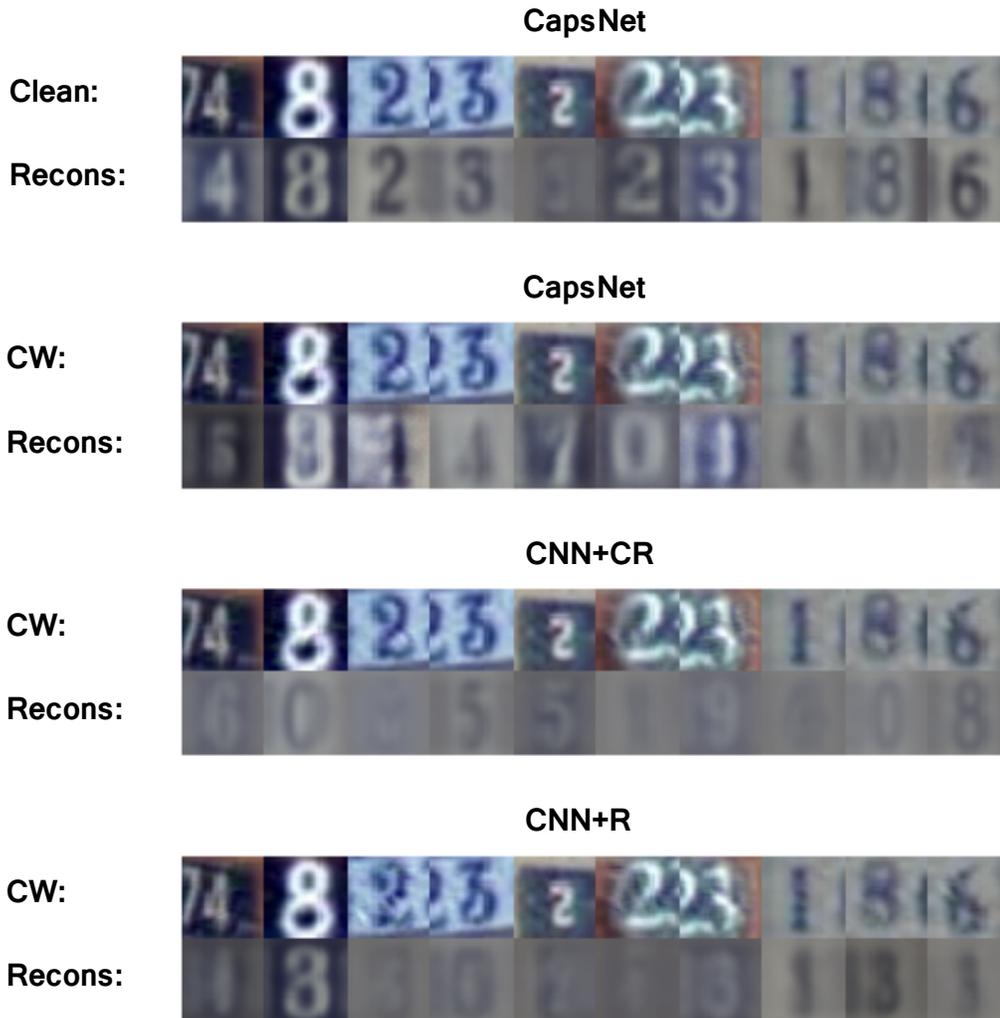


Figure 11: In the first set of images, the source clean image is presented in the first row with its reconstruction by Capsule Network presented in the second row. In the following three sets of images, the top row is the targeted CW adversarial examples and the bottom is the corresponding reconstructed image processed by each model on the SVHN dataset.

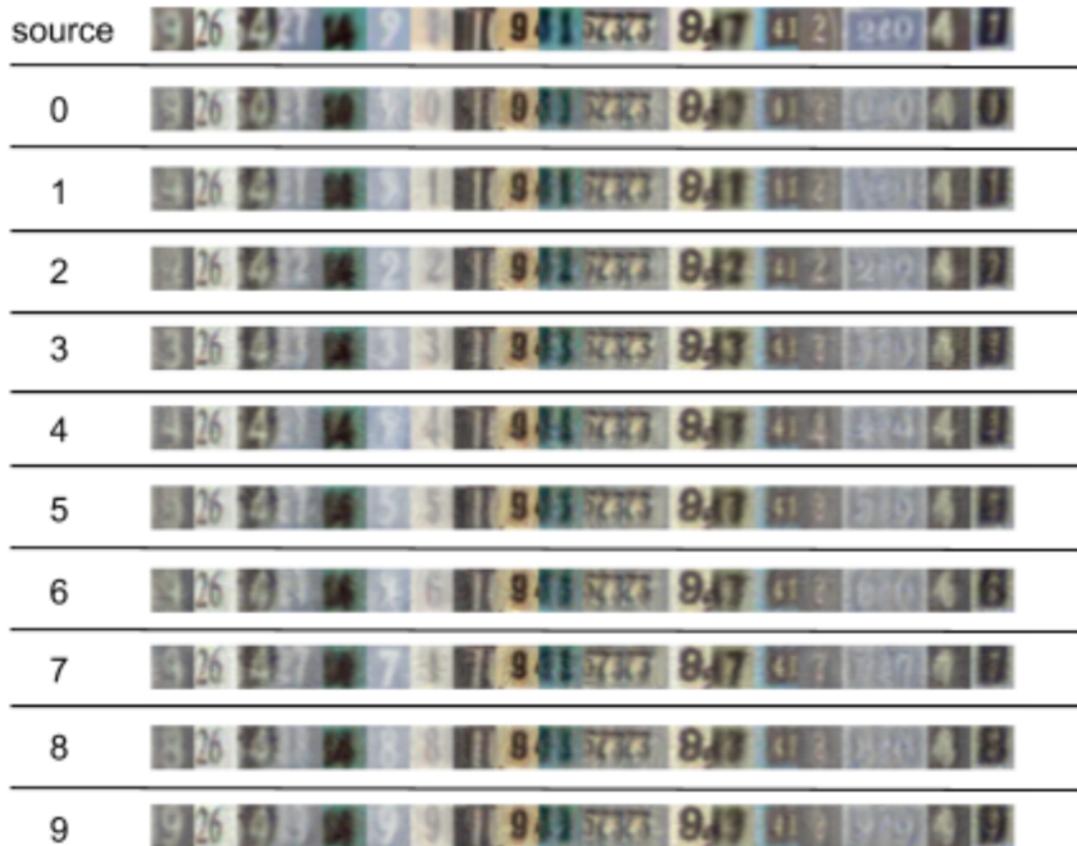


Figure 12: These are randomly sampled (not cherry picked) inputs (top row) and the result of adversarial perturbing them with targeted R-PGD against the CapsNet model with an epsilon bound of 25/255 based on the pixel range (other rows). Many of these attacks are not successful. Note the visual similarity between many of the attacks and the target class.