

1      **CMPT 489 Project: Adversarial Machine Learning in Computer Vision**

2  
3      SHRAY KHANNA, Simon Fraser University

4  
5      HAIPENG LI, Simon Fraser University

6  
7      JOHN SWEENEY, Simon Fraser University

8  
9      Adversarial machine learning is a very important and cutting-edge topic in the field of machine learning. Adversarial machine learning  
10     is a technique that attempts to fool machine learning models through malicious input. This paper conducts a survey, implementation,  
11     and defense against some of the most popular theoretical adversarial machine learning attacks. In the review section, there is an  
12     analysis and explanation of how an attack can theoretically be implemented. Implementation and the defense against implemented is  
13     conducted for selected attack methods. The methods analyzed and implemented demonstrate how differently models can act from  
14     what their designers intended under adversarial situations.

15  
16     Additional Key Words and Phrases: adversarial machine learning, neural networks, computer vision, system security

17  
18     **ACM Reference Format:**

19     Shrav Khanna, Haipeng Li, and John Sweeney. 2019. CMPT 489 Project: Adversarial Machine Learning in Computer Vision. 1, 1  
20     (December 2019), 11 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

21  
22     **1 INTRODUCTION**

23  
24     State-of-the-art deep neural networks have achieved major breakthroughs in many research areas such as Computer  
25     Vision[13][12][9], Natural Language Processing[10][21], and Fraud Detection. [18] Machine learning techniques were  
26     initially designed for controlled environments in which the training and test data are from the same statistical distribution.  
27     In the real world, adversaries can manipulate input data to exploit vulnerabilities of machine learning algorithms thus  
28     compromising the security of the machine learning system. [7] Adversarial machine learning is the technique to fool  
29     models with malicious input.

30  
31     Many implementations of deep neural networks in the real world are vulnerable to adversarial machine learning  
32     attacks. This paper focuses on adversarial machine learning attacks in the field of Computer Vision. Computer vision  
33     models are increasingly deployed in real-world environments like autonomous vehicles and next-generation employee-  
34     less stores, where robustness and security are crucial. In the case of a successful adversarial machine learning attack on  
35     an autonomous vehicle, there may be fatalities. This paper analyzes and implements published adversarial attacks on  
36     object classifications. In addition, novel attack defenses are implemented and explored.

37  
38     Authors' addresses: Shrav Khanna, Simon Fraser University, skhanna@sfu.ca; Haipeng Li, Simon Fraser University, haipengl@sfu.ca; John Sweeney,  
39     Simon Fraser University, sweeney@sfu.ca.

40  
41  
42     Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not  
43     made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components  
44     of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to  
45     redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

46  
47     © 2019 Association for Computing Machinery.

48     Manuscript submitted to ACM

49  
50     Manuscript submitted to ACM

## 53      2 REVIEWS OF ATTACKS ON OBJECT CLASSIFICATION

### 54      2.1 Fast Gradient Sign Method (FGSM)

55  
 56 Machine learning models misclassify inputs even on a high test-score when small worst-case perturbations are fed to  
 57 the model. The primary cause of Neural network behaving this way is the linearity in them. Szegedy et al [8] suggested  
 58 a method to efficiently compute an adversarial perturbation for a given image. Generally, machine learning models  
 59 misclassify examples that are only a little different from the data trained on them. A linearity in high dimension space  
 60 can cause adversarial examples.  
 61

62  
 63 Adversarial training gives a bonus to the regularization to control overfitting by adding a layer on top of strategies like dropout, pretraining and model averaging. A good performance on test set does not mean the model performs  
 64 naturally well, when exposed to a fake example these models generally fail. This is upsetting as the most common way  
 65 in computer vision is to use CNN features where Euclidean distance gives the perceptual distance. This resemblance is  
 66 flawed as when images with immensely small distance lies in the wrong classes in network's representation.  
 67

68  
 69 Linear explanation of the adversarial examples comes from the linear models. The precision of the features are  
 70 limited due to which a linear classifier responds same way to an input  $x$  as that of an adversarial input  $x' = X + N$ ,  
 71 where  $N$  is smaller than the precision of features. The classifier then assigns the same class to input  $x$  as that of input  
 72  $x'$ . Now, when weight vector gives dot product with the input to give a score for classifier we get:  
 73

$$74 \quad W^T X' = W^T X + W^T \cdot \eta$$

75 The perturbation now grows by  $w^T \cdot \eta$ . To maximize this assign  $\eta = sign(w)$ . This way the adversarial example grows  
 76 linearly with  $N$  ( $N$  being the input size). Szegedy et al[8] gives a way to compute these perturbations:  
 77

$$78 \quad \eta = \epsilon \cdot sign(\nabla_x J(\theta, x; y))$$

79 Where  $\eta$  is the perturbation,  $x$  the input to the model,  $y$  the targets associated with  $x$  (for machine learning tasks that  
 80 have targets) and  $J(\theta; x; y)$  be the cost used to train the neural network. They suggested to linearize the cost function  
 81 around the current value of  $\theta$ . This method is referred as fast gradient sign method.

82 This method generates the adversarial examples by FGSM which exploit the “linearity” of deep neural network models  
 83 which were thought as non-linear in high dimensional space.  
 84

85 The authors also developed a “one-step target class” where instead of using true label  $l$  they used  $l_t$  which is the  
 86 output class with least probability  $I_c$ . The authors also showed that the robustness of DNN is improved with use of  
 87 FGSM and its variants. Training with adversarial objective function was found to be a good regularization (Yuan et al  
 88 [22]). These results show that linearity is present even if it is not seen clearly and also models which optimize easily  
 89 can be perturbed easily.  
 90

### 91      2.2 Basic Iterative Method (BIM) and Iterative Least Likely Class Method (ILCM)

92 Kurakin et al [2] proposed two methods for adversarial attacks on images which are basically an improvement of Fast  
 93 Gradient Sign Method.  
 94

95 The idea behind Basic Iterative Method (BIM)[2] is to perturb images by taking a single large step in the direction that  
 96

increases the loss of the classifier. This gradient method is more similar to gradient ascent as compared to descent. The one-step methods perturb images by taking a single large step in the direction that increases the loss of the classifier (i.e. one-step gradient ascent). As presented by Chakraborty et al[4] this method is an intuitive extension to iteratively take multiple small steps while adjusting the direction after each step. The Basic Iterative Method does exactly that, and iteratively computes the following:

$$x_0^{adv} = x, \quad x_{N+1}^{adv} = Clip_{x,\epsilon}[x_N^{adv} + \alpha \cdot sign(\nabla_x J(x_N^{adv}, y_{true}))]$$

where:  $x$  is Clean Input Image;  $x_i^{adv}$  - Adversarial Image at  $i^{th}$  step;  $J$  is the Loss Function;  $y_{true}$  is the Model Output for  $x$ ;  $\epsilon$  is the Tunable Parameter and  $\alpha$  is the Step Size

The number of iterations are  $\min(\epsilon + 4, 1.25\epsilon)$ . This method does not typically rely on any approximation of the model and produces additional harmful adversarial examples when run for more iterations.

Similar to extending the FGSM[8] to its ‘one-step target class’ variation, Kurakin et al.[2] also extended BIM to Iterative Least-likely Class Method (ILCM). The adversarial examples generated by the ILCM method has been shown to seriously affect the classification accuracy of a modern deep architecture Inception v3, even for very small values of  $\epsilon$ , e.g.  $\epsilon < 16$ .

The equation of ILCM is given by:

$$x_0^{adv} = x, \quad x_{N+1}^{adv} = Clip_{x,\epsilon}[x_N^{adv} - \alpha \cdot sign(\nabla_x J(x_N^{adv}, y_{LL}))]$$

where:  $x$  is Clean Input Image;  $x_i^{adv}$  is the Adversarial Image at  $i^{th}$  step;  $J$  is the Loss Function;  $y_{LL}$  is for Least Likely Class;  $\epsilon$  is the Tunable Parameter and  $\alpha$  is the Step Size

For desired class, authors chose least likely class denoted by  $y_{LL}$  and is given by:

$$y_{LL} = \arg \min_y [P(y/x)]$$

For a well-trained classifier, the least-likely class is usually highly dissimilar from the true class, so ILCM attack results in more interesting mistakes, such as mistaking a cat for a flower.

Loss in FGSM, BIM and ILCM is given by:

$$Loss = \frac{1}{(m - k) + \lambda \cdot k} \left[ \sum_{i \in CLEAN} L(X_i/y_i) + \lambda \sum_{i \in ADV} L(X_i^{adv}/y_i) \right]$$

where:

$L(X|y)$  is a loss on a single example  $X$  with true class  $y$ ;  $m$  is total number of training examples in mini-batch;  $k$  is the number of training examples in mini-batch;  $\lambda$  is the parameter to control the relative weight of adversarial examples in the loss; The values we take are:  $\lambda = 0.3$ ,  $m = 32$  and  $k = 16$

### 157 2.3 One pixel attack for fooling deep neural networks

158 Su et al [20] proposed a novel black-box method for generating one-pixel adversarial perturbations that work on  
 159 any neural network. The attack relies on an algorithm called differential evolution (DE), which is a population based  
 160 stochastic function minimizer for solving complex multi-modal optimization problems. [19]

161 The idea behind the one pixel attack is to maximize the classifier's score for an adversary's target class given the  
 162 ability to only change a single pixel. The equation could also be written as an N pixel attack. The single pixel attack is  
 163 described by the following equation where  $f$  is the image classifier,  $adv$  is the target class and vector  $e(x)$  is an additive  
 164 adversarial perturbation according to  $x$ . [20]

$$\begin{aligned} 168 \quad & \max_{e(x)} f_{adv}(x + e(x)) \\ 169 \quad & \text{subject to } \|e(x)\| = 1 \end{aligned}$$

170 The perturbation is encoded into an array for the candidate solution, which is optimized by differential evolution.  
 171 [20] At each iteration the set of candidate solutions are produced by the DE formula:

$$174 \quad x_i(g+1) = x_{r1}(g) + F(x_{r2}g - x_{r3}g), r1 \neq r2 \neq r3$$

176 At each time step, the best candidate solutions will be chosen up to a defined value.

177 One pixel attacks for fooling deep neural networks have the advantage of being difficult or even impossible for a  
 178 human to discern the difference between the original and adversarial image. In addition, there does not currently exist  
 179 a true defense that does not have a computational cost. [20]

### 182 2.4 Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models

184 Brendel et al [3] propose the Boundary Attack, a decision-based attack that begins from a large adversarial perturbation  
 185 then to reduces the perturbation over multiple iterations while remaining adversarial. The attack is simple and unlike  
 186 FGSM [8], DeepFool [15] and the One Pixel Attack [20] are easy to deploy against real-world examples. In addition,  
 187 unlike these attacks, the boundary attack is a true black-box attack requiring no input from inside the model. Most  
 188 prominent examples are web services that run user uploaded-images through a neural network, such as Google Images.  
 189 The boundary attack is a decision-based attack. [3] There are no complete defenses against decision-based attacks  
 190 because the attack can exploit all types of vulnerabilities; it is an extremely flexible attack as long as an initial adversarial  
 191 image is locatable. Locating in initial adversarial image for the majority of models is trivial.

192 The boundary attack as its name describes tests the model for its boundary; how close can the the algorithm get to  
 193 the original image, using a variety of methods, while the image remains adversarial.

194 The attack can be described by the following algorithm: [3]

199 Data: original image  $o$ , adversarial criterion  $c(\cdot)$ , decision of model  $d(\cdot)$

200 Result: adversarial example  $\tilde{o}$  such that the distance  $d(o, \tilde{o}) = \|o - \tilde{o}\|_2^2$

201 while  $k <$  maximum number of steps do

203 draw random perturbation from proposal distribution  $n_k \sim P(\tilde{o}^{k-1})$ ;

204 if  $\tilde{o}^{k-1} + n_k$  is adversarial then

205     set  $\tilde{o}^k = \tilde{o}^{k-1} + n_k$

206 else

```

209     set tildeok = ok-1
210     end
211     k = k + 1
212   end
213
214

```

215     The attack needs to be initialized from a state that is already adversarial. In addition, at every time-step we must  
 216     efficiency draw possibilities from a proposal distribution. These possibilities much be within the input domain and  
 217     reduce the distance of the perturbed image from the original image.

## 219     2.5 Universal Adversarial Perturbations

220     So far we have discussed methods like FGSM [8], DeepFool [15] and One-pixel Attack [20] compute perturbation on a  
 221     single image. Moosavi-Dezfooli et al proposed Universal Adversarial Perturbations [14] to compute a "small universal  
 222     perturbation" being able to fool the state-of-the-art neural networks on a set of images with high fooling rate. The  
 223     perturbation is what the authors call "quasi-imperceptible" for human vision system.

224     Moosavi-Dezfooli et al proposes an iterative algorithm to compute the universal perturbation. For the formal definition  
 225     of perturbation, we assume that the clean image is sampled from a distribution  $\mu$ . A perturbation  $v$  is "universal" if it  
 226     satisfies the following constraints:

$$231 \quad P_{x \sim \mu}(\hat{k}(x + v) \neq \hat{k}(x)) \geq 1 - \delta \quad s.t. \quad ||v||_p \leq \xi \quad (1)$$

232     where  $P$  denotes the probability,  $\delta \in (0, 1]$  quantifies the desired fooling rate for all images sampled from the  
 233     distribution  $\mu$ ,  $||.||_p$  denotes the  $l_p$  norm and  $\xi$  controls the magnitude of the perturbation vector  $v$ .

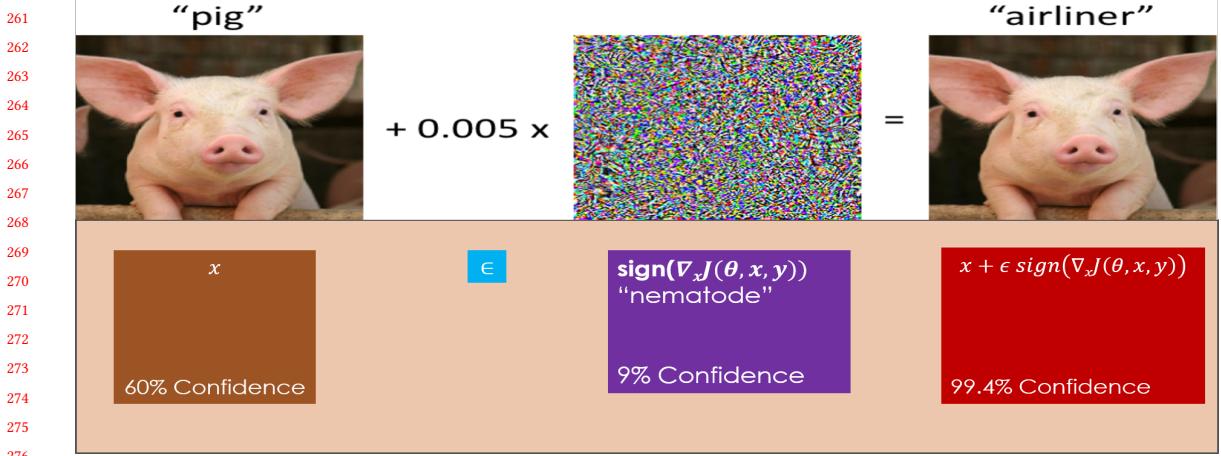
234     The authors compute the universal perturbation by limiting its  $l_2$  norm and  $l_\infty$  norm, and show that the upper bound  
 235     of its norm is 4% perturbation of the respective image norm, which has achieved a significant spoofing ratio of 0.8  
 236     or higher for the state-of-the-art image classifier. Their iterative method of calculating perturbations is related to the  
 237     DeepFool [15], which incrementally pushes a data point (or an image) to the decision boundary of its class. However,  
 238     in this case, all training data points are pushed in turn to their respective decision boundaries, and the perturbations  
 239     calculated on all images gradually accumulate by each time the accumulator is back-projected onto the desired sphere  
 240     of radius  $p$ .

## 244     3 IMPLEMENTATING FGSM[8]

245     Fast Gradient Sign Method has the lowest time complexity as compared to other methods that we have surveyed. The  
 246     reason behind that is because it deals with the sign of gradient and does not care about magnitude. This makes it faster  
 247     at execution. The scaling factor  $\epsilon$  (epsilon) is responsible for increasing or decreasing the noise (perturbation) in the  
 248     input image. We tested this method on various scenarios with a wide variety of classes. Through our experiments,  
 249     we've found out that this method has a limitation and performs differently for every state-of-the-art model. Fig 1 shows  
 250     how this attack works on an input image with a small perturbation.

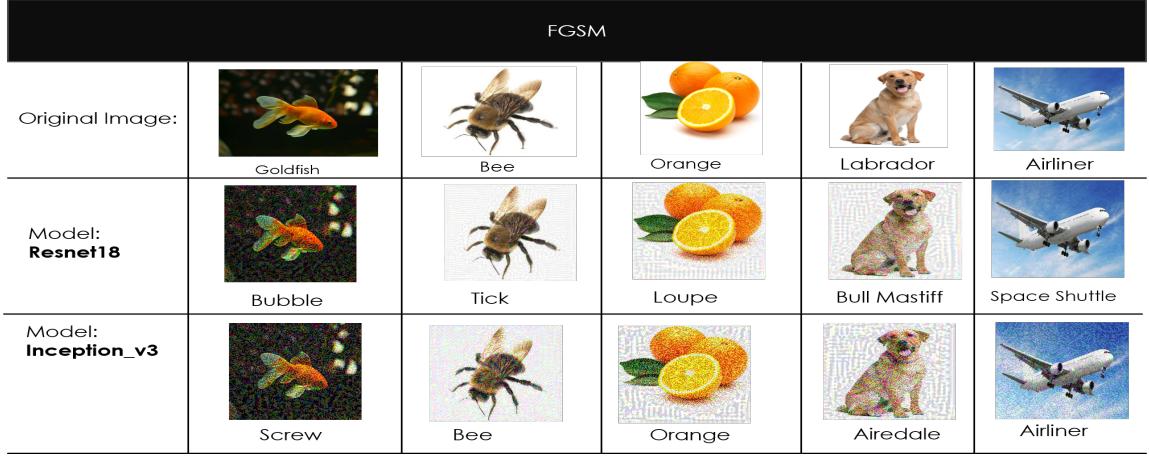
251     As you see that original image of Pig was reported with actual class (Pig) with a confidence of 60% but when we add  
 252     a small perturbation, the model assigns it with the class Airliner with a confidence level of 99.4%. The model in the  
 253     above example used is Resnet18.

254     We tested this method against Resnet18, Resnet101, Alexnet and Inception\_v3 which are currently state-of-the-art  
 255     models used for object classification and detection. After running a variety of test images we found that model has



277 Fig. 1. FGSM attack on an input image of Pig. FGSM depends on the sign of the gradient of cost function and not the magnitude.  
 278 Scale your perturbations using  $\epsilon$  (Epsilon).  
 279

280 limitations to Deeper Layered Networks such as Resnet101. One reason behind this is that it computes perturbation for  
 281 the entire input at once and thus, is not able to learn anything. We've shown a basic comparison of it's working on test  
 282 images in Fig2.  
 283



302 Fig. 2. FGSM attack on Resnet18 (model trained on 18 Residual layers) and Inception\_v3 (better and deeper layered model).  
 303

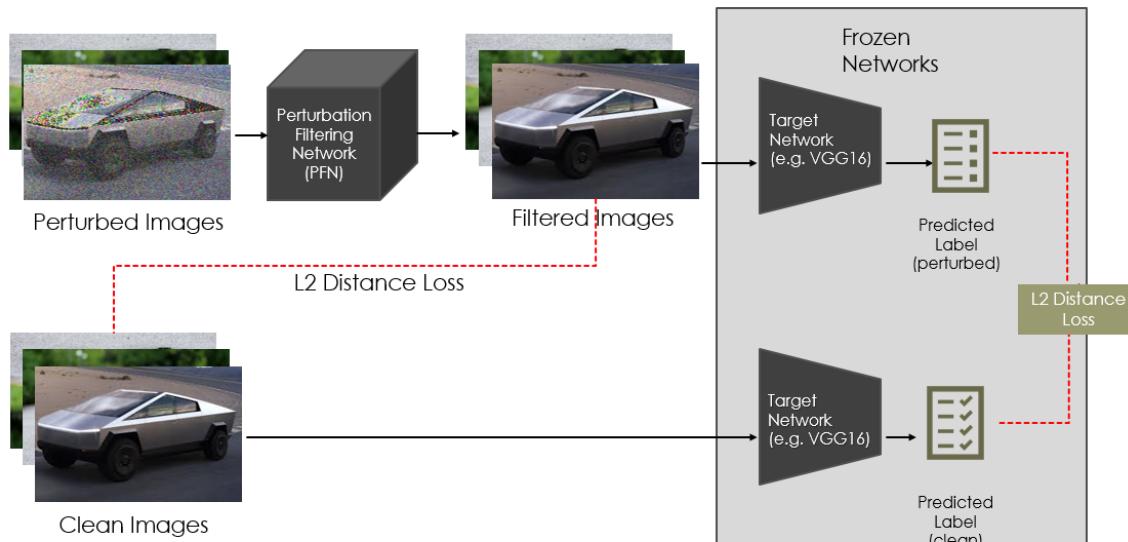
304 Fig 2 shows that even the perturbation added to the input image to show the noisy images which fooled the model.  
 305 As you can see we had to add a lot of perturbation for Inception\_v3 as compared to Resnet18. When we used Resnet18,  
 306 the amount of perturbation needed was relatively low for instance images of Bee and Airliner are almost the same. We  
 307 have used the maximum possible value for  $\epsilon$  in images of Orange and Airliner, after this value the image becomes very  
 308 hazy. The Inception\_v3 model is still predicting the correct class even after having the maximum perturbation.  
 309  
 310

## 313 4 UNIVERSAL PERTURBATION ATTACK AND DEFENSE

### 314 4.1 Overview

315 "Universal Perturbations" can be applied to "any" image to trick a state-of-the-art network classifier, seriously threatening  
 316 the success of deep learning in practice. We first implement a Universal Perturbation Attack and then propose a defense  
 317 framework to effectively protect the network from such attack. Our approach learns a Perturbation Filtering Network  
 318 (PFN) as filtering layers to a targeted model, such that no modification needs to be done on the target model. The PFN  
 319 is trained using real and synthetic image-agnostic perturbations. Our experiments show that our framework can reduce  
 320 the fooling rate of unseen adversarial perturbations from 81.1% to as low as 2.0%.

### 321 4.2 Proposed Approach



347 Fig. 3. We propose PFN: From the clean data, image-agnostic perturbations are computed and augmented with the synthetic  
 348 perturbations. Both the clean and perturbed images are fed into the Perturbed Filtering Network (PFN). PFN is trained by connecting  
 349 it to the first layer of the target network in order to keep the parameters of the target network unchanged during training.  
 350

351 We propose Perturbation Filtering Network (PFN), which is trained as perturbation filtering layers attached to the  
 352 first layer of the targeted network. We define the combination of PFN and target network as joint network. Fig. 3 shows  
 353 the network structure of joint network. We train the PFN using both clean and adversarial examples to ensure that the  
 354 image transformation learned by our network is not biased towards the adversarial examples.

355 We use a combined loss function:  $L_{rec}$  and  $L_{lab}$ .  $L_{rec}$  is the loss between the output image  $\hat{I}_i$  of PFN (we call it  
 356 filtered image) and the clean image  $I_i$  in pixel-level L2 distance.  $L_{lab}$  is the L2 loss between the predicted probability  
 357 label vector  $p_i, \hat{p}_i$  by target model and joint model respectively. These losses are defined as follows:

$$361 \quad 362 \quad 363 \quad 364 \quad L_{rec} = \frac{1}{N} \sum_{i=1}^N \sqrt{(I_i - \hat{I}_i)^2} \quad (2)$$

$$L_{lab} = \frac{1}{N} \sum_{i=1}^N \sqrt{(\mathbf{p}_i - \hat{\mathbf{p}}_i)^2} \quad (3)$$

Unlike the defense network proposed by N Akhtar [1], which minimizes the cross-entropy loss [6] of the predicted labels, we minimize the sum of  $L_{rec}$  and  $L_{lab}$  during training using weight parameters  $\alpha$  and  $\beta$  to control the contribution of two losses:

$$L_{total} = \alpha L_{rec} + \beta L_{lab} \quad (4)$$

We train the PFN using both clean and adversarial examples to ensure that the image transformation learned by our network is not biased towards the adversarial examples. The PFN is implemented as 10-ResNet blocks [9] sandwiched by convolution layers. The 224x224x3 input image is fed to Conv 3x3, stride = 1, feature maps = 64, "same" convolution; followed by 10 ResNet blocks, where each block consists of two convolution layers with ReLU activations [16]. We minimize Eq. 4 with ADAM optimizer [11]. We set the learning rate to 0.005. We used mini-batch size of 16, and trained the PFN for a given targeted network for at least 3 epochs.

## 5 ATTACK & DEFENSE EXPERIMENTS

### 5.1 Dataset

Our dataset is based on Natural Images [17] - a compiled dataset of 6899 images from 8 distinct classes (airplane, car, cat, dog, flower, fruit, motorbike, person), which is used for training a classifier as target model. When we do the attack and defense, we use a subset of Natural Images dataset which contains 1600 images for training (200 images per class) and around 640 images (around 80 images per class) for testing. Note that this is a relatively small dataset considering our limited computing power. In the future we plan to use the original ImageNet dataset or its subset for the experiments.

### 5.2 Train a Classifier as Target Model

We conduct experiments on two target models: In experiment 1 we attack then the Resnet18  $M_1$  [9] pre-trained on ImageNet [5], which has 1000 classes. In experiment 2 we trained Resnet18  $M_2$  on our dataset - Natural Images which has 8 classes. We are able to train  $M_2$  with 99% accuracy on Natural Images dataset. Note that  $M_1$  is not trained on Natural Images dataset, so each image will be classified as one of 1000 classes from ImageNet, which is not be as accurate as  $M_2$ . We show the difference in next section.

### 5.3 Compute Universal Perturbation

We compute 12 different Universal Perturbation on a subset of Natural Images dataset with 1600 training images against both  $M_1$  and  $M_2$  using our implementation of [14]. We notice that the attack on  $M_1$  is much easier than  $M_2$  with 92% fooling rate in 3 epochs for the former but only 81.1% fooling rate in 8 epochs for the latter. This implies that the Universal Perturbation Attack has better success rate for a model with more output classes. We show the perturbed images for experiment 2 in the second column of Fig 4.

### 5.4 Defense Against Universal Perturbation

We train the PFN against the attack on model  $M_1$  and  $M_2$  by minimizing Eq. 4 where  $\alpha = 40$  and  $\beta = 1$ . 10 different Universal Perturbations are used for training and 2 for testing. We evaluated the performance by *fooling rate*, which is

	Original	Perturbed	Trained with $L_{total}$	Trained with $L_{rec}$	Trained with $L_{lab}$
Fooling Rate	0%	81.1%	2.0%	3.3%	5.1%
Images, predicted labels and confidence					
	Airplane 0.97	Fruit 0.54	Airplane 0.75	Airplane 0.73	Airplane 0.62
Images, predicted labels and confidence					
	Person 0.99	Fruit 0.67	Person 0.99	Person 0.94	Person 0.90

Fig. 4. Representative examples of the defense results. First, second columns show the original images and images with universal perturbation. The last three columns compare the filtered image quality, fooling rates and prediction by our PFN trained with  $L_{total}$ ,  $L_{rec}$ ,  $L_{lab}$  respectively.

the percentage of perturbed images successfully fool the target network by making it predict wrong labels. In experiment 1 we are able to reduce the fooling rate from 92% to 33%, and in experiment 2 from 81.1% to 2.0% in 3 epochs. We can see our methods are particularly effective on experiment 2 with the ResNet trained on our small dataset. As mentioned in previous section, this dataset is hard to fool using Universal Perturbation. So it is not surprising that it is easy to achieve such a good defense result. Some representative examples are demonstrated in the first three columns in Fig 4. Note that the fooling rates are tested on testing dataset, which means the network never saw the perturbation nor the images before during testing phase, which shows that our defense framework has very good generalizability.

We conduct an ablation study to validate each term of our loss function in Eq. 4, shown in last two columns in Fig 4. Comparing the three methods with different loss function, we can clearly see the superiority of our loss function  $L_{total}$ , which produces better image quality which is the closest to the original image, lower fooling rate and higher confidence predicted by the target model. When trained without  $L_{rec}$  (5th column in Fig 4), the color filtered images are brighter than original images as a result of lacking the constraint of recovering the original images in the loss function. When trained with only  $L_{rec}$  without  $L_{lab}$  (4th column in Fig 4), the color distribution of filtered images are much closer to the original images, however, it still gives worse fooling rate with lower prediction confidence. By combining  $L_{rec}$  and  $L_{lab}$  we are able to get the best results. See our supplementary material for more results.

## 6 CONCLUSION

The examples that were presented in this paper show how many modern machine learning algorithms are vulnerable in surprising ways. The failure of simple machine learning models, when subjected to the attacks, demonstrate how differently the models act from what their designers intended. Fortunately, there are excellent solutions for the specific attacks described in this paper that make the risk of a particular type of attack negligible or substantially lower the

chance of success. However, the defenses we proposed are not completely adaptive and doesn't perform well on certain dataset; it may leave another vulnerability open to attackers who know the defense being applied.

## CONTRIBUTION

Here are the contributions of each group member in the project:

**Shray Khanna:** Surveyed two papers: Fast Gradient Sign Method (FGSM[8]) and Basic Iterative Method (and Iterative Least Likely Class Method [2]). Implemented the first method (FGSM) on several models to test the results. Made and presented the final presentation. Wrote partial documentation.

**Haipeng Li:** Surveyed and implemented the Universal Perturbation Attack [14]. Proposed and implemented the defense against Universal Perturbation Attack: Perturbation Filtering Network (PFN). Made slides for the initial presentation and presented in both presentations. Wrote partial documentation.

**John Sweeney:** Surveyed two papers: the One Pixel Attack[20] and Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. [3] Partially implemented the One Pixel Attack[20], before it was decided as a group to remove One Pixel Attack Attack due to the defense being redundant. Presented on One Pixel Attack for final presentation. Wrote the abstract, introduction and conclusion for the paper.

## REFERENCES

- [1] Naveed Akhtar, Jian Liu, and Ajmal Mian. 2018. Defense against universal adversarial perturbations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 3389–3398.
- [2] Ian J. Goodfellow Alexey Kurakin and Samy Bengio. 2016. Adversarial examples in the physical world. *CoRR* abs/1607.02533 (2016). <http://arxiv.org/abs/1607.02533>
- [3] Wieland Brendel, Jonas Rauber, and Matthias Bethge. 2017. Decision-Based Adversarial Attacks: Reliable Attacks Against Black-Box Machine Learning Models. *arXiv:stat.ML/1712.04248*
- [4] Anirban Chakraborty, Manaa Alam, Vishal Dey, Anupam Chattopadhyay, and Debdeep Mukhopadhyay. 2018. Adversarial Attacks and Defences: A Survey. *CoRR* abs/1810.00069 (2018). <http://arxiv.org/abs/1810.00069>
- [5] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. Ieee, 248–255.
- [6] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. Deep learning. MIT press.
- [7] Ian Goodfellow, Patrick McDaniel, and Nicolas Papernot. 2018. Making Machine Learning Robust Against Adversarial Inputs. In *Communications of the ACM*. 56–66.
- [8] Ian J Goodfellow, Jonathon Shlens, and Christian Szegedy. 2014. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572* (2014).
- [9] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 770–778.
- [10] Sepp Hochreiter and Jürgen Schmidhuber. 1997. Long short-term memory. *Neural computation* 9, 8 (1997), 1735–1780.
- [11] Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980* (2014).
- [12] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. 2012. Imagenet classification with deep convolutional neural networks. In *Advances in neural information processing systems*. 1097–1105.
- [13] Yann LeCun, Koray Kavukcuoglu, and Clément Farabet. 2010. Convolutional networks and applications in vision. In *Proceedings of 2010 IEEE International Symposium on Circuits and Systems*. IEEE, 253–256.
- [14] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, Omar Fawzi, and Pascal Frossard. 2017. Universal adversarial perturbations. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1765–1773.
- [15] Seyed-Mohsen Moosavi-Dezfooli, Alhussein Fawzi, and Pascal Frossard. 2016. Deepfool: a simple and accurate method to fool deep neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2574–2582.
- [16] Vinod Nair and Geoffrey E Hinton. 2010. Rectified linear units improve restricted boltzmann machines. In *Proceedings of the 27th international conference on machine learning (ICML-10)*. 807–814.
- [17] Prasun Roy, Subhankar Ghosh, Saumik Bhattacharya, and Umapada Pal. 2018. Effects of Degradations on Deep Neural Network Architectures. *arXiv preprint arXiv:1807.10108* (2018).

- 521 [18] SamanehSorournejad, Zahra Zojaji, Reza Ebrahimi Atani, and Amir Hassan Monadjemi. 2016. A Survey of Credit Card Fraud Detection Techniques:  
522 Data and Technique Oriented Perspective. arXiv:1611.06439v1  
523 [19] Rainer Storn and Kenneth Price. 1997. Differential Evolution – A Simple and Efficient Heuristic for global Optimization over Continuous Spaces.  
524 *Journal of Global Optimization* 11, 4 (01 Dec 1997), 341–359. <https://doi.org/10.1023/A:1008202821328>  
525 [20] Jiawei Su, Danilo Vasconcellos Vargas, and Kouichi Sakurai. 2019. One Pixel Attack for Fooling Deep Neural Networks. *IEEE Transactions on*  
526 *Evolutionary Computation* 23, 5 (Oct 2019), 828–841. <https://doi.org/10.1109/tevc.2019.2890858>  
527 [21] Ilya Sutskever, Oriol Vinyals, and Quoc V Le. 2014. Sequence to sequence learning with neural networks. In *Advances in neural information processing*  
528 *systems*. 3104–3112.  
529 [22] Xiaoyong Yuan, Pan He, Qile Zhu, Rajendra Rana Bhat, and Xiaolin Li. 2017. Adversarial Examples: Attacks and Defenses for Deep Learning. *CoRR*  
530 abs/1712.07107 (2017). arXiv:1712.07107 <http://arxiv.org/abs/1712.07107>
- 531  
532  
533  
534  
535  
536  
537  
538  
539  
540  
541  
542  
543  
544  
545  
546  
547  
548  
549  
550  
551  
552  
553  
554  
555  
556  
557  
558  
559  
560  
561  
562  
563  
564  
565  
566  
567  
568  
569  
570  
571  
572