

It GAN DO Better: GAN-based Detection of Objects on Images with Varying Quality

Charan D. Prakash
Arizona State University
cprakash@asu.edu

Lina J. Karam
Arizona State University
karam@asu.edu

Abstract

In this paper, we propose in our novel generative framework the use of Generative Adversarial Networks (GANs) to generate features that provide robustness for object detection on reduced quality images. The proposed GAN-based Detection of Objects (GAN-DO) framework is not restricted to any particular architecture and can be generalized to several deep neural network (DNN) based architectures. The resulting deep neural network maintains the exact architecture as the selected baseline model without adding to the model parameter complexity or inference speed. We first evaluate the effect of image quality not only on the object classification but also on the object bounding box regression. We then test the models resulting from our proposed GAN-DO framework, using two state-of-the-art object detection architectures as the baseline models. We also evaluate the effect of the number of re-trained parameters in the generator of GAN-DO on the accuracy of the final trained model. Performance results provided using GAN-DO on object detection datasets establish an improved robustness to varying image quality and a higher mAP compared to the existing approaches.

1. Introduction

Deep learning and neural networks have been a popular choice for computer vision based applications such as multi-class object detection and classification. The parameters θ of a neural network that is designed for object detection are learned by training the network to fit a function $f(x; \theta) = y$ where x represents an image in the training dataset \mathbf{X} and $y \in \mathbf{Y}$ represents the output bounding boxes and their corresponding class labels for the image x . The assumption in the above training approach is that the images in the test dataset are also drawn from the same distribution as that of the training dataset, \mathbf{X} . However, in practice, the test data might not lie in the distribution \mathbf{X} , due to various factors such as defocus and camera shake blur or images af-



Figure 1: A reduced quality image with ground-truth object shown using a dotted bounding box is used as an input for object detection for different SSD300 [21] based DNN models: baseline SSD300 model trained on high quality images, fine-tuned SSD300 model trained on images with varying quality, and SSD300 model trained on images with varying quality using our proposed GAN-DO framework. Each text box below an image specifies the class of the corresponding object bounding box along with the predicted class confidence score.

ected by noise during sensor acquisition and compression. In such cases, the neural network trained on the distribution \mathbf{X} leads to errors in object detection. Such errors in object detection could lead to several issues ranging from social controversies (e.g., person misclassified as animal) to fatal accidents (e.g., autonomous vehicles failing to detect pedestrians). This behaviour is of concern in the light of recent studies showing that variations in image quality deteriorates the performance of DNNs significantly [6, 28].

In this paper, we propose a framework for GAN-based detection of objects (GAN-DO) that learns an adversarial objective. The proposed GAN-DO framework guides the

learning in a direction such as to maximize the similarity between the features that are computed by the GAN’s generator for the reduced quality images and the features that are computed by the baseline model for the original high quality images. This ability of the GAN-DO framework leads to robust object detection, as shown in Fig. 1. Our proposed GAN-DO framework is not restricted to a particular architecture and is designed to accommodate the use of several DNN based object detection models and to improve the accuracy and the robustness of the selected baseline model to images with varying quality. We believe this to be the first work to adopt a GAN framework for training object detection models in order to obtain robust object detection on reduced quality images.

Three contributions are made in this paper. Firstly, we evaluate the effect of varying image quality on the object classification and the object bounding box regression of object detection models. Secondly, we propose a novel generative GAN-DO framework that consists of two neural networks collectively working as a GAN that learns an adversarial objective through training, to add robustness to the object detection model. We show that the proposed GAN-DO framework outperforms the widely used fine-tuning framework in improving the accuracy and the robustness of the baseline object detection model to varying image quality while maintaining the identical model architecture, complexity and inference speed of the baseline model. Finally, we investigate the effect of the number of re-trained parameters using the proposed GAN-DO framework, on the object detection accuracy.

2. Related Work

Neural networks are known to be susceptible to variations in image quality for the task of image classification. Recent work has shown that the performance of neural networks decreases on reduced quality images [6, 13, 28, 33]. Tests [6, 13] showed that some architectures such as VGG16 [30] were more robust to variations in image quality than other architectures, such as GoogleNet [32]. Prior work mainly concentrated on the performance of DNNs for image classification and not for object detection.

One of the well-known approaches of improving the accuracy of the pre-trained models is to “fine-tune” the pre-trained model on reduced quality images. Previous work showed that fine-tuning can improve the accuracy for the task of image classification on images with distortions such as noise and blur [1, 33]. Vasiljevic *et al.* also showed that fine-tuning on a uniform mix of sharp and blurred images produced improved accuracy [33].

In the literature, an additional constraint has been previously imposed during fine-tuning in the form of stability loss [34] for achieving robustness in applications such as image classification and similar-image ranking. How-

ever, recent work [31] has shown that, while stability training works well for some types of distortions such as JPEG bit-rate compression, the performance degrades severely for other types of distortions such as noise and blur [1]. This characteristic limits the abilities of stability training since existing networks are already known to be robust to JPEG distortions [6]. Furthermore, manually choosing or designing effective stability loss terms for each distortion and task is challenging. For example, the authors of [34] propose K-L divergence as the stability term for the task of image classification and L_2 distance as the stability term for the task of feature embedding and similar-image ranking. No loss term is proposed for the task of object detection that includes bounding box regression.

Feature quantization is another approach for image classification that was proposed to improve the robustness of DNNs to varying image quality without changing the DNN model complexity. Sun *et al.* propose different types of additional non-linearities such as flooring and exponential power operations on the features for feature quantization [31]. However, tests in [31] show that there is no single non-linearity that performs better on all types of distortions. In many cases the model trained using quantization is seen to be more susceptible to distortion than the selected baseline model. Moreover, there are distortions such as defocus blur where the baseline model performs better than all of the variations proposed using feature quantization [31].

Dodge and Karam proposed the architecture of Mix-QualNet [7], an ensemble method based on mixture of experts. Each expert in the model is trained on a particular quality degradation and the gating network predicts the type and level of distortion. Borkar and Karam proposed Deep-Correct [1] to identify and rank filters that are more susceptible to image quality reductions than other filters. An additional stack of convolutional layers are added to these filters to improve the network performance while the other filters remain unchanged. Diamond *et al.* [5] proposed an architecture that prepended a network to the classification network to produce a task-oriented intermediary image that is optimized for image classification. The model complexity and the inference speed of the above mentioned frameworks increase due to the presence of additional network layers and the prepended network, respectively. Moreover, as mentioned previously, all the aforementioned approaches are focused on the task of image classification alone and do not consider the task of object detection, which not only includes object classification at different scales but also includes object localization and bounding box regression.

While there are image denoising [4] and image deblurring [8, 16, 23] methods that can be used as a pre-processing stage for object detection, these methods add significant computational overhead during inference. As an example, the work in [16] requires the image to be fed-forward

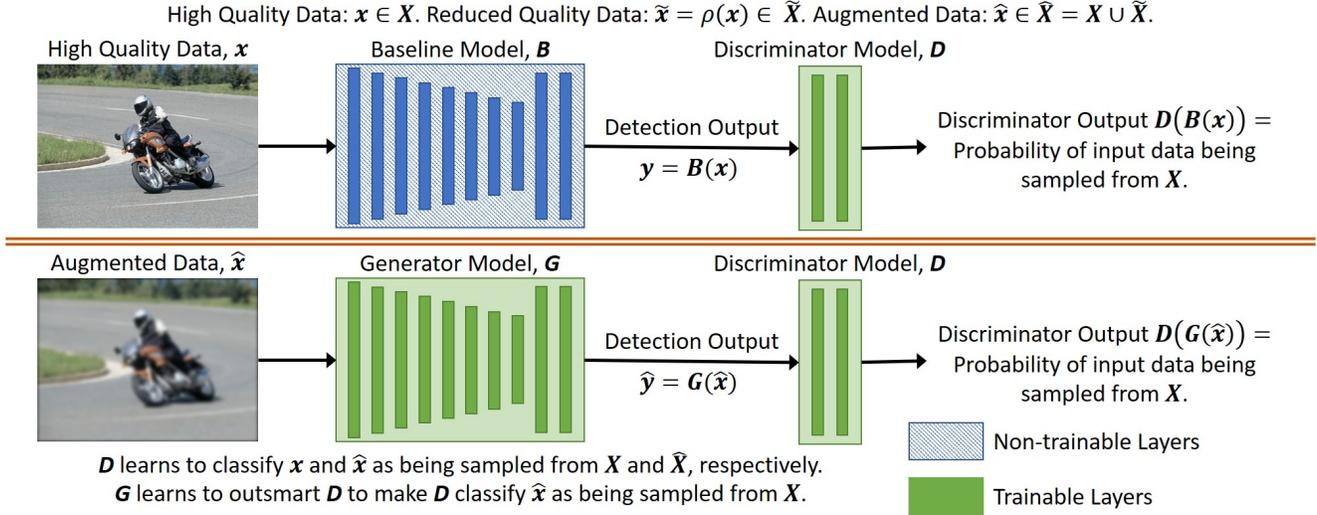


Figure 2: Block diagram of the proposed GAN-based Detection of Objects (GAN-DO) framework. The discriminator learns to distinguish between the pre-trained baseline model output y of the high quality original data x and the generator output \hat{y} of the augmented data \hat{x} with varying quality. The generator G learns to outsmart the discriminator D to make D classify the generator’s output of the augmented data as the baseline model’s output of the original data.

through 2 strided convolutions, 9 residual blocks and 2 transposed convolutional layers to generate a deblurred image. This additional computational overhead results in additional model capacity and complexity, and a latency in computing the object detection output which might not be acceptable in critical real-time applications. Moreover, recent work [1, 34] show that denoised and deblurred images generated for improved visual perception do not translate to better performance in image understanding tasks. Additionally, one would require another blur/noise detection module to decide whether or not to use the pre-processing stage during inference in scenarios that tackle both sharp and blurred/noisy images. In the context of object detection, prior work [14] has shown that fine-tuning can outperform domain adaptation methods when there is availability of sufficient annotated training data. GANs have been previously used to achieve resilience in object scale variations [17] by ensuring the features for object detection of small objects look similar to the features of large objects. However, no prior work has investigated the effect of reduced quality images on object detection. Hence, there is a need to produce DNN object detection models that are robust to varying image quality.

3. Proposed Framework

As described in the previous section, existing methods [34] employ an additional loss term along with the loss for the intended task in order to introduce robustness to the network. However, manually choosing or designing the stabil-

ity loss for each desired task and/or for each image quality level is challenging and none of the aforementioned methods propose an additional loss term for the task of object detection. Therefore, in our proposed framework, we design a network that learns an adversarial objective through training, to add robustness to the object detection model. Our proposed training framework, GAN-DO, consists of two neural networks, namely a generator and a discriminator collectively working as a GAN. However, only the generator needs to be retained during deployment of the object detection model at test time.

For every high quality image $x \in X$ in the clean original dataset X , we construct the reduced quality dataset \tilde{X} , using the image quality distortion $\rho_j(\cdot)$ such that $\tilde{x} = \rho_j(x)$, where $\rho_j(\cdot)$ is a randomly picked quality reduction level from a pool of J levels. We create an augmented dataset \hat{X} that combines the original dataset X and the reduced quality dataset \tilde{X} . This augmented dataset \hat{X} is used for training in the proposed framework. For comparison, the same augmented dataset is also used to train a fine-tuned model.

Fig. 2 shows the block diagram of the proposed GAN-DO framework. The framework consists of a generator and a discriminator learning to outsmart each other. During training, the object detection output $y \in Y$ for each of the high quality images x in the training set X , is computed using the pre-trained model (pre-trained on high quality images), referred henceforth as the baseline model B . Similarly, the object detection output $\hat{y} \in \hat{Y}$ for each of the varying quality images \hat{x} in the augmented dataset \hat{X}

is computed using the generator G . The objective of the discriminator D is to accurately classify y and \hat{y} to be originating from $B(x)$ and $G(\hat{x})$, respectively. The objective of the generator G is to outsmart the discriminator to make D classify \hat{y} as originating from $B(x)$.

In other words, instead of specifying an explicit stability loss term, we train the discriminator D to learn an adversarial objective that distinguishes $B(x)$ from $G(\hat{x})$. Based on this adversarial objective learned by D , the generator G learns to minimize the distance between y from \hat{y} such that the generator’s output due to \hat{x} is similar to y . Upon successful completion of training, the discriminator D can be discarded and the generator G can be used instead of the baseline model to provide increased robustness to variations in image quality. The details of the loss functions, architectures and training methodologies are provided in the following subsections.

3.1. Loss Function

GAN networks are known to train faster and more effectively when they are combined with the task-oriented loss [12, 22, 23, 25]. Therefore the total loss L_{total} of the proposed framework is given by:

$$L_{total} = L_{OD} + \lambda L_{GAN} \quad (1)$$

where L_{OD} is the object detection loss of the baseline model and L_{GAN} is the adversarial objective of the GAN. For convenience but without loss of generality, in order to consider L_{total} as a loss function to be minimized, λ is a weighting factor that is set to a positive value when training the generator G and to a negative value when training the discriminator D . The terms L_{OD} and L_{GAN} are described in more detail below.

The object detection loss L_{OD} is formulated as:

$$L_{OD} = \frac{1}{N} (L_{class} + \alpha L_{bb}) \quad (2)$$

where N is the number of predicted bounding boxes, L_{class} is the classification loss, L_{bb} is the bounding box regression loss and α is a hyper-parameter. The specific choice of L_{class} , L_{bb} and α are model-specific loss terms and depends on the selected baseline model.

We consider the SSD [21] and the RetinaNet [19] models for evaluating our GAN-DO framework in this paper. Both the above mentioned models employ smooth L1 loss [27] for bounding box regression. SSD uses categorical cross-entropy and hard negative mining [21] for L_{class} and N represents only a small subset of predicted bounding boxes using hard negative mining. RetinaNet uses Focal Loss [19] for L_{class} and uses all the predicted boxes for loss calculation. In this paper, we use the object detection loss terms as proposed for their respective models for both fine-tuning

and for our GAN-DO framework, in order to provide a fair comparison.

In order to make the network robust to variations in image quality, we propose to combine the object detection loss L_{OD} with an adversarial objective L_{GAN} during training. Consider a dataset of original data $x \in X$ and the corresponding set of augmented data $\hat{x} \in \hat{X}$ with probability distributions P_x and $P_{\hat{x}}$, respectively. Let $y \in Y$ and $\hat{y} \in \hat{Y}$ be the corresponding outputs of x and \hat{x} computed from the baseline model B and generator G , respectively, as follows (Fig. 2):

$$\begin{aligned} y &= B(x) \\ \hat{y} &= G(\hat{x}) \end{aligned} \quad (3)$$

The generator G generates an output $G(\hat{x})$ using the input \hat{x} . The discriminator D distinguishes the “real” original image data x from the “fake” augmented image data \hat{x} using the object detection outputs, $B(x)$ and $G(\hat{x})$. The objective function of such a GAN is given by:

$$\begin{aligned} L_{GAN} &= \mathbb{E}_{x \sim P_x(x)} (\log D(B(x))) \\ &+ \mathbb{E}_{\hat{x} \sim P_{\hat{x}}(\hat{x})} (\log(1 - D(G(\hat{x})))) \end{aligned} \quad (4)$$

where the discriminator output $D(\cdot)$ represents the discriminator’s predicted probability of the input image belonging to the original data X . The goal of the discriminator is to maximize this objective function and the goal of the generator is to minimize this objective function. L_{GAN} is implemented as a binary cross-entropy function as discussed in [10].

3.2. Network Architecture

We consider two architectures in this paper for the task of object detection, SSD [21] and RetinaNet [19]. The SSD architecture [21] uses VGG16 [30] as the feature extractor of the object detection network. In this paper, we use SSD300 [21], which corresponds to the resizing of input images to size 300 x 300 pixels. The RetinaNet architecture consists of a ResNet [11] together with a Feature Pyramid Network (FPN) [18] as the feature extractor. In this paper, we use RetinaNet50-400 [19], which uses ResNet-50 [11] as the feature extractor. The input images are resized such that the longest side of the image is resized to 400 pixels while maintaining the aspect ratio of the original input image. We import the baseline model architecture as the generator architecture in order to retain the model complexity and inference speed. The generator computes the object detection output in a way similar to that of the baseline model, for the input images from the augmented dataset \hat{X} .

The output from the object detection models serve as input to the discriminator. The output of the discriminator is the probability of the input image being sampled from X . The discriminator of our implementation contains a single

fully connected layer. In our tests we observed that increasing the model capacity of the discriminator with more layers made the discriminator so strong that the generator could not outsmart the discriminator. Based on the size, complexity and capacity of the baseline model, the discriminator model size can be varied for other tasks. Although additional complexity is added to the framework during training by the use of the discriminator, it should be noted that the discriminator is discarded during inference. Therefore, during inference, the model complexity and speed of the model resulting from the GAN-DO framework remains identical to the baseline model.

3.3. Training and Inference

Similar to fine-tuning, the weights for the generator in the GAN-DO framework are initialized with the pre-trained weights of the baseline model (SSD300 or RetinaNet50-400). The augmented dataset $\hat{\mathbf{X}}$ is created such that it contains a uniform mix (1:1 ratio) of high quality images \mathbf{x} and reduced quality images $\hat{\mathbf{x}}$ in order for the network to perform well on both high quality and reduced quality images. The pseudo-code of the proposed framework that is used to train an object detection model is given in Algorithm 1.

For each iteration of training, an image $\hat{\mathbf{x}}_s$ in a random augmented mini-batch $\hat{\mathbf{X}}_S$ of size S ($s \sim [1, 2, \dots, S]$), is selected on the fly using the equation:

$$\hat{\mathbf{x}}_s = \begin{cases} \mathbf{x}_s, & \text{if } s \leq S/2 \\ \rho_j(\mathbf{x}_s), j \sim U[1, 2, \dots, J], & \text{otherwise} \end{cases} \quad (5)$$

where ρ_j is a randomly picked j -th quality distortion kernel from a pool of J distortion levels. Details about the distortions are presented in Section 4.2.

The original mini-batch with high quality images \mathbf{x}_s , and the augmented mini-batch with varying quality images $\hat{\mathbf{x}}_s$, are used to compute the object detection outputs from the baseline and the generator models, $\mathbf{B}(\mathbf{x}_s)$ and $\mathbf{G}(\hat{\mathbf{x}}_s)$, respectively, as in Eq. (3). The discriminator \mathbf{D} is first trained to predict $\mathbf{x}_s \in \mathbf{X}$ (or equivalently $\mathbf{y}_s \in \mathbf{Y}$) and $\hat{\mathbf{x}}_s \in \hat{\mathbf{X}}$ (or equivalently $\hat{\mathbf{y}}_s \in \hat{\mathbf{Y}}$), for $\mathbf{B}(\mathbf{x}_s)$ and $\mathbf{G}(\hat{\mathbf{x}}_s)$, respectively. The generator \mathbf{G} is trained to predict $\mathbf{G}(\hat{\mathbf{x}}_s)$ such that \mathbf{D} predicts $\hat{\mathbf{x}}_s \in \mathbf{X}$ (or equivalently $\hat{\mathbf{y}}_s \in \mathbf{Y}$).

All networks in this paper are trained using the Adam optimizer [15]. For fine-tuning, the decay rates of the first and second moments of gradients (β_1 and β_2 in [15]) are set to 0.9 and 0.99, respectively. In the proposed GAN framework, the decay rates of the first and second moments of gradients are set to 0.5 and 0.99, respectively, for both \mathbf{G} and \mathbf{D} as these values are shown [26] to stabilize adversarial training.

Algorithm 1 Proposed training methodology using GAN-DO framework

Input: Training dataset with original images \mathbf{X} . Distortion kernels $\rho_j(\cdot)$, $j = 1, \dots, J$. Baseline model \mathbf{B} , Generator model \mathbf{G} initialized with pre-trained weights of \mathbf{B} , Discriminator model \mathbf{D} with Normal initialized weights, number of training iterations L and mini-batch size S .

- 1: **for** $l = 1$ to L **do**
 - 2: Draw a random mini-batch of images \mathbf{x}_s , $s = 1, \dots, S$ from training dataset \mathbf{X} .
 - 3: Create the augmented mini-batch $\hat{\mathbf{X}}_S$ with images $\hat{\mathbf{x}}_s$ on the fly using Eq. (5).
 - 4: Use $\mathbf{B}(\mathbf{x}_s)$ to train \mathbf{D} to predict $\mathbf{x}_s \in \mathbf{X}$. Update weights of \mathbf{D} to minimize Eq. (1) with $\lambda < 0$.
 - 5: Use $\mathbf{G}(\hat{\mathbf{x}}_s)$ to train \mathbf{D} to predict $\hat{\mathbf{x}}_s \in \hat{\mathbf{X}}$. Update weights of \mathbf{D} to minimize Eq. (1) with $\lambda < 0$.
 - 6: Train \mathbf{G} to predict $\mathbf{G}(\hat{\mathbf{x}}_s)$ such that $\mathbf{D}(\mathbf{G}(\hat{\mathbf{x}}_s))$ predicts $\hat{\mathbf{x}}_s \in \mathbf{X}$. Update weights of \mathbf{G} to minimize Eq. (1) with $\lambda > 0$.
 - 7: **end for**
 - 8: Discard \mathbf{D} . Use \mathbf{G} to perform object detection with improved robustness to variations in image quality.
-

4. Experimental Results

4.1. Datasets and Evaluation Metrics

Since previous work [21] has shown that model performance can be improved by including both the PASCAL VOC2007 [9] and PASCAL VOC2012 [24] training images, the union of the PASCAL VOC2007 `train` and VOC2012 `trainval` images were used for training the models for both the proposed framework and fine-tuning. The PASCAL VOC2007 `val` images were used for validation (more details provided in Section 4.4) and the PASCAL VOC2007 `test` images were used for testing and both Average Precision (AP) and mean Average Precision (mAP) were computed across all 20 annotated classes as per the PASCAL VOC object detection evaluation metric [9].

The Common Objects in Context (COCO) dataset [20] evaluates object detection methods over 80 object categories. All images from COCO `trainval35k` were used for training and all images from COCO `minival` were used for testing as COCO does not publicly provide labels for evaluation on their test dataset. We compute $\text{AP}^{\text{IoU}=0.5}$, $\text{AP}^{\text{IoU}=[0.50:0.05:0.95]}$, $\text{AP}^{\text{IoU}=0.75}$, AP^{small} , $\text{AP}^{\text{medium}}$ and AP^{large} according to the COCO object detection evaluation metrics [3].

4.2. Image Quality Distortions

Blur and noise are among the most commonly encountered image quality distortions in many popular applications

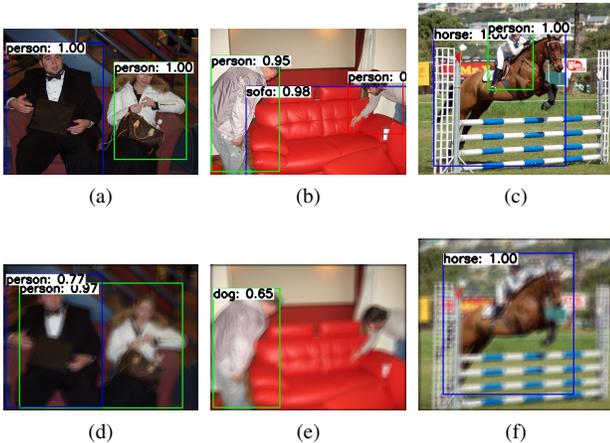


Figure 3: Effect of image quality on object detection. Top row contains object detection results using the SSD300 baseline model on high quality images from the PASCAL VOC2007 test dataset. Bottom row contains object detection results using the SSD300 baseline model on corresponding reduced quality images (affected by camera shake blur).

such as social media, cellphones and autonomous driving. In this paper, we evaluate our framework using three different types of blur, namely, camera shake, Gaussian, and uniform circular average kernels referred to as defocus blurs in [31, 33]. We choose six levels of blur for Gaussian and defocus blur. The radius of the blur kernels in pixels is varied in the range [2,12] in incremental steps of 2 pixels for Gaussian blur and defocus blur. The standard deviation of each kernel of the Gaussian blur is set to one half of the radius of the blur kernel. For camera shake blur, we generate 50 random camera shake kernels as described in [2]. In this paper, we also evaluate our proposed framework on images affected by additive white Gaussian noise (AWGN). We choose five levels of AWGN with the standard deviation in the range [20,100] in incremental steps of 20. During inference, the quality of each image x_s in the mini-batch was reduced by applying blur or noise of a randomly chosen level to create a reduced quality mini-batch with images \tilde{x}_s . The models were tested on the reduced quality mini-batches for object detection accuracy.

4.3. Effect of Image Quality on Bounding Box Regression and Object Classification

The effect of image quality on DNN’s classification accuracy was discussed in prior work [1, 5, 6, 33] for the image classification task using the ImageNet 2012 dataset (ILSVRC2012) [29]. The focus of this paper is on object detection (which includes both object classification and object localization) rather than the task of image classification

Table 1: Effect of image quality on object classification loss and bounding box regression loss for the SSD300 baseline model on the PASCAL VOC2007 test dataset. Bold numbers show best accuracy.

Method	Distortion	Classification Loss	Regression Loss
SSD300	None	1.954	0.973
	Gaussian Blur	4.012	1.461
	Defocus Blur	4.250	1.544
	Camerashake Blur	4.305	1.580
	AWGN	4.137	1.532

Table 2: Effect of image quality on object classification loss and bounding box regression loss for the baseline models with different architectures on the COCO minival dataset. Bold numbers show best accuracy.

Method	Distortion	Classification Loss	Regression Loss
SSD300	None	2.492	1.238
	Gaussian Blur	3.347	1.663
	Defocus Blur	3.518	1.755
	Camerashake Blur	3.447	1.732
	AWGN	3.917	1.872
RetinaNet50-400	None	0.910	0.266
	Gaussian Blur	1.272	0.480
	Defocus Blur	1.364	0.515
	Camerashake Blur	1.353	0.471
	AWGN	1.383	0.526

considered in prior work. Hence, as a contribution, in addition to the object classification loss L_{class} , we also investigate the effect of image quality on object bounding box regression L_{bb} (as defined in Section 3.1) for the object detection task. Examples of incorrect object detection on reduced quality images due to errors in bounding box regression and object classification are shown in Fig. 3. Tables 1 and 2 quantify the effect of image quality on both object classification loss and object bounding box regression loss of the baseline object detection models on the PASCAL VOC2007 test and COCO minival datasets, respectively¹.

From Tables 1 and 2, we notice that the reduced image quality results in an increase in both L_{class} and L_{bb} , thereby decreasing the object detection accuracy. Camera shake blur was found to affect most the model performance for the SSD300 on the PASCAL VOC2007 test dataset and results in up to a 120% increase in L_{class} and up to a 60% increase in L_{bb} , compared with the corresponding losses of the SSD300 baseline model using high quality input images. AWGN was found to affect most the model performance for both SSD300 and RetinaNet50-400 on the COCO minival dataset and results in up to a 57% increase in L_{class} and 98% increase in L_{bb} . Therefore, there is a need for more robust DNN models for the task of object

¹Results on the PASCAL VOC2007 dataset are provided only for the SSD baseline model as the authors of RetinaNet [19] do not provide a RetinaNet baseline model trained on the PASCAL VOC dataset.

Table 3: mAP of SSD300 based DNN models on the PASCAL VOC2007 *test* dataset using images of varying quality. Bold numbers show best accuracy.

Distortion	Approach		
	Baseline	Fine-tuning	GAN-DO
Gaussian Blur	44.70	63.03	64.40
Defocus Blur	40.77	61.83	63.08
Camerashake Blur	39.69	64.63	67.19
AWGN	42.20	65.84	67.47

detection with both classification and bounding box regression features that are resilient to variations in image quality.

4.4. Performance of the Proposed GAN-DO Framework

For the SSD300 trained using the GAN-DO framework, the weight of the adversarial objective $|\lambda|$ (in Eq. (1)) was set to 1. The generator was trained every iteration while the discriminator was trained every other iteration. For training SSD300 on the PASCAL VOC, the learning rate for both fine-tuning and the GAN-DO framework (generator and the discriminator) were set to 10^{-5} . The learning rate of all the models (fine-tuning model, generator and discriminator) was divided by 10 when the validation loss failed to improve for τ consecutive epochs. The process of decaying the learning rate was repeated twice for the fine-tuning model and the GAN-DO framework before terminating the training process. The parameter τ was empirically determined to provide best results when set to 4 for fine-tuning and 10 for the GAN-DO framework.

The performance improvement that is obtained using the GAN-DO framework on the PASCAL VOC2007 *test* dataset is shown in Table 3. It can be seen that the GAN-DO framework results in the highest mAP as compared to the baseline and fine-tuned models across all the considered distortion types. Furthermore, the GAN-DO framework results in the highest AP across most of the object classes.

The ability of the GAN-DO framework to perform better at higher levels of blur is shown in Table 4. Table 4 shows that, while the baseline model and the GAN-DO model exhibit comparable performance for the lowest blur level (blurs with 2 pixel radius), the proposed GAN-DO based model achieves the highest performance as compared to both the baseline and the fine-tuned models at all other blur levels. Table 4 also shows that the proposed GAN-DO based model performs better than fine-tuning on high-quality images ($r=0$). Therefore, training using our GAN-DO framework results in a DNN model that is more robust to varying image quality as compared to fine-tuning and the baseline model.

In order to check the robustness of the GAN-DO framework to unseen distortions, tests were conducted where the

Table 4: mAP of SSD300 based DNN models on the PASCAL VOC2007 *test* dataset using different levels of blurs. ‘r’ specifies the radius of the distortion kernel. Bold numbers show best accuracy.

Distortion	Approach	r=0	r=2	r=4	r=6	r=8	r=10	r=12
		Baseline	77.47	75.67	65.96	51.45	36.07	23.84
Gaussian Blur	Fine-tuning	74.73	73.93	70.16	65.47	60.58	55.51	51.00
	GAN-DO	76.17	75.29	71.5	66.79	61.87	57.17	52.47
Defocus Blur	Baseline	77.47	74.68	61.60	44.91	30.07	20.17	13.60
	Fine-tuning	74.96	73.32	68.63	63.98	59.44	54.85	50.42
	GAN-DO	75.94	74.20	69.83	65.53	61.41	57.05	52.61

Table 5: mAP comparison of SSD300 based DNN models on the PASCAL VOC2007 *test* dataset using unseen blur distortions. Bold numbers show best accuracy.

Trained on	Framework	Tested on		
		Gaussian Blur	Defocus Blur	Camerashake Blur
Gaussian Blur	Fine-tuning	63.03	59.97	54.39
	GAN-DO	64.40	62.14	55.62
Defocus Blur	Fine-tuning	61.03	61.83	58.04
	GAN-DO	61.66	63.08	58.18
Camerashake Blur	Fine-tuning	57.28	57.22	64.63
	GAN-DO	57.91	59.67	67.19

models trained on one type of blur were tested on other types of blurs. It was observed in [33] that fine-tuning generalizes well on unseen blur distortions for the task of image classification. However, it can be seen from Table 5 that models trained with our GAN-DO framework perform better than fine-tuning in terms of mAP across all unseen blur distortions. Consequently, DNN models trained using our GAN-DO framework generalize better than fine-tuned models and result in an increased robustness and improved mAP accuracy on images affected by unseen distortions.

Table 6 shows the effect of the re-trainable parameters using the GAN-DO framework on the accuracy of the SSD300 based object detection. Experiments were conducted on the PASCAL VOC2007 *test* dataset with camera shake blur to compute mAPs of different configurations of the proposed framework, represented by GAN- k . Only filters from the first layer up to the k layer (checkpoint) were trained while the filters in the other subsequent layers remained unchanged. Instead of updating all the filters in the baseline model (represented by GAN-*all_layers*), the checkpoint k was set at different layers in the SSD300 architecture. The model’s performance results in terms of mAP when the checkpoint k was set at layers `pool1`, `pool2`, `pool3` and `pool4` are shown in Table 6.

It can be seen from Table 6 that the GAN-DO framework achieves the highest mAP if the parameters are re-trained up to `pool3` (GAN-*pool3*) with a comparable but slightly lower performance for GAN-*pool4*. Consequently it can be observed that a higher level of robustness is achieved by re-training only the feature extractor layers instead of the entire network which includes task-specific layers for classification and localization. However, it was also observed

Table 6: mAP comparison of SSD300 based DNN models on the PASCAL VOC2007 *test* dataset for camera shake blur. Bold numbers show best performance.

Model	High Quality images	Reduced Quality images	Epochs to converge
Baseline	77.47	39.73	-
Fine-tuning	74.90	64.63	35
GAN- <i>pool1</i>	76.56	56.13	62
GAN- <i>pool2</i>	76.51	63.70	58
GAN- <i>pool3</i>	76.46	68.12	66
GAN- <i>pool4</i>	74.93	68.08	49
GAN- <i>all_layers</i>	76.12	67.19	29

that training all layers in the model (GAN-*all_layers*) with the proposed framework converges in fewer epochs as compared to the other configurations of GAN-*k*.

We also test the proposed GAN-DO framework with SSD300 on a dataset larger than the PASCAL VOC such as COCO. All models (fine-tuning, generator and discriminator) were trained with a learning rate of 10^{-5} for the first 20 epochs and 10^{-6} for the next 10 epochs. Table 7 shows the performance results that are obtained by training SSD300 with the GAN-DO framework across varying image quality. It can be seen that the GAN-DO framework results in the highest accuracy for all the types of tested quality reductions across all IoU thresholds. It can also be seen that training using the proposed framework achieves up to 12% performance improvement in terms of AP on large objects as compared to fine-tuning. This characteristic can be highly useful in scenarios like autonomous driving where closer objects that appear larger need to be detected with higher accuracy.

In order to show that the GAN-DO framework works on different architectures, we train RetinaNet50-400 [19] with the proposed framework on COCO. The weight of the adversarial objective $|\lambda|$ (in Eq. (1)) was set to 0.5. Both the generator and the discriminator were trained at each iteration since the RetinaNet50-400 has more parameters than SSD300 and can adapt quicker based on the feedback from the discriminator. The fine-tuned model and the generator of the proposed framework were trained with a learning rate of 10^{-5} for the first 20 epochs and 10^{-6} for the next 10 epochs. The learning rate of the discriminator was set at 10^{-4} throughout the training process.

Table 8 shows the performance results that are obtained by training RetinaNet50-400 with the GAN-DO framework across different distortions. Compared to SSD300, RetinaNet50-400 is a larger network with more trainable parameters than SSD300. Fine-tuning benefits from the model capacity of RetinaNet50-400 to recover most of the lost accuracy in terms of AP. However, compared to the baseline and fine-tuned models, models trained using the GAN-DO framework achieve the best performance in terms of AP across all types of image quality reductions except defocus blur, while remaining comparable to fine-tuning on images affected by defocus blur. Furthermore, similar to SSD300,

Table 7: Object detection accuracy (AP) comparison of SSD300 based DNN models on the COCO *minival* dataset for varying image quality. Bold numbers show best accuracy.

Distortion	Method	Avg. Precision, IoU:			Avg. Precision, Area:		
		0.50:0.95	0.50	0.75	small	medium	large
None	Baseline	24.7	42.4	25.3	5.9	26.4	41.4
	Baseline	14.9	26.0	15.1	2.2	13.2	30.3
Gaussian Blur	Fine-tuning	16.6	30.4	16.3	2.9	17.9	31.5
	GAN-DO	17.6	31.7	17.2	3.0	17.8	33.3
	Baseline	13.8	24.0	13.9	1.8	12.1	28.2
Defocus Blur	Fine-tuning	15.8	29.1	15.4	2.6	16.8	30.3
	GAN-DO	17.1	31.0	17.0	2.3	16.8	32.7
	Baseline	13.6	24.5	13.0	1.2	11.6	28.5
Camerashake Blur	Fine-tuning	16.4	31.2	15.5	2.9	18.4	30.7
	GAN-DO	18.2	33.3	17.7	2.8	18.6	33.8
	Baseline	11.9	21.3	11.9	2.1	12.2	22.6
AWGN	Fine-tuning	15.9	30.4	14.9	3.4	18.1	28.6
	GAN-DO	17.8	32.7	17.3	3.4	18.6	32.3

Table 8: Object detection accuracy (AP) comparison of RetinaNet50-400 based DNN models on the COCO *minival* dataset for varying image quality. Bold numbers show best accuracy.

Distortion	Method	Avg. Precision, IoU:			Avg. Precision, Area:		
		0.50:0.95	0.50	0.75	small	medium	large
None	Baseline	29.9	46.1	32.0	10.7	32.9	47.5
	Baseline	16.7	26.3	17.6	5.8	16.4	31.6
Gaussian Blur	Fine-tuning	24.1	37.9	25.4	8.3	25.2	40.4
	GAN-DO	24.6	38.4	26.1	8.0	25.9	42.1
	Baseline	15.1	24.2	15.7	4.4	14.3	29.6
Defocus Blur	Fine-tuning	24.6	38.7	25.6	7.6	25.4	42.2
	GAN-DO	24.5	38.5	25.6	7.1	25.2	42.4
	Baseline	14.6	24.4	14.8	2.9	13.5	29.2
Camerashake Blur	Fine-tuning	24.8	39.4	25.7	6.9	26.2	42.1
	GAN-DO	25.4	40.4	26.4	7.1	26.8	43.6
	Baseline	13.4	21.7	13.9	4.1	13.8	23.0
AWGN	Fine-tuning	24.2	38.6	25.1	7.5	25.4	40.6
	GAN-DO	24.5	38.9	25.4	7.6	25.7	40.7

RetinaNet50-400 models trained using the GAN-DO framework achieve a higher AP on larger objects as compared to fine-tuning, across all types of image quality reductions.

5. Conclusion

In this paper, we show that the accuracy of object detection networks is sensitive to images with varying quality. We propose a novel framework, called the GAN-DO framework, for re-training the parameters of the baseline model in order to increase the robustness of the model to varying image quality. The model resulting from our GAN-DO framework is identical to the baseline model in terms of model complexity and inference speed while achieving robustness to varying image quality. The GAN-DO framework outperforms the fine-tuned and baseline models across different types of tested image quality reductions and over different baseline DNN models.

References

- [1] Tejas S Borkar and Lina J Karam. DeepCorrect: Correcting DNN models against image distortions. *IEEE Transactions on Image Processing*, 28(12):6022–6034, 2019. 2, 3, 6
- [2] Ayan Chakrabarti. A neural approach to blind motion deblurring. In *European Conference on Computer Vision*, pages 221–235, 2016. 6
- [3] COCO object detection evaluation metric. <http://cocodataset.org/#detection-eval>. 5
- [4] Kostadin Dabov, Alessandro Foi, Vladimir Katkovnik, and Karen Egiazarian. BM3D image denoising with shape-adaptive principal component analysis. In *Signal Processing with Adaptive Sparse Structured Representations*, 2009. 2
- [5] Steven Diamond, Vincent Sitzmann, Stephen Boyd, Gordon Wetzstein, and Felix Heide. Dirty pixels: Optimizing image classification architectures for raw sensor data. *arXiv preprint arXiv:1701.06487*, 2017. 2, 6
- [6] Samuel Dodge and Lina Karam. Understanding how image quality affects deep neural networks. In *International Conference on Quality of Multimedia Experience*, pages 1–6, 2016. 1, 2, 6
- [7] Samuel F Dodge and Lina J Karam. Quality robust mixtures of deep neural networks. *IEEE Transactions on Image Processing*, 27(11):5553–5562, 2018. 2
- [8] Weisheng Dong, Lei Zhang, Guangming Shi, and Xin Li. Nonlocally centralized sparse representation for image restoration. *IEEE Transactions on Image Processing*, 22(4):1620–1630, 2012. 2
- [9] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes (VOC) challenge. *International Journal of Computer Vision*, 88(2):303–338, 2010. 5
- [10] Ian Goodfellow. NIPS 2016 tutorial: Generative adversarial networks. *arXiv preprint arXiv:1701.00160*, 2016. 4
- [11] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–778, 2016. 4
- [12] Phillip Isola, Jun-Yan Zhu, Tinghui Zhou, and Alexei A Efros. Image-to-image translation with conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 5967–5976, 2017. 4
- [13] Samil Karahan, Merve Kilinc Yildirim, Kadir Kirtac, Ferhat Sukru Rende, Gultekin Butun, and Hazim Kemal Ekenel. How image degradations affect deep CNN-based face recognition? In *2016 International Conference of the Biometrics Special Interest Group*, pages 1–5. IEEE, 2016. 2
- [14] Mehran Khodabandeh, Arash Vahdat, Mani Ranjbar, and William G Macready. A robust learning approach to domain adaptive object detection. In *IEEE International Conference on Computer Vision*, 2019. 3
- [15] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5
- [16] Orest Kupyn, Volodymyr Budzan, Mykola Mykhailych, Dmytro Mishkin, and Jiří Matas. Deblurgan: Blind motion deblurring using conditional adversarial networks. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 8183–8192, 2018. 2
- [17] Jianan Li, Xiaodan Liang, Yunchao Wei, Tingfa Xu, Jiashi Feng, and Shuicheng Yan. Perceptual generative adversarial networks for small object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1222–1230, 2017. 3
- [18] Tsung-Yi Lin, Piotr Dollár, Ross B Girshick, Kaiming He, Bharath Hariharan, and Serge J Belongie. Feature pyramid networks for object detection. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2117–2125, 2017. 4
- [19] Tsung-Yi Lin, Priyal Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 2018. doi: [10.1109/TPAMI.2018.2858826](https://doi.org/10.1109/TPAMI.2018.2858826). 4, 6, 8
- [20] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. Microsoft COCO: Common objects in context. In *European Conference on Computer Vision*, pages 740–755, 2014. 5
- [21] Wei Liu, Dragomir Anguelov, Dumitru Erhan, Christian Szegedy, Scott Reed, Cheng-Yang Fu, and Alexander C Berg. SSD: Single shot multibox detector. In *European Conference on Computer Vision*, pages 21–37, 2016. 1, 4, 5
- [22] Michael Mathieu, Camille Couprie, and Yann LeCun. Deep multi-scale video prediction beyond mean square error. *arXiv preprint arXiv:1511.05440*, 2015. 4
- [23] Thekke Madam Nimisha, Akash Kumar Singh, and Ambasamudram N Rajagopalan. Blur-invariant deep learning for blind-deblurring. In *IEEE International Conference on Computer Vision*, pages 4762–4770, 2017. 2, 4
- [24] The PASCAL Visual Object Classes Challenge 2012 (VOC2012). <http://host.robots.ox.ac.uk/pascal/VOC/voc2012/>. 5
- [25] Deepak Pathak, Philipp Krahenbuhl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 2536–2544, 2016. 4
- [26] Alec Radford, Luke Metz, and Soumith Chintala. Unsupervised representation learning with deep convolutional generative adversarial networks. *arXiv preprint arXiv:1511.06434*, 2015. 5
- [27] Shaoqing Ren, Kaiming He, Ross Girshick, and Jian Sun. Faster R-CNN: towards real-time object detection with region proposal networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 39(6):1137–1149, 2017. 4
- [28] Erik Rodner, Marcel Simon, Robert B Fisher, and Joachim Denzler. Fine-grained recognition in the noisy wild: Sensitivity analysis of convolutional neural networks approaches. *arXiv preprint arXiv:1610.06756*, 2016. 1, 2
- [29] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision*, 115(3):211–252, 2015. 6

- [30] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. [2](#), [4](#)
- [31] Zhun Sun, Mete Ozay, Yan Zhang, Xing Liu, and Takayuki Okatani. Feature quantization for defending against distortion of images. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 7957–7966, 2018. [2](#), [6](#)
- [32] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1–9, 2015. [2](#)
- [33] Igor Vasiljevic, Ayan Chakrabarti, and Gregory Shakhnarovich. Examining the impact of blur on recognition by convolutional networks. *arXiv preprint arXiv:1611.05760*, 2016. [2](#), [6](#), [7](#)
- [34] Stephan Zheng, Yang Song, Thomas Leung, and Ian Goodfellow. Improving the robustness of deep neural networks via stability training. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 4480–4488, 2016. [2](#), [3](#)