

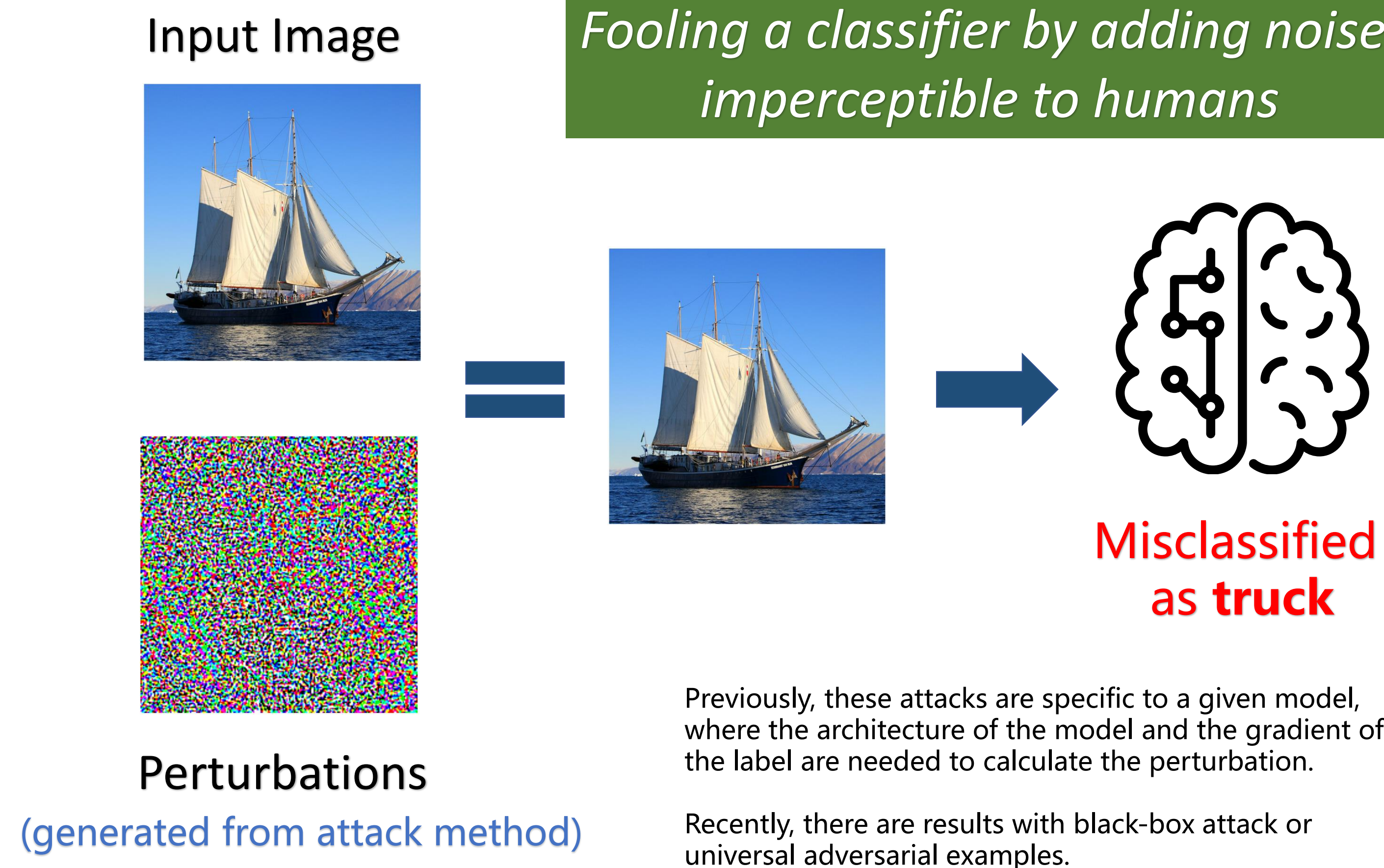
PIXELDEFEND: LEVERAGING GENERATIVE MODELS TO UNDERSTAND AND DEFEND AGAINST ADVERSARIAL EXAMPLES

Authors: Yang Song, Taesup Kim, Sebastian Nowozin, Stefano Ermon, Nate Kushman

Reproducer: Ruifan Yu

David R. Cheriton School of Computer Science, University of Waterloo

Adversarial Attack



Attack Methods

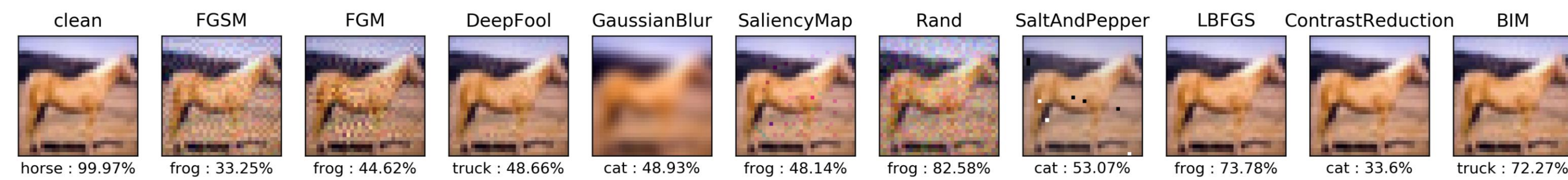
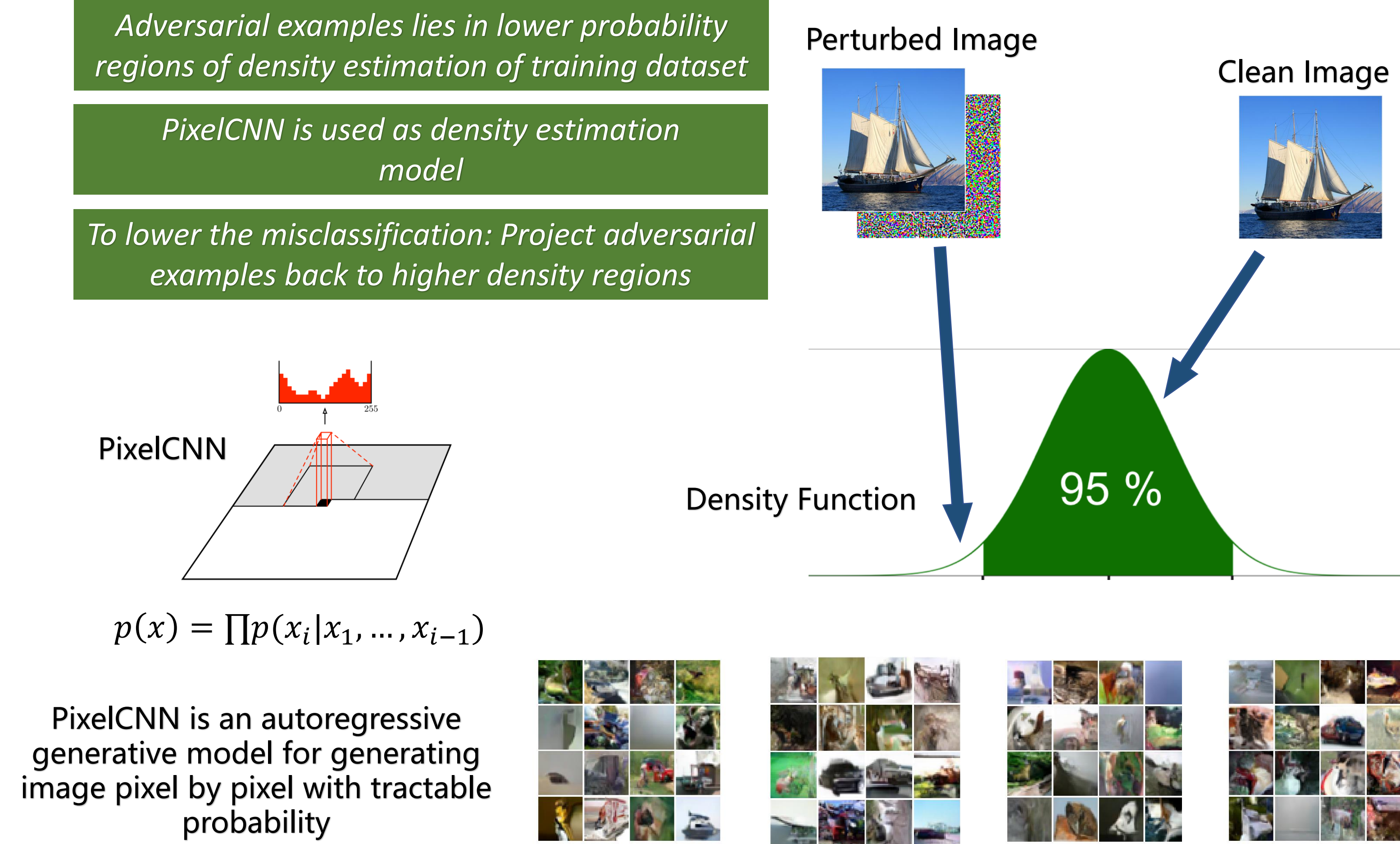


Table 1: Ten different attack methods in this reproduction project

	Attack Method	Description
1	Gradient Sign Attack (FGSM)	Adds the sign of the gradient to the image, gradually increasing the magnitude until the image is misclassified.
2	Gradient Attack (FGM)	Perturbs the image with the gradient of with respect to the loss of label.
3	DeepFool Attack (DeepFool)	Attack the model by finding the nearest decision boundary and making linear iteration to perturb images
4	GaussianBlur Attack (GaussianBlur)	Blurs the image until it is misclassified.
5	Saliency Map Attack (SaliencyMap)	Attack the model by calculating the saliency map of target label.
6	Uniform Noise Attack (Rand)	Adds uniform noise to the image, gradually increasing the standard deviation until the image is misclassified.
7	Salt And Pepper Noise Attack (SaltAndPepper)	Increases the amount of salt and pepper noise until the image is misclassified.
8	LBFGS Attack (LBFGS)	Uses L-BFGS-B to minimize the distance between the image and the adversarial as well as the cross-entropy between the predictions for the adversarial and the one-hot encoded target class.
9	Contrast Reduction Attack (ContrastReduction)	Reduces the contrast of the image until it is misclassified.
10	Iterative FGSM (BIM)	Like GradientSignAttack but with several steps for each epsilon.

Motivation



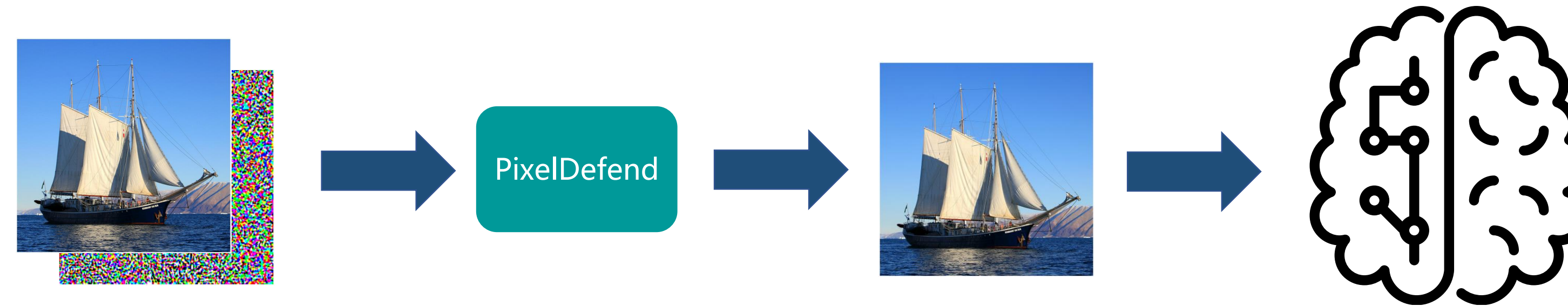
PixelDefend

Algorithm 1 PixelDefend

Input: Image X , Defense parameter ϵ_{defend} , Pre-trained PixelCNN model p_{CNN}

Output: Purified Image X^*

```
1:  $X^* \leftarrow X$ 
2: for each row  $i$  do
3:   for each column  $j$  do
4:     for each channel  $k$  do
5:        $x \leftarrow X[i, j, k]$ 
6:       Set feasible range  $R \leftarrow [\max(x - \epsilon_{\text{defend}}, 0), \min(x + \epsilon_{\text{defend}}, 255)]$ 
7:       Compute the 256-way softmax  $p_{\text{CNN}}(X^*)$ .
8:       Update  $X^*[i, j, k] \leftarrow \arg \max_{z \in R} p_{\text{CNN}}[i, j, k, z]$ 
9:     end for
10:   end for
11: end for
```



The purpose of this paper is to adopt a PixelCNN generative model for detecting and defending against adversarial perturbed images as it is observed that most adversarial examples lie in low probability regions.

The output from PixelCNN is the exact probability estimation of the pixel value. By training the PixelCNN model on training dataset, the model is predicting the probability of given images that belongs to the same distribution of training data.

$$p = \frac{1}{N+1} \left(\sum_{i=1}^N \mathbb{I}[T_i \leq T] + 1 \right) = \frac{T+1}{N+1} = \frac{1}{N+1} \left(\sum_{i=1}^N \mathbb{I}[p_{\text{CNN}}(\mathbf{X}_i) \leq p_{\text{CNN}}(\mathbf{X}')] + 1 \right)$$

Intuitively, the p-value is telling how many images in training data are less likely to be sampled from training data distribution. For normal images, the rank of its p value should be random since every image is equally likely sampled from training data distribution.

$$\text{BPD}(\mathbf{X}) \triangleq -\log p_{\text{CNN}}(\mathbf{X}) / (I \times J \times K \times \log 2)$$

Here, Bits Per Dimension is another measurement of the probability, which is same as negative log-likelihood.

Reproducible?

Yes

Claim #1

Adversarial examples mainly lie in the low probability regions of the training distribution, regardless of attack. The inputs are more outside of the training distribution if their p-value distribution has a larger deviation from uniform.

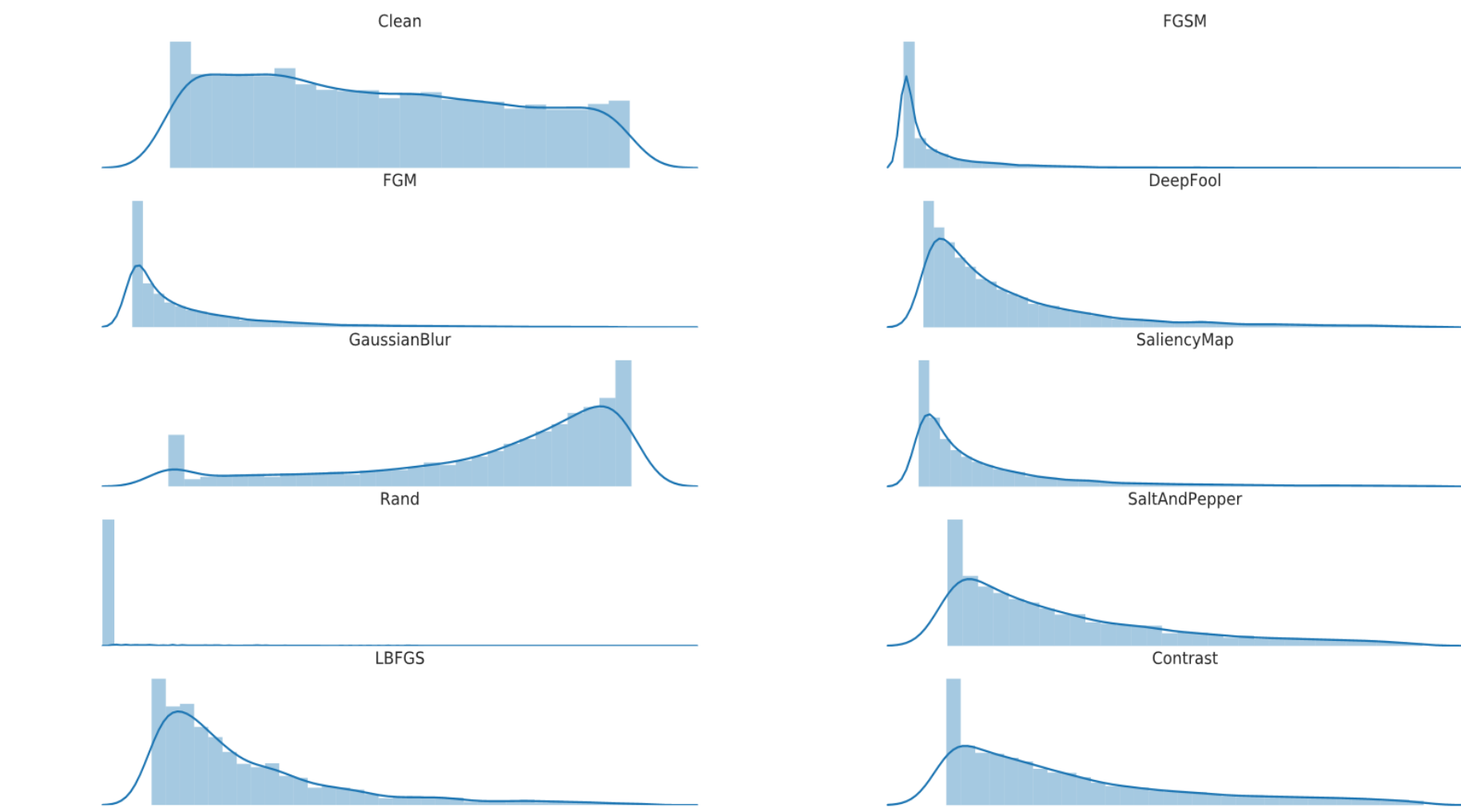


Figure 1: The distribution of the p-value of ten attacks on ResNet18 on CIFAR-10 dataset

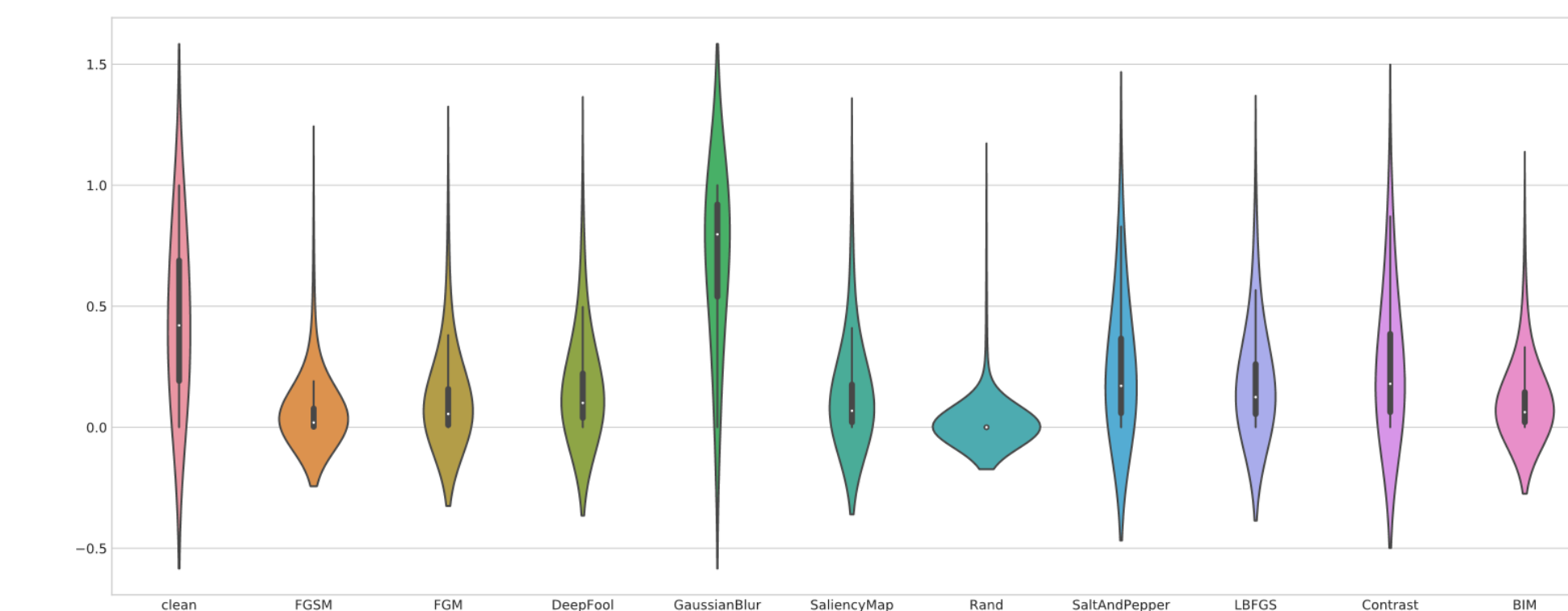


Figure 2: The violin plot of the Bits Per Dimension(BPD) Distributions over ten attack methods

The distribution of log-likelihoods show considerable difference between perturbed images and clean images. The log-likelihoods from PixelCNN also provide a quantitative measure for detecting adversarial examples.

* PixelCNN model has no information about the attacking methods for producing those adversarial examples, and no information about the ResNet model either



Claim #3

PixelDefend is a model-attack-agnostic defense strategy that can make model predict the correct label.

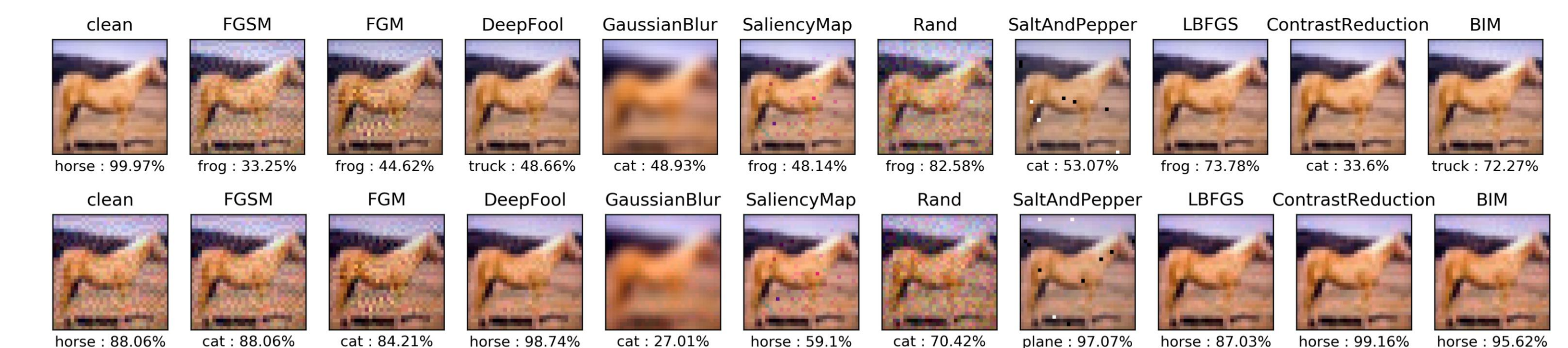


Figure 3: the images after the purification by PixelDefend
The texts below each image are the predicted label and confidence given by ResNet

Paper Reviews



- This paper introduces a novel family of methods for defending against adversarial attacks based on the idea of purification.
- PixelDefend can be used to protect already deployed models and be combined with other model-specific defenses.
- Adversary to PixelCNN. It is not clear why a PixelCNN may not be adversarially attacked, nor if such a model would be able to guard against an adversarial attack.. Are there adversarial examples for PixelCNN model?

In our opinion, it would be interesting to be what is the performance after the model is fine-tuned on the PixelDefend outputs. However, due to the slow generation speed, we were not able to generate enough purified images for further experiments. Since this project is a reproduction challenge, we focused on the PixelCNN model. As future works, we will try other generative models and make a comparison.

