

Guided Perturbations: Self Corrective Behavior in Convolutional Neural Networks

Swami Sankaranarayanan *
 University of Maryland
 College Park, MD
 swamiviv@umiacs.umd.edu

Arpit Jain
 GE Global Research Center
 Niskayuna, NY
 Arpit.Jain@ge.com

Ser Nam Lim
 GE Global Research Center
 Niskayuna, NY
 limser@ge.com

Abstract

Convolutional Neural Networks have been a subject of great importance over the past decade and great strides have been made in their utility for producing state of the art performance in many computer vision problems. However, the behavior of deep networks is yet to be fully understood and is still an active area of research. In this work, we present an intriguing behavior: pre-trained CNNs can be made to improve their predictions by structurally perturbing the input. We observe that these perturbations - referred as Guided Perturbations - enable a trained network to improve its prediction performance without any learning or change in network weights. We perform various ablative experiments to understand how these perturbations affect the local context and feature representations. Furthermore, we demonstrate that this idea can improve performance of several existing approaches on semantic segmentation and scene labeling tasks on the PASCAL VOC dataset and supervised classification tasks on MNIST and CIFAR10 datasets.

1. Introduction

Convolutional Neural Networks (CNNs) have achieved state of the art results on several computer vision benchmarks such as ILSVRC [13] and PASCAL VOC [3] over the past few years. Despite their overwhelming success, recent results have highlighted that they can be sensitive to small adversarial noise in the input [4] or can be easily fooled using structured noise patterns [12]. To understand how a CNN can learn complex and meaningful representations but at the same time be easily fooled by simple and imperceptible perturbations still remains an open research problem. The work of Goodfellow *et. al.* [4] and Szegedy *et al.* [15] among others, bring out the intriguing properties of neural networks by introducing perturba-

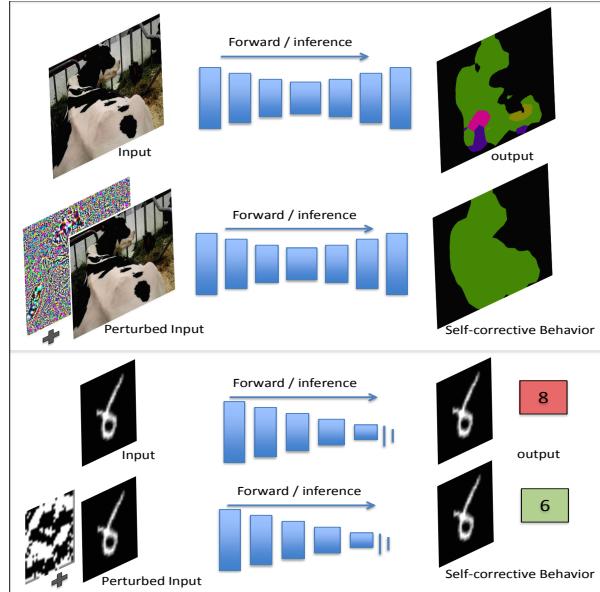


Figure 1: Self-corrective behavior due to Guided Perturbations (GP) for segmentation and classification tasks.

tions in either the hidden layer weights or the input image. While these approaches have focused on understanding the effect of adversarial noise in deep networks, in this work we present an interesting observation: input perturbations can enable a CNN to correct its mistakes. We find that these perturbations exploit the local neighborhood information from network's prediction which in turn results in contextually smooth predictions.

In almost all the CNN based approaches, the output is obtained using a single forward pass during the prediction time. In the proposed approach, we use the prediction made by the network during the forward pass to generate perturbations at the input. Specifically, we backpropagate the gradient of the prediction error all the way to the input. We would like to emphasize that the error gradients are generated purely based on the network's prediction without any knowledge of ground truth. We *perturb* the input image by

*Work performed while interning at GE Global Research

adding to it a scaled version of the gradient signal. This is fed back to the network again for prediction. Figure 1 shows an example of the self-corrective behavior of the generated perturbations for segmentation and classification tasks. This example shows that these perturbations of the input image could be viewed as a form of structured distortion that is added to the input such that the context gets amplified in each pixel’s neighborhood which enables the network to correct its own mistakes. The proposed approach is simple and easy to implement and does not require retraining or modification in network’s architecture.

Existing approaches to improve performance on segmentation and classification tasks have been geared towards novelties in network architecture or using large amount of training data or both. While these are valid ways to improve the networks performance, the proposed approach highlights an inherent behavior of CNNs that can be used to improve their prediction without requiring additional learning or training data. We would like to note here that while the behavior of Guided Perturbations (GP) is similar to Conditional Random Fields based approaches, the difference in our case is that there is no explicit modeling of context or neighborhood interactions. Since our approach is network independent, this doesn’t preclude networks which model context explicitly and we show improvements in such networks too.

To the best of our knowledge, this is the first approach to show existence of a self-corrective behavior in CNNs and use of such behavior for improvement in performance on segmentation and classification tasks. To summarize, the major contributions of this paper are:

- We present a novel and intriguing observation: there exist structured perturbations which when used to perturb the input leads to a corrective behavior in CNNs.
- We propose a generalized framework to improve the performance of any pretrained CNN model that is architecture independent and requires no learning assuming the network is trained end-to-end.

2. Background

In recent years, there have been several approaches that attempt to analyze the behavior of CNNs for classification problems. Mahendran *et al.* [10] proposed an approach to invert the function learned by the CNN in order to generate as faithful a reconstruction of the input as possible. This is performed by minimizing a regularized energy function that approximates the representation function that is learned by the deep network. Another interesting work in this direction is the Fooling Images work of Nguyen *et al.* [12] that is further extended by Yosinski *et al.* [17]. The main objective in both the approaches is to synthesize images to confuse CNN by maximizing the activation of individual neurons

from different layers of a deep network. This leads to interesting results such as images that look like random noise but which the CNN classifies into different classes with high confidence. The approaches that are closer in spirit to our proposed approach are the ones that predate these recent ones: Szegedy *et al.* [15] and Goodfellow *et al.* [4]. Their study shows that there exist a lot of adversarial examples which are the result of minor perturbations of the input that causes the CNN to misclassify input images on classification tasks; these examples can be generated by adding a fraction of the gradient that is generated by wiggling the classifier output in the direction of the target class.

One of the applications that this paper focuses on is semantic segmentation. A lot of research have gone into understanding the expressive ability of CNNs for such problems. Recent methods for image segmentation such as Fully Convolutional Networks (FCN) by Long *et al.* [14] have provided an easy framework that casts the image segmentation problem as a pixelwise label classification problem. The major difference in their work was the image level output generation and backpropagation which was made possible by the work of Zeiler *et al.* [19]. This image level back propagation provides a simple way to learn a discriminative representation of classes at the pixel level. Several recent approaches such as CRFasRNN by Zheng *et al.* [20], DeepLab by Chen *et al.* [2] and GCRF by Vemulapalli *et al.* [16] have improved the FCN framework by explicitly modeling context. CRFasRNN casts the CRF iterations, which has been traditionally used as a post processing function in image segmentation problem to ensure label compatibility, as a Recurrent Neural Network. They formulate the steps required to perform a mean field iteration in a CRF including message passing and learning a label compatibility transform as a layer in a CNN, which is unrolled in time over T iterations. The unary potentials are computed using the FCN-8s network which is then refined using the RNN structure. By casting this as a CNN layer they perform end-to-end training. More recently, Yu *et al.* [18] propose to train a multiscale context aggregation module on top of a modified FCN-8s network. This context module improves the performance of the base network on its own or in combination with CRF based approaches.

In this paper, we describe an interesting property of deep networks about how it can change its predictions for the better using minor perturbations of the input. Furthermore, we provide useful applications of our approach by showing how it can be used to improve prediction performance on challenging computer vision tasks.

3. Our Approach

In this section, we describe our approach to generate guided perturbations by using the gradient information obtained from the network’s output. We perform experiments

to study different aspects of these perturbations and how they affect the network representations. Since our approach to generate guided perturbation is different for segmentation and classification tasks, we discuss them separately.

3.1. Semantic Segmentation

Figure 2 illustrates our approach for semantic segmentation task. Given an input image we perform a forward pass to compute the output - which is usually the output of a softmax function that gives a class probability vector for each pixel. The prediction output is then binarized by setting the probability of the most confident class to one and the others to zero. This is done for each pixel and the error gradient is computed at the softmax layer by setting this modified output as ground truth. Let $\mathbb{X} \in \mathcal{R}^{M \times N \times C_{in}}$ represent the input image to the deep network, $\mathbb{Y} \in \mathcal{R}^{M \times N \times C_{out}}$ represent ground truth labeling, where C_{in} is the number of input channels, C_{out} is the number of classes and $M \times N$ is the dimensionality of the input image. Let θ represent the parameters of the network and $\mathcal{J}(\theta, \mathbb{X}, \mathbb{Y})$ represent the loss function that is optimized during training. During prediction time, let \mathbb{Y}_{pred} be the predicted labeling. In order to generate an error gradient for backpropagation, we create a pseudo ground truth labeling \mathbb{Y}_{pseudo} by modifying \mathbb{Y}_{pred} as follows: We initialize \mathbb{Y}_{pseudo} with \mathbb{Y}_{pred} . Let the k^{th} component of \mathbb{Y}_{pseudo} be represented as $\mathbf{y}_k = [y_{k_1}, \dots, y_{k_{C_{out}}}]$, which is a C_{out} -dimensional score vector. We modify \mathbf{y}_k to be a 1-hot encoded vector with the maximally confident class set to 1 and others to zero. Then, the error gradient signal is computed based on the loss function $\mathcal{J}(\theta, \mathbb{X}, \mathbb{Y}_{pseudo})$ and backpropagated through the network up to the input. Let the backpropagated error gradient signal at the input be represented as: $\nabla_X \mathcal{J}(\theta, \mathbb{X}, \mathbb{Y}_{pseudo})$.

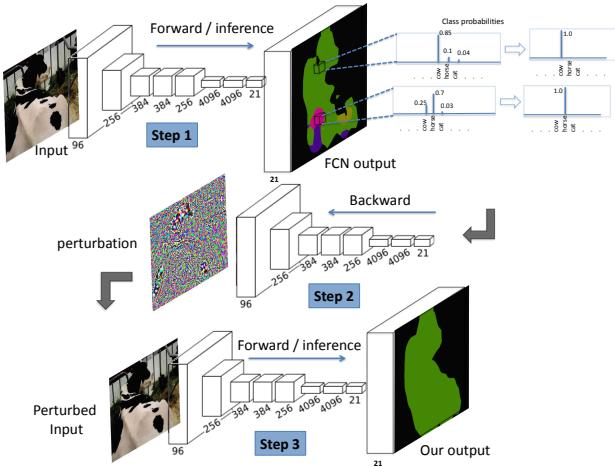


Figure 2: Processing pipeline for the proposed approach for semantic segmentation

The perturbed input is then generated as follows:

$$\mathbb{X}_{per} = \mathbb{X} + \epsilon \cdot sign(\nabla_X \mathcal{J}(\theta, \mathbb{X}, \mathbb{Y}_{pseudo})), \quad \epsilon > 0 \quad (1)$$

where ϵ is a non negative scaling factor that is model dependent and $sign(\cdot)$ represents the signum function taken elementwise. \mathbb{X}_{per} is then fed into the network for a forward pass to generate the final output.

It can be argued that the above method of generating gradients using the network's prediction can lead to inaccurate gradient information propagated through the network especially in cases where the network's output contains many misclassified pixels. The key insight we provide in this work is that despite misclassifications by the deep network, the gradients at the input obtained from the network's prediction, in general, improve the final output of the deep network.

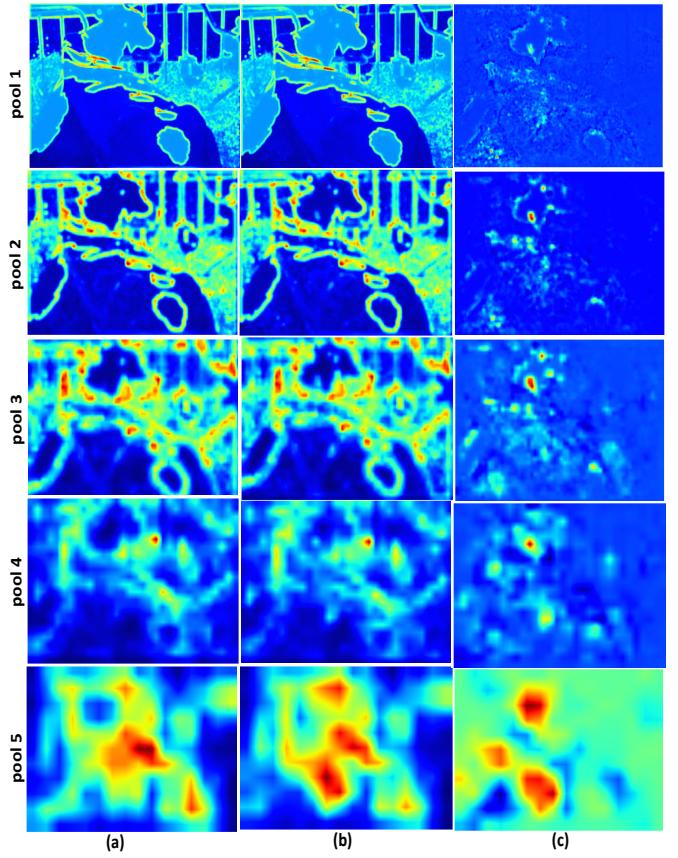


Figure 3: Visualization of filter responses showing how the correct context is propagated along the FCN-32s network. Column (a): filter responses during the forward pass using the original input. Column (b): filter responses during the forward pass using the perturbed input. Column (c): difference between (a) and (b)

3.1.1 Understanding Guided Perturbations

In this section, we perform several experiments to provide insight into different aspects of guided perturbations. Please refer to Figure 2 for the steps (Step 1, Step 2, Step 3) mentioned in this section.

Impact of perturbations on filter responses: To get a clear understanding of what happens during the forward pass in Step 3 that vastly changes the network’s prediction, we visualize the filter responses for the FCN-32s network in Figure 3. This model was chosen due to its simpler architecture but we observed similar behavior in other deep architectures too. In Figure 3, we plot the average filter responses at different layers through the deep network after upsampling them bilinearly to image size. As can be observed, the influence of the added perturbations are not visually explicit until the *pool5* layer but the difference of the filter responses in Column (c) indicate that the information propagates from layers as early as *pool2*.

Next, we analyze the pixels for which the network predictions changed from Step 1 to Step 3. Figure 4(c) shows the pixels that were classified wrongly during the forward pass in Step 1 but were correctly classified at the final output. On the other hand, 4(d) shows the pixels that were correctly classified in Step 1 but were incorrect at the final output. Observe that, the correctly classified pixels between Step 1 and Step 3 are mostly internal to the image where additional contextual information is available for the network to switch its prediction whereas the small number of misclassified pixels are largely concentrated along the boundary regions of the image where the context is ambiguous. We present more of such visualization examples in the Appendix (6).

Approximating ideal gradient direction: In this experiment, we would like to answer the question: what are the ideal perturbations that can be generated at the input? The best one can do is to use the ground truth to generate error gradients at the softmax layer which is then backpropagated to generate the perturbed input. When this perturbed input is fed back to the network, the result is a vastly improved prediction as shown in Figure 5 (c). While perturbations from ground truth significantly improve the performance, this information is not available during prediction time. The



Figure 4: (a) Output of FCN-32s network (b) Output from the proposed approach (c) Pixels that were incorrectly classified by FCN-32s corrected by our approach (d) Pixels that were incorrectly classified by our approach that FCN-32s classified correctly.

novelty of the current work is that the ground truth gradient direction is being approximated well enough by the predicted gradient directions that are computed using only the network’s prediction.

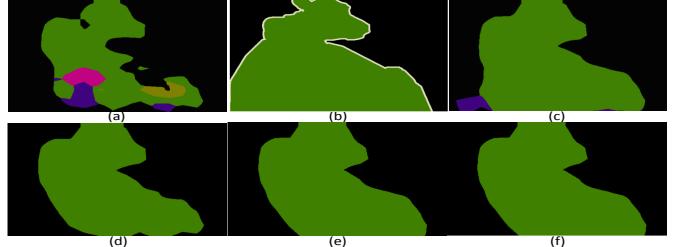


Figure 5: (a) Output of FCN-32s (b)Ground truth labeling (c) Output of perturbed input using ground truth gradients (d)-(f) Output of perturbed input using guided perturbations for iteration 1, 2 and 3 respectively.

To understand the extent of usefulness of the predicted gradients, we performed an experiment where the three steps outlined in our approach (Figure 2) is applied over successive iterations Figure 5 (d) shows the output of our approach obtained in the first iteration and Figure 5 (e)-(f) show the output over successive iterations. This shows that the most significant improvement happens at the first iteration and the subsequent iterations yield little improvement. We observe similar behavior on average over the PASCAL VOC2012 validation set.

Intuition based on overlapping receptive fields: In a CNN, the receptive fields of neighboring pixels define a context for their interactions. The advantage of having overlapping receptive fields is that the neighborhood connectivity is established automatically without explicitly specifying it. As long as the errors made by the CNN are sparse with respect to each pixel’s receptive fields, the error gradients when accumulated over the entire network and used to perturb the input image exhibit a corrective behavior. The effect of GP can be seen as a type of residual information that is propagated through the network which results in contextual smoothing. This is evident by looking at the filter responses in Figure 3, more specifically in Column 3, which shows the difference in responses with and without GP. It can be observed from the *pool5* responses that the peak activations occur around neighborhoods where there are competing classes. GP perturb these neighborhoods the most thus resulting in contextually smooth predictions in those regions.

Analyzing GP in depth: In figure 6, we show how guided perturbations impact the decisions made by the deep network by considering a local region in the input image and tracking its classification scores at *score-fr* layer (before upsampling layer) across different values of ϵ . In the top half of figure 6a, the patch of interest in the RGB image is marked by a red box and its corresponding region in the

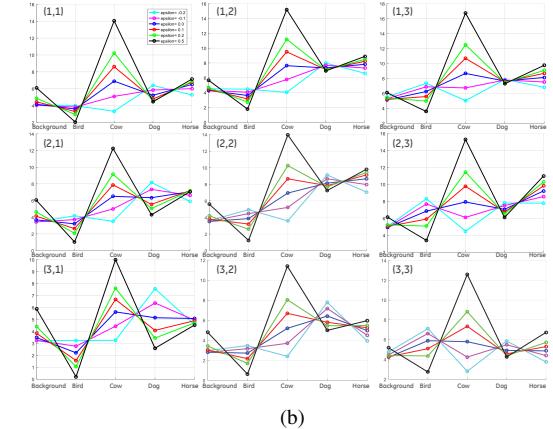
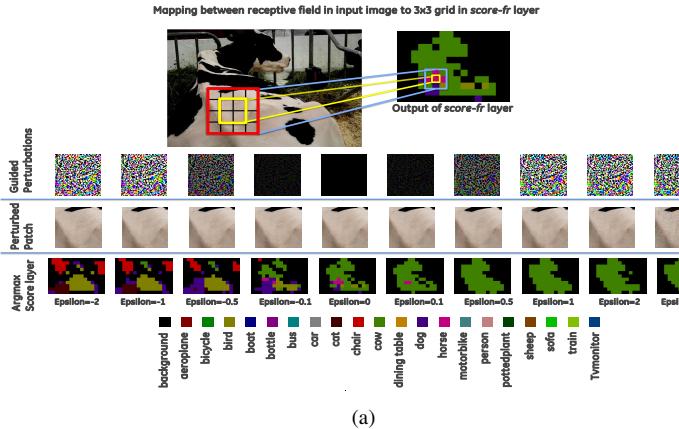


Figure 6: (a) The top half shows an RGB patch in the input image and its corresponding patch in the *score-fr* layer output of the FCN-32s network, before upsampling to image size. In the bottom half, we show, for different values of ϵ the guided perturbations, the perturbed RGB patches and the *score-fr* output. Notice how the the scores become contextually smoother for $\epsilon > 0$. (b) The actual score values of the top-5 predicted classes for the 3x3 grid marked in blue in figure (a) are plotted. Observe that for a range of positive values of ϵ , the correct class score (*cow*) dominates the others across the entire neighborhood. The legend in (1,1) applies to all the plots. Best viewed in screen. Please zoom for clarity.

score-fr output is marked by a blue box. Immediately below, the following are shown for different values of ϵ : (1) Guided perturbations generated at the input (2) perturbed RGB patches (3) output of the *score-fr* layer. Important observations that can be made from figure 6a are:

- Input perturbations corresponding to positive ϵ improve the score output over a vast range of values. This visualization shows how guided perturbations are able to operate at a local level by leveraging neighborhood contextual information as can be directly observed from the images of score layer shown in the bottom row.
- Even a small negative value of ϵ results in a large adverse effect on the score output, without any perceptible change in the perturbed RGB patch. This shows that a negative ϵ corresponds to an adversarial setting.

This discussion motivates our choice of using $\epsilon > 0$ to generate GP. To further analyze how these input perturbations affect the actual classifier score, we show in figure 6b, the predictions of the deep network for the 3x3 grid in the *score-fr* output from figure 6a, for different values of ϵ . For clarity, we only show the predicted scores of the top-5 classes. From the score values of the grid position (2,2), we can observe that as ϵ increases, the score of true class (*cow*) keeps increasing while the scores of the confusing classes do not vary much. The other plots show that this trend is observed across the entire neighborhood of the 3x3 grid. Thus, it can be inferred that perturbations at the input affect the decision of the deep network in a contextually consistent manner. We again observe that the score of the true class drops significantly even for a small negative ϵ which is consistent with our earlier observations.

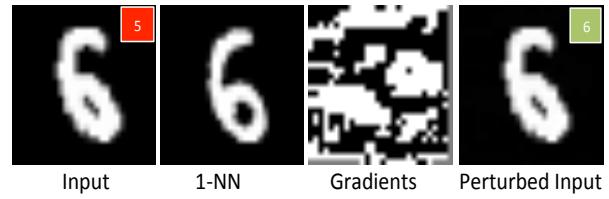


Figure 7: The input image is classified as ‘5’. By perturbing the input from the gradients generated using the class of top nearest neighbor the network changes its prediction to ‘6’

3.2. Image Classification

The method described in Section 3.1 for semantic segmentation cannot be applied directly for classification tasks. In fact, as shown by the work on adversarial examples [4], gradients in directions other than the ground truth severely affects the network’s performance. Then the question arises, what constitutes the context of a test image for the classification task. A natural choice is the nearest neighbors in the learned feature space. To test this hypothesis, we perform the following experiment: Given an input image, we first extract the feature from the deep network and use it to select top k nearest neighbors from the training set using euclidean distance metric. We then perturb the test image with the weighted average of gradients generated using the class of the i^{th} nearest neighbor. Following the notation established in Section 3.1, the equation for perturbed image is given as follows: $\mathbb{X}_{per} = \mathbb{X} + \epsilon \sum_{i=1}^k (w_i \text{sign}(\nabla_X \mathcal{J}(\theta, \mathbb{X}_i, \mathbb{Y}_{nn_i})))$, where \mathbb{X}_i is the i^{th} nearest neighbor; k is the number of nearest neighbors and w_i is weight associated with each

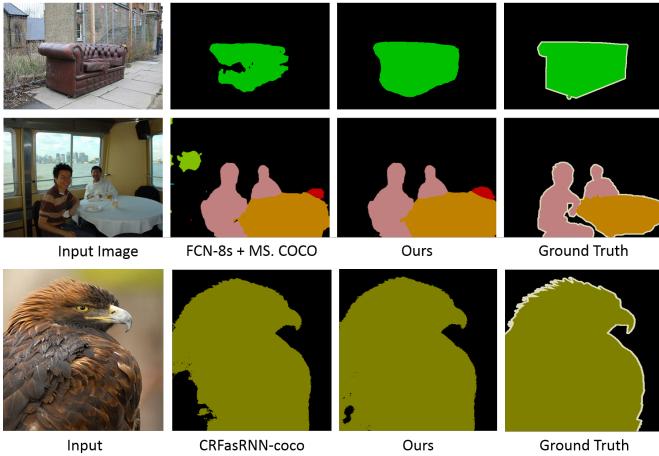


Figure 8: Qualitative results on the PASCAL VOC2012 reduced validation set. In the top two rows, we compare our result with the FCN-8s part of CRFasRNN that has been trained on MS.COCO dataset [9] and publicly released by [20]. In the bottom row, we compare with the complete CRFasRNN framework[20]. More results can be found in the Appendix(6).

nearest neighbor i and $\mathcal{J}(.)$ corresponds to the loss function. Figure 7 shows an example where the network correctly classifies the perturbed input generated using this procedure.

4. Experiments

In this section, we perform several experiments showing how our approach could be seamlessly applied on top of several pretrained deep networks. We test our method on the semantic segmentation task on PASCAL VOC2012 dataset [3], scene labeling task on the PASCAL Context 59-class dataset [11] and classification tasks on the MNIST and CIFAR10 datasets [7]. These results support how our approach is able to generalize across different types of problems in computer vision and highlights the advantage that it can be used with any pretrained model.

4.1. Evaluation Metrics

We evaluate our approach using the mean Intersection over Union (mIoU) metric commonly used for semantic segmentation as reported in [14]. Let n_{ij} be the number of pixels of class i predicted to belong to class j , N_{cl} be number of classes, and $t_i = \sum_j n_{ij}$ be the total number of pixels of class i . It is then formulated as $\text{mean IoU} = \frac{1}{N_{cl}} \sum_i \frac{n_{ii}}{t_i + \sum_j n_{ji} - n_{ii}}$. For MNIST and CIFAR-10, we use classification accuracy as a metric to compare against the baseline.

Table 1: Results on the reduced VOC2012 validation set with 346 images. '-coco' denotes that the model was trained on MS COCO data in addition to the SBD dataset. Numbers in brackets show the magnitude of change compared to the corresponding base models.

Method	mean IU
FCN-32s	62.10
FCN-32s + GP	64.71(+2.61)
FCN-8s	63.97
FCN-8s + GP	66.97(+3.0)
FCN-8s-coco	69.85
FCN-8s-coco + GP	71.99(+2.14)
CRFasRNN-coco	72.95
CRFasRNN-coco + GP	73.75(+0.8)

4.2. Semantic Segmentation

We use PASCAL VOC2012 dataset for evaluating our approach for semantic segmentation task. It consists of 21 classes including background. We use the following pretrained models as baselines and show the improvement that can be obtained using our approach for each of them:(1) FCN-32s and FCN-8s [14]: these models are trained using the SBD dataset[5] that consists of 9,600 images. (2) FCN-8s-coco and CRFasRNN [20]: these are trained using the images from MS COCO[8] and the SBD dataset using a total of 77,784 images.

For all these methods, we use the publicly available models at the time of submission. We use a single NVIDIA TitanX GPU for our experiments and CAFFE library[6] for implementation. The pretrained models used in this section are obtained from the CAFFE Model Zoo [1] at the time of submission. All the reported results are computed with 1 iteration of our approach unless mentioned otherwise. Table 1 shows the results of applying the proposed approach to the different pretrained models during prediction time over a reduced validation set of 346 images as done in [20]. As can be observed, the proposed approach results in increased performance over all the listed pretrained models. This reiterates the fact that our approach is indeed architecture independent and can be easily integrated even with complex feedforward architectures like CRFasRNN. Table 2 shows the evaluation of our approach on PASCAL VOC2012 *test set* using FCN-8s pretrained network as the base model to demonstrate the improvement shown by our method in an unbiased setting. The ϵ value used for the test set was tuned on the validation set.

4.3. Scene Labeling

The scene labeling task is a dense pixel labeling task that is evaluated on the PASCAL Context dataset. While there are more than 400 classes defined, the challenge entails evaluating on the 59 classes that are specified as most frequent [11]. The labeled classes contain scene elements

Table 2: Results on the PASCAL VOC2012 test set consisting of 1456 images using FCN-8s as the base network.¹ Use of Guided Perturbations improves the performance of the base network on 19 out of 21 classes.

Method	bkg	aero	bicycle	bird	boat	bottle	bus	car	cat	chair	cow	table	dog	horse	mbk	person	plant	sheep	sofa	train	tv	mean
FCN-8s [14]	92.0	82.4	36.1	75.6	61.4	65.4	83.3	77.2	80.1	27.9	66.8	51.5	73.6	71.9	78.9	77.1	55.3	73.4	44.3	74.0	63.2	67.2
FCN-8s + GP	92.4	84.4	35.9	79.3	62.6	70.5	86.2	80.0	82.8	28.0	71.9	55.2	74.6	75.6	80.2	77.4	56.9	75.6	45.8	77.4	63.18	69.3

Table 3: Results on the PASCAL-Context 59-classes validation set.

Method	mean IU
FCN-8s	39.12
FCN-8s + GP	40.44

in addition to objects that appear in the PASCAL VOC segmentation challenge, making this a much harder benchmark. To evaluate our approach on this task, we use the FCN-8s model from [14] as our baseline that was trained on the standard training split of 10,000 images provided with the dataset. The results, which were generated on the validation set consisting of 5105 images are shown in Table 3. We improve the performance of the FCN-8s network by 1.3% which is significant given the large size of the validation set. Please note that the ϵ value was not tuned to fit this dataset rather the best performing ϵ from Table 1 was used.

4.4. Classification

As explained in Section 3.2, our approach can be used with pretrained classification networks as well using the perturbations generated from the nearest neighbors. To evaluate the performance, we tested the method on two standard classification datasets: MNIST and CIFAR10. MNIST consists of grayscale images of digits while CIFAR10 consists of more realistic images of object classes. We follow the standard training/testing split for both the cases with 50,000 images used for training the model and 10,000 images used for testing. We use 3 nearest neighbor with equal weights for all our experiments. For MNIST, we use a CNN with 2 conv. layers and 2 fully-connected layers with a 20-50-500-10 architecture and for CIFAR10, we use a CNN with 5 conv. and 2 fully-connected layers with a 64-64-128-128-128-10 architecture.

Table 4: Results on the classification task on MNIST and CIFAR10 datasets.

Dataset	Baseline	Proposed
MNIST	98.92	99.15
CIFAR10	76.31	76.95

Table 4 shows the results of our classification experiments. GP improves performance over the baseline on both the datasets. However, the improvement in performance is not as high as in the segmentation case which could be attributed to two reasons: (1) the base networks themselves have learned a very strong representation and (2) the context information in the classification task is relatively weak

compared to the segmentation task.

5. Ablative Experiments

For all the experiments in this section, we use FCN-32s network and the validation set used in section 4.2.

Speed-Performance trade-off The guided perturbations generated at the input layer of a deep network improves the performance of the base model. However, there is a computational overhead due to performing an additional backward and forward pass. As an alternative, the backward pass could be performed up to an intermediate layer in the deep network instead of the input layer. In this section, we provide results addressing the trade off between computational time and resulting performance due to perturbing layers other than the input.

Table 5: Trade-off between performance and computation times obtained by truncating guided perturbations over different layers across the deep network. Original time taken is 0.12s per image. The baseline performance is 62.1%

layer	input	pool2	pool3	pool4
Time	0.33s	0.27s	0.24s	0.22s
mIOU	64.71	64.61	64.55	64.3

It can be observed from Table 5 that even using the perturbed input from as late as *pool4* layer the improvement in performance remains almost constant while computation time drops significantly. This experiment shows that effect of GP is not only observed at the input but also in the intermediate layers of the deep network and hence can be leverages for reducing the computational cost.

Guided Perturbations (vs) other strategies In this section, we perform an ablative experiment where we perturb the input image in different ways in order to distinguish them from Guided Perturbations and show that the GP yields the most improvement in performance. As explained in Section 3.1, to generate a guided perturbation, we replace the softmax output with a one-hot encoded vector for the class of maximum confidence. We consider different methods to modify the label distribution that is obtained from the softmax function as follows:

- *random-onehot*: The class label is chosen in an uniformly random manner and used as ground truth instead of the maximum probability class.

¹Anonymous link: <http://host.robots.ox.ac.uk:8080/anonymous/8LV1YS.html>

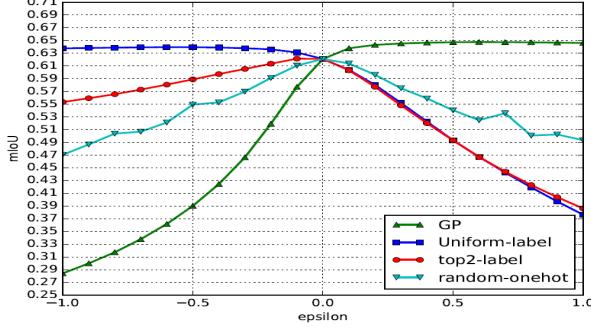


Figure 9: Mean IOU values for several perturbations generated by using different types of label distributions on the validation set over the range $\epsilon = [-1, 1]$ with FCN-32s as the base network. Please refer to section 5 for details.

- *Uniform-label*: An uniform label distribution is produced by assigning equal probability to all the classes and used as encoding to generate the error gradient.
- *top2-label*: Modified label distribution contains equal probability to top two predicted classes and used as encoding to generate the error gradient.

Figure 9 shows the effect of different types of label distribution on the segmentation performance. At the outset, it can be observed that GP gives the best quantitative performance of 64.7% compared to the second best case, which is the uniform setting with negative ϵ which scores 63.8%. We can also observe that when we perform GP, $\epsilon < 0$ corresponds to the adversarial setting. Intuitively, this setting is equivalent to maximizing the loss of the softmax classification function during training. Hence, the backpropagated gradient always moves away from the correct class. In our approach, GP is always generated by setting $\epsilon > 0$ as mentioned in sections 3.1 and 3.2. The setting involving choosing a random label to generate the one-hot vector at the softmax output results in poor performance across all values of ϵ since gradient directions become random and the resulting perturbations adversely affect the performance of the deep network on the perturbed input image.

The interesting case to analyze from Figure 9 is the performance of the *Uniform-label* setting for $\epsilon < 0$. To understand this effect, Figure 10 illustrates a toy example showing the difference between the error gradients generated using GP and *Uniform-label* setting for a different possible output score distributions from the CNN. In this toy example, the CNN is trained to classify among 5 classes. Observe that, for the unimodal case, the gradient signal generated for a uniform output label distribution has the same relative magnitude as the gradient signal generated for GP but the dominant gradient direction is exactly the opposite. This offers an explanation as to why the performance of the *Uniform-label* setting for $\epsilon < 0$ tends towards GP. However, GP still gives a better performance compared to the uniform label distribution and this can be understood by looking at

bottom half of the figure 10. In this case, the score vector is bimodal and hence there are two dominant directions in the gradient signal. Notice that the top gradient direction in the case of GP still points towards the correct class and all other directions move away from the correct class, as expected. But in the case of uniform label distribution, there are two competing directions and hence there is higher probability in this case for the gradient to move in the wrong direction.

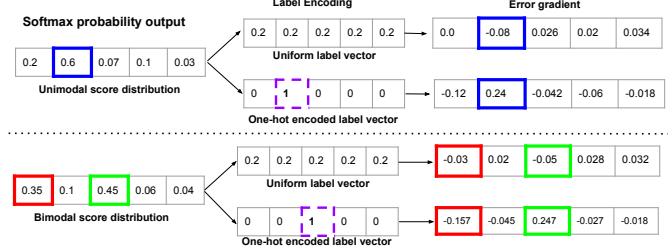


Figure 10: Difference in the gradient signal generated between *Uniform-label* setting and GP for the case of unimodal output score distribution (top) and bimodal output score distribution (bottom). The dominant gradient direction in both cases is shown in the colored boxes. The exact derivation for computing these gradient values is given in the Appendix (6).

We also performed an experiment where the input image was perturbed directly without backpropagation using random gaussian noise of different standard deviations. We did not observe any improvement in performance compared to the base network and performance dropped for large values of standard deviation.

Effect of Scaling parameter We evaluate the performance of our approach using FCN-32s and FCN-8s networks over a range of scaling parameter ϵ on the validation set. Figure 11 shows how the performance varies based on the scaling factor. It can be observed that improvement in performance is generally obtained over a wide range of values of ϵ . This indicates that network’s behavior is not very sensitive to the value of ϵ though there seems to be an optimal value for best performance that depends on the deep model. We use $\epsilon = 0.55$ for FCN-32s, $\epsilon = 0.7$ for FCN-8s network and $\epsilon = 0.22$ for CRFasRNN network for our experiments.

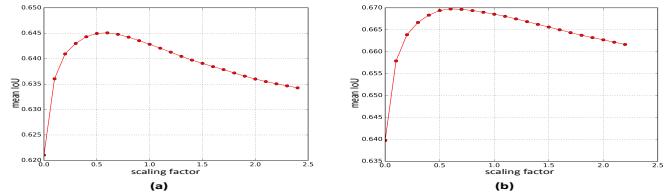


Figure 11: Effect of scaling factor ϵ on performance of FCN-32s (left) and FCN-8s (right) networks evaluated on the reduced PASCAL VOC2012 validation set. Best viewed in screen. Please zoom for clarity.

6. Conclusion

In this paper, we have shown novel self-corrective behavior of CNNs for segmentation and classification tasks. We showed that guided perturbations can improve the network’s performance without additional training or network modification. We have demonstrated this effect on several publicly available datasets and using different network architectures. We have presented several experiments that try to understand and explain different aspects of guided perturbations. We believe that this behavior can lead to novel network designs and better end-to-end training procedures. Another interesting research direction would be to model these perturbations using a generative framework and explore their universality across different tasks.

References

- [1] Caffe model zoo. <https://github.com/BVLC/caffe/wiki/Model-Zoo>. Accessed: 2010-09-30. [6](#)
- [2] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, and A. L. Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *arXiv preprint arXiv:1606.00915*, 2016. [2](#)
- [3] M. Everingham, S. A. Eslami, L. Van Gool, C. K. Williams, J. Winn, and A. Zisserman. The pascal visual object classes challenge: A retrospective. *International Journal of Computer Vision*, 111(1):98–136, 2015. [1, 6](#)
- [4] I. J. Goodfellow, J. Shlens, and C. Szegedy. Explaining and harnessing adversarial examples. *arXiv preprint arXiv:1412.6572*, 2014. [1, 2, 5](#)
- [5] B. Hariharan, P. Arbeláez, L. Bourdev, S. Maji, and J. Malik. Semantic contours from inverse detectors. In *2011 International Conference on Computer Vision*, pages 991–998. IEEE, 2011. [6](#)
- [6] Y. Jia, E. Shelhamer, J. Donahue, S. Karayev, J. Long, R. Girshick, S. Guadarrama, and T. Darrell. Caffe: Convolutional architecture for fast feature embedding. In *Proceedings of the 22Nd ACM International Conference on Multimedia, MM ’14*, pages 675–678. ACM, 2014. [6](#)
- [7] A. Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009. [6](#)
- [8] T. Lin, M. Maire, S. J. Belongie, L. D. Bourdev, R. B. Girshick, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft COCO: common objects in context. *CoRR*, abs/1405.0312, 2014. [6](#)
- [9] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick. Microsoft coco: Common objects in context. In *European Conference on Computer Vision*, pages 740–755. Springer, 2014. [6](#)
- [10] A. Mahendran and A. Vedaldi. Understanding deep image representations by inverting them. In *2015 IEEE conference on computer vision and pattern recognition (CVPR)*, pages 5188–5196. IEEE, 2015. [2](#)
- [11] R. Mottaghi, X. Chen, X. Liu, N.-G. Cho, S.-W. Lee, S. Fidler, R. Urtasun, and A. Yuille. The role of context for object detection and semantic segmentation in the wild. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 891–898, 2014. [6](#)
- [12] A. Nguyen, J. Yosinski, and J. Clune. Deep neural networks are easily fooled: High confidence predictions for unrecognizable images. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 427–436. IEEE, 2015. [1, 2](#)
- [13] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein, A. C. Berg, and L. Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3):211–252, 2015. [1](#)
- [14] E. Shelhamer, J. Long, and T. Darrell. Fully convolutional networks for semantic segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, PP(99):1–1, 2016. [2, 6, 7, 10](#)
- [15] C. Szegedy, W. Zaremba, I. Sutskever, J. Bruna, D. Erhan, I. Goodfellow, and R. Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013. [1, 2](#)
- [16] R. Vemulapalli, O. Tuzel, M.-Y. Liu, and R. Chellappa. Gaussian conditional random field network for semantic segmentation. *CVPR*, 2016. [2](#)
- [17] J. Yosinski, J. Clune, A. Nguyen, T. Fuchs, and H. Lipson. Understanding neural networks through deep visualization. *arXiv preprint arXiv:1506.06579*, 2015. [2](#)
- [18] F. Yu and V. Koltun. Multi-scale context aggregation by dilated convolutions. *arXiv preprint arXiv:1511.07122*, 2015. [2](#)
- [19] M. D. Zeiler, D. Krishnan, G. W. Taylor, and R. Fergus. Deconvolutional networks. In *Computer Vision and Pattern Recognition (CVPR), 2010 IEEE Conference on*, pages 2528–2535. IEEE, 2010. [2](#)
- [20] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1529–1537, 2015. [2, 6, 10](#)

Appendix

Appendix contains additional material and examples to support our paper. The explicit formula used in the error gradient computation for the toy example in Figure 10 from the paper is derived in section 7. Figures 12, 13 and 14 show additional examples of improved performance when the proposed approach is used with the FCN-8s [14], FCN8s-coco [20] and CR-FasRNN [20] pretrained models respectively. Figure 15 shows examples to supplement the claim made in Figure 4 in the paper that the pixels that are predicted correctly by our approach are more internal to the image whereas the small number of pixels that are predicted wrongly tend to occur towards the boundaries. These examples are generated using the FCN-32s deep network [14]. Finally, Figure 16 shows additional results of using our approach for the MNIST classification task.

7. Error gradient computation for Figure 10

Let the score output of the deep network be: $\mathbf{z} \in \mathbb{R}^{N_c}$. To get a probability distribution over classes, this is passed through a softmax operator whose output is given as: $\mathbf{y} = \left\{ \frac{e^{z_i}}{\sum_i e^{z_i}} \right\}_{i=1}^{N_c}$, where N_c is the number of classes. If $k \in [1, N_c]$ is the correct class, then the error gradient computed at the softmax output with respect to its input \mathbf{z} is given as follows: Let \sum_C denote $\{\sum_{i=1}^{N_c} e^{z_i}\}$, then

$$\text{if } i = k : \quad \frac{\partial y_i}{\partial z_i} = \frac{\sum_C e^{z_i} - e^{z_i} e^{z_i}}{\sum_C^2} = \frac{e^{z_i}}{\sum_C} \left(1 - \frac{e^{z_i}}{\sum_C} \right) = y_i(1 - y_i) = y_k(1 - y_i) \quad (2)$$

$$\text{if } i \neq k : \quad \frac{\partial y_i}{\partial z_k} = -\frac{0 - e^{z_i} e^{z_k}}{\sum_C^2} = -\frac{e^{z_i}}{\sum_C} \frac{e^{z_k}}{\sum_C} = -y_i y_k = y_k(0 - y_i) \quad (3)$$

(2)-(3) could be summarized in the following single equation:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = y_k(\ell - \mathbf{y}) \quad (4)$$

where $\ell \in \mathbb{R}^{N_c}$ is the label distribution, which in this case is a one hot vector with $\ell_k = 1$ and others zero. For a more general case, where ℓ defines a distribution among classes, this formula generalizes in a straight forward manner as follows:

$$\frac{\partial \mathbf{y}}{\partial \mathbf{z}} = (\ell \cdot \mathbf{y})(\ell - \mathbf{y}) \quad (5)$$

It can be observed that (5) is a general version of (4) since the maximum probability value y_k is replaced by the dot product between the label distribution and the output of softmax operation.

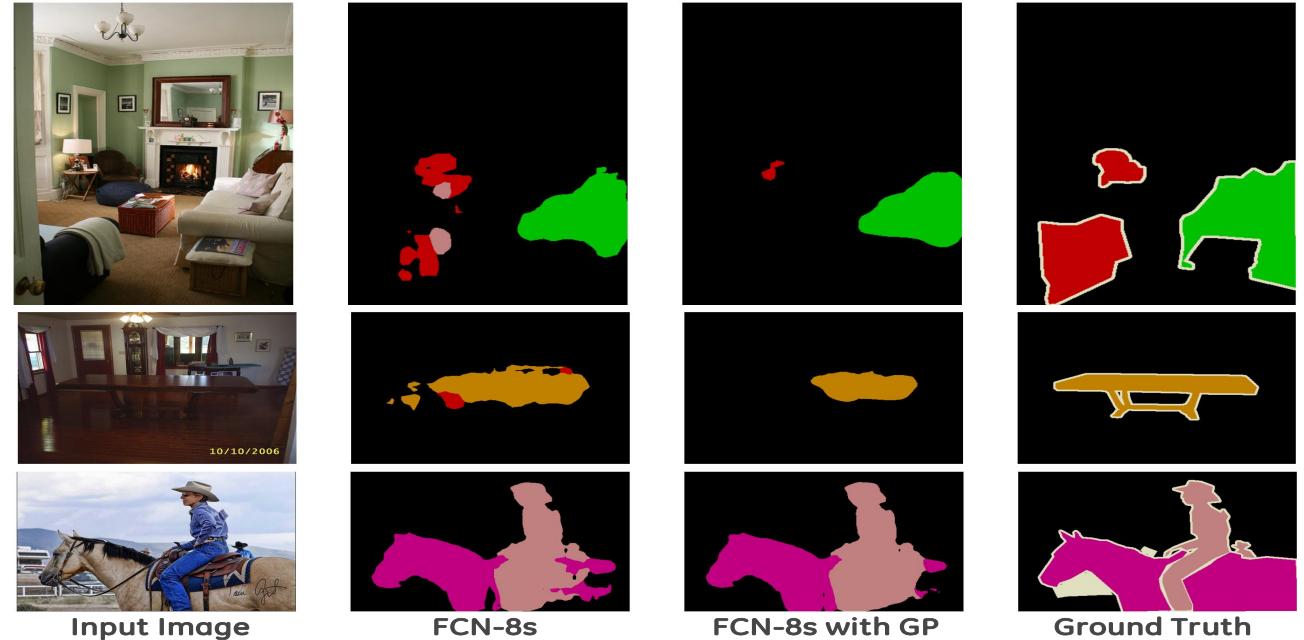
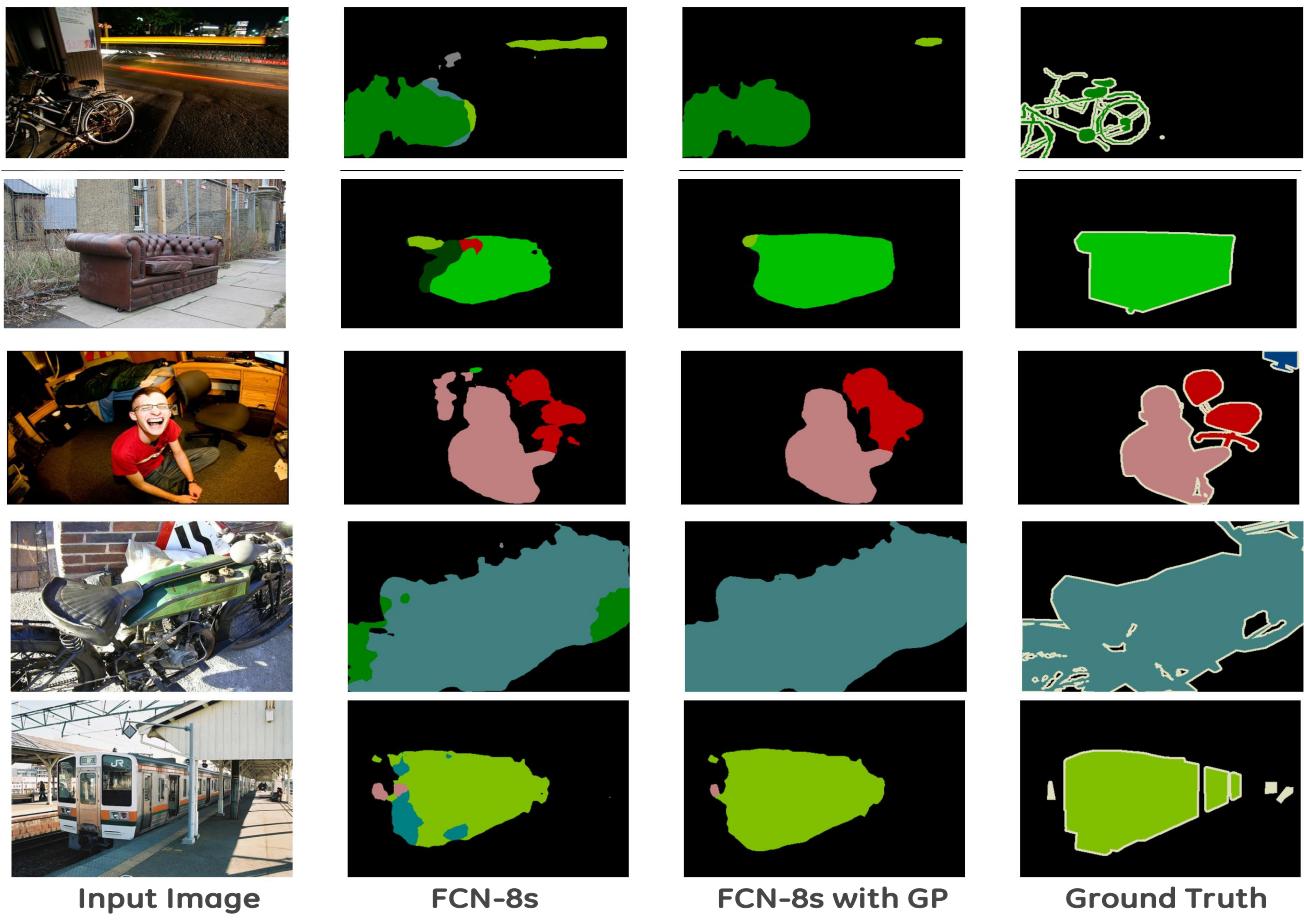
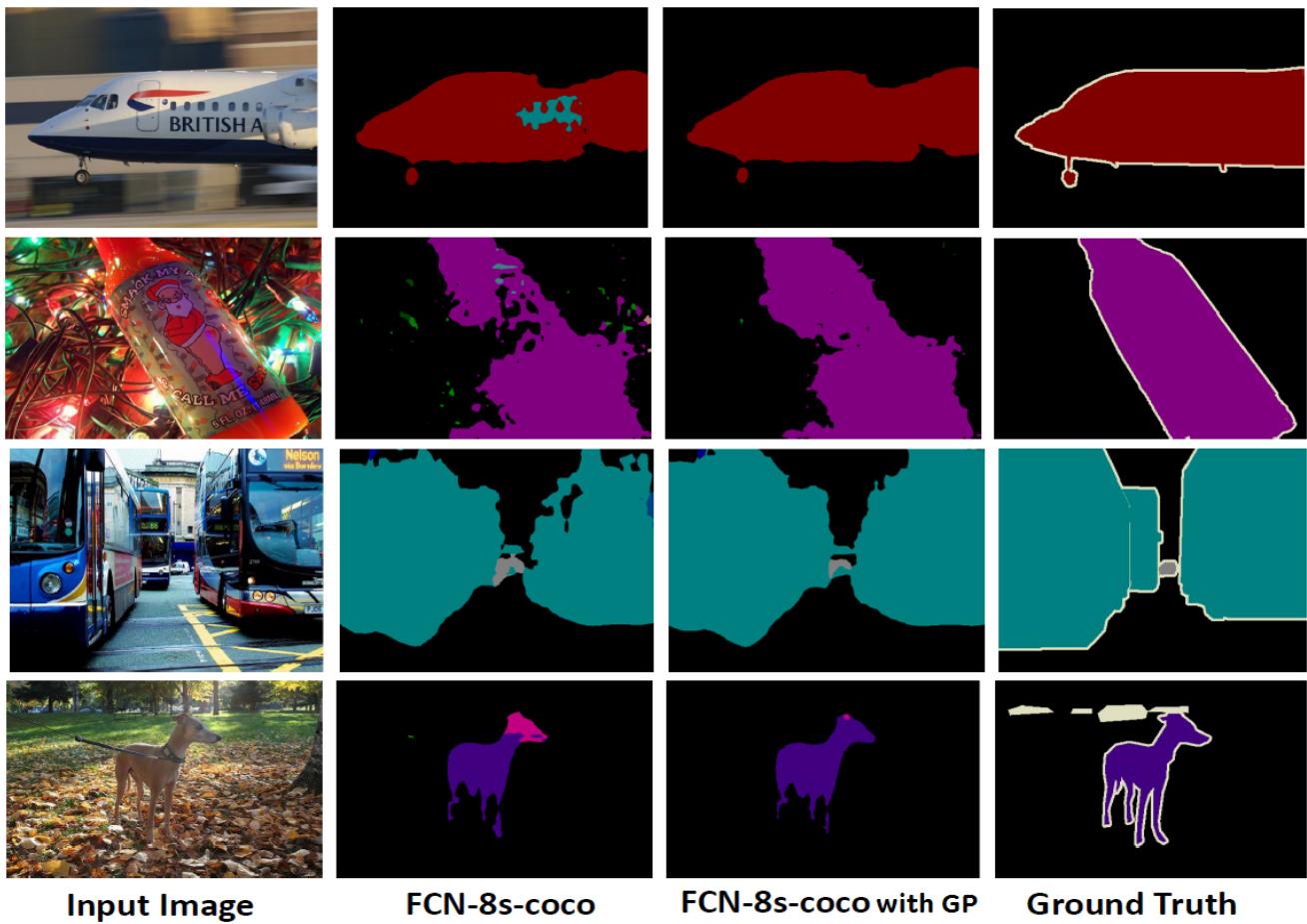


Figure 12: Qualitative results on the PASCAL VOC2012 reduced validation set - Comparison with FCN-8s pretrained model. Top half shows the successful outputs, Bottom half shows the failure cases.

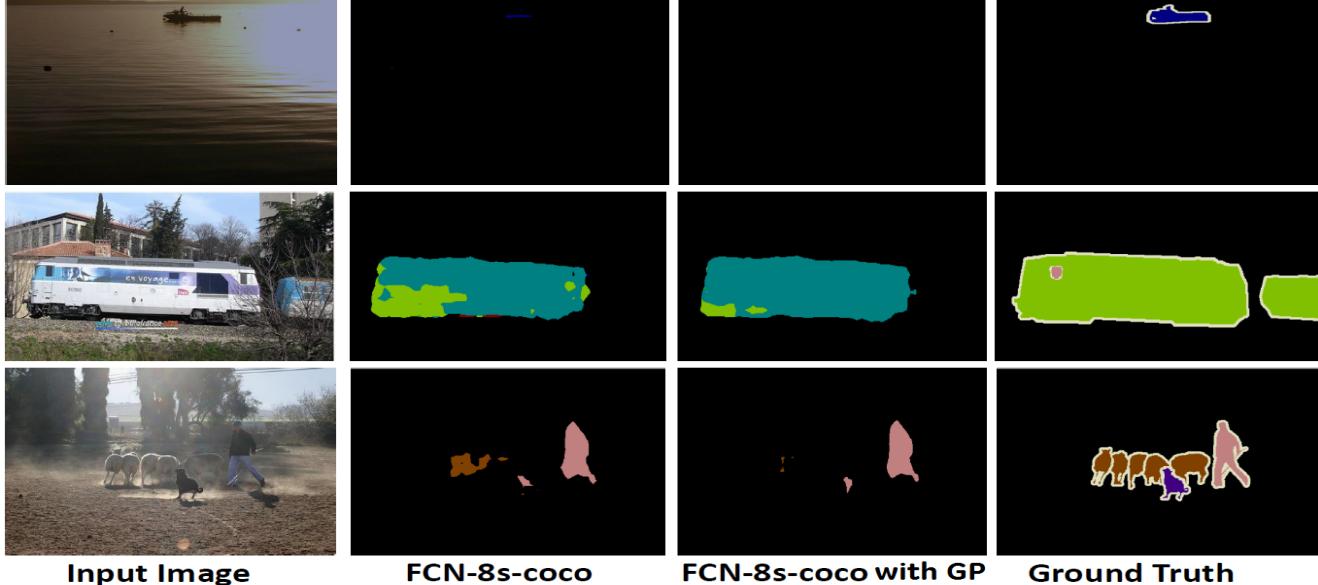


Input Image

FCN-8s-coco

FCN-8s-coco with GP

Ground Truth



Input Image

FCN-8s-coco

FCN-8s-coco with GP

Ground Truth

Figure 13: Qualitative results on the PASCAL VOC2012 reduced validation set - Comparison with FCN-8s-coco pretrained model. Top half shows the successful outputs, Bottom half shows the failure cases.

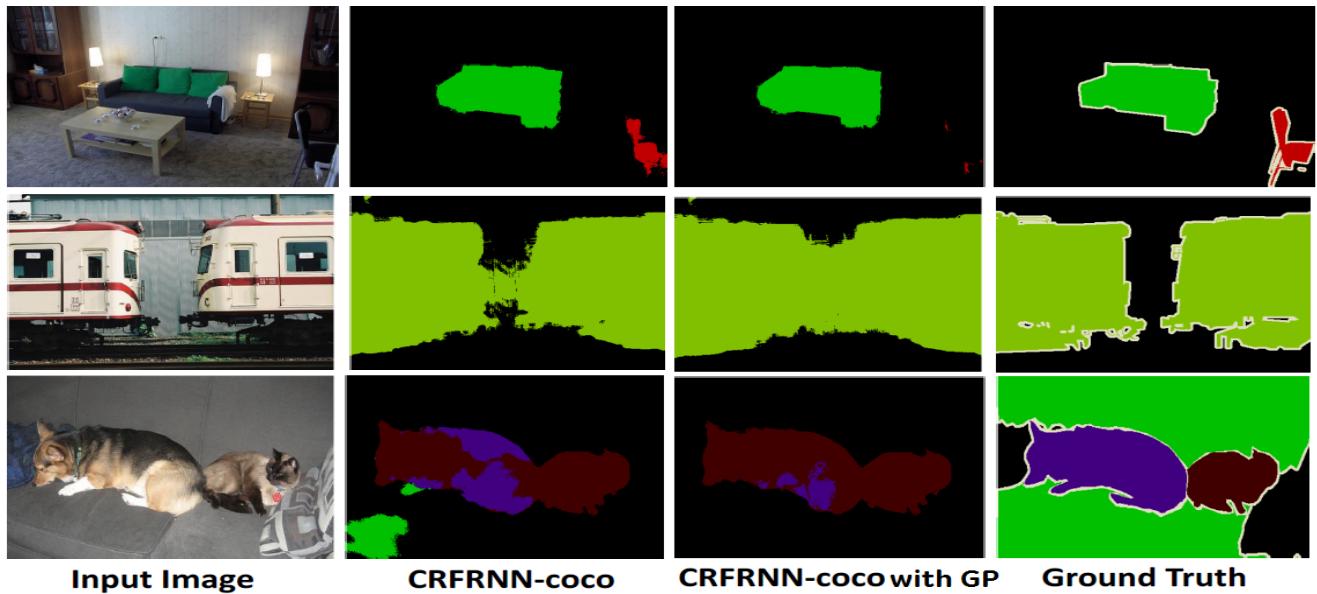
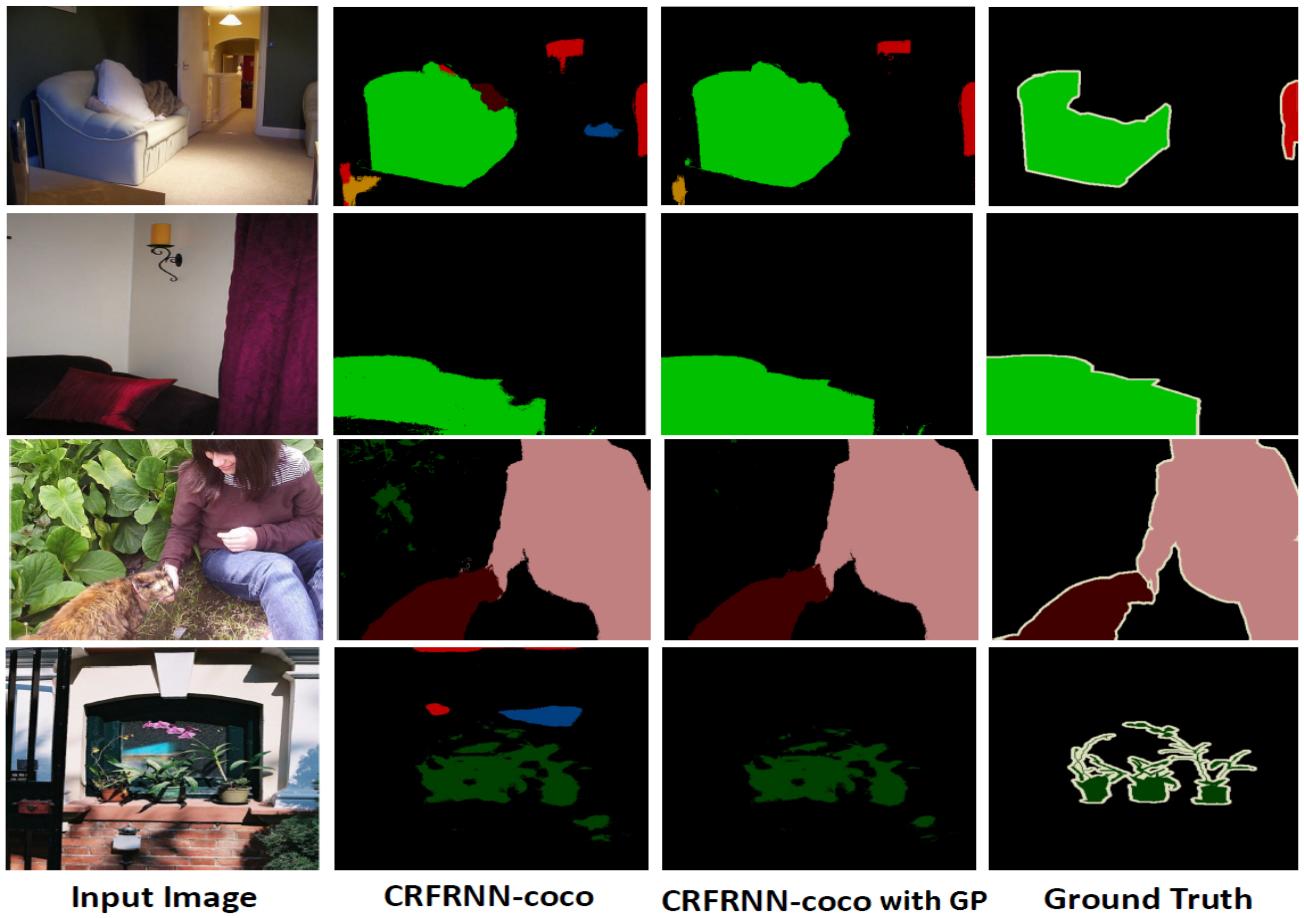


Figure 14: Qualitative results on the PASCAL VOC2012 reduced validation set - Comparison with CRFRNN-coco pretrained model. Top half shows the successful outputs, Bottom half shows the failure cases.

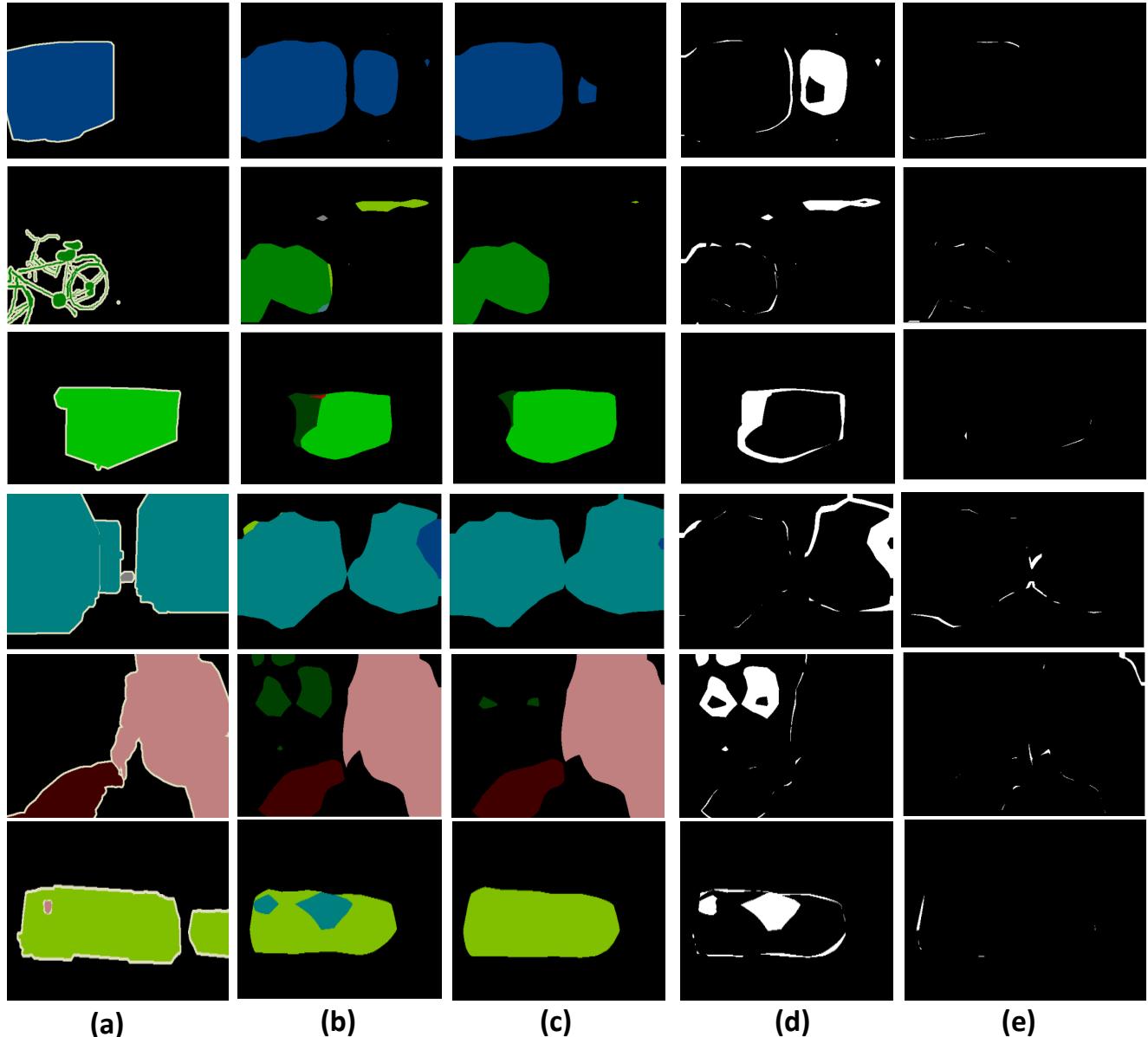


Figure 15: (a) Ground truth (b) Output of FCN-32s network (c) Output from the proposed approach (d) Pixels that were incorrectly classified by FCN-32s corrected by our approach (e) Pixels that were incorrectly classified by our approach that FCN-32s classified correctly.

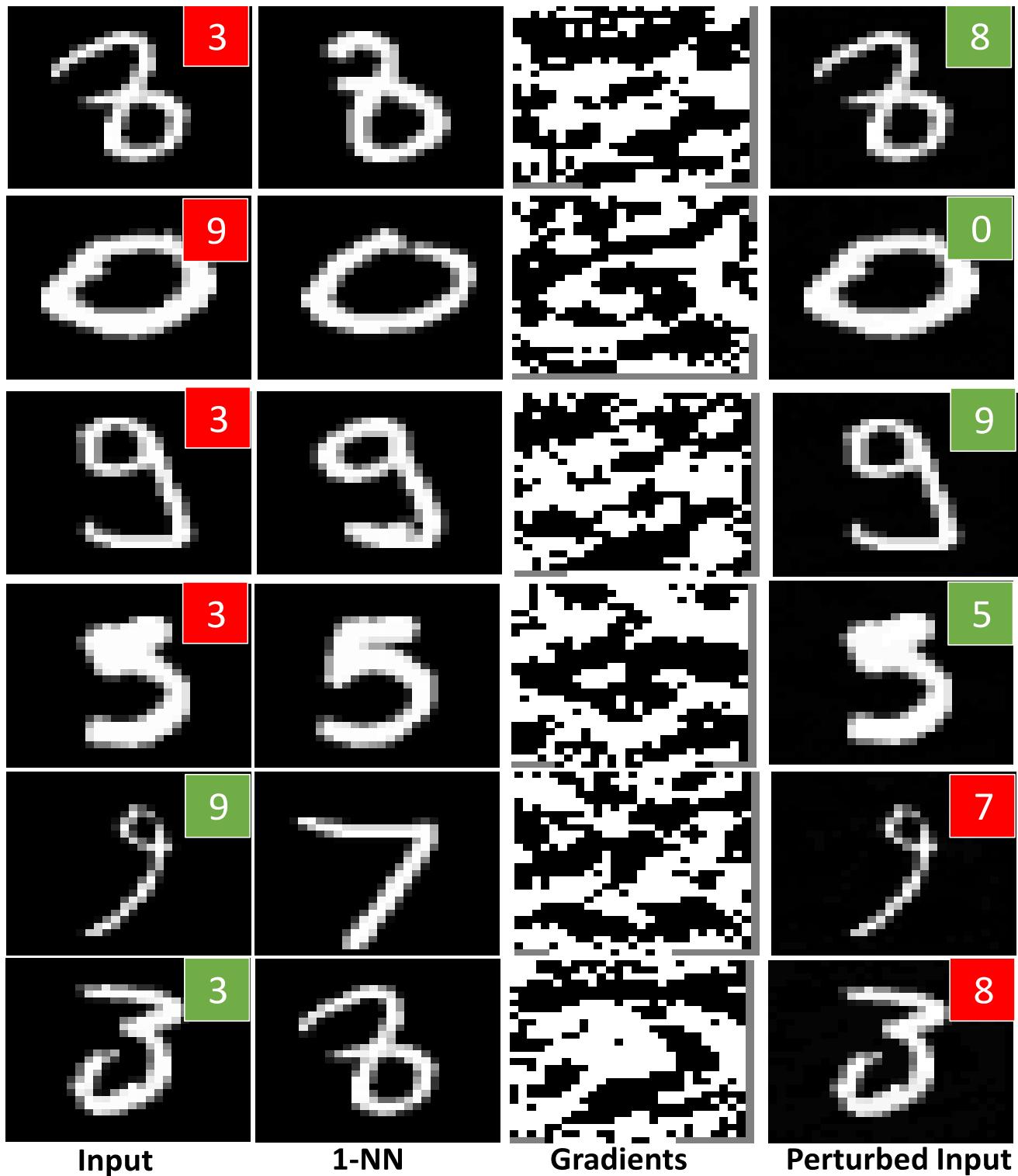


Figure 16: Example results of using the proposed approach for MNIST digits classification task. Top four rows shows situations where our approach was successful in correcting the classifier errors while bottom two rows showcase the failures. The red and green labels show the final deep network output: red indicates a mistake and green indicates correct prediction.