

Not Using the Car to See the Sidewalk — Quantifying and Controlling the Effects of Context in Classification and Segmentation

Rakshith Shetty¹ Bernt Schiele¹ Mario Fritz²

¹Max Planck Institute for Informatics, Saarland Informatics Campus

²CISPA Helmholtz Center i.G., Saarland Informatics Campus

¹firstname.lastname@mpi-inf.mpg.de

²firstname.lastname@cispa.saarland

Abstract

Importance of visual context in scene understanding tasks is well recognized in the computer vision community. However, to what extent the computer vision models for image classification and semantic segmentation are dependent on the context to make their predictions is unclear. A model overly relying on context will fail when encountering objects in context distributions different from training data and hence it is important to identify these dependencies before we can deploy the models in the real-world. We propose a method to quantify the sensitivity of black-box vision models to visual context by editing images to remove selected objects and measuring the response of the target models. We apply this methodology on two tasks, image classification and semantic segmentation, and discover undesirable dependency between objects and context, for example that “sidewalk” segmentation relies heavily on “cars” being present in the image. We propose an object removal based data augmentation solution to mitigate this dependency and increase the robustness of classification and segmentation models to contextual variations. Our experiments show that the proposed data augmentation helps these models improve the performance in out-of-context scenarios, while preserving the performance on regular data.

1. Introduction

Visual context of an object in an image is an important source of information for scene understanding tasks in both human and computer vision [21, 15]. Contextual cues such as presence of frequently co-occurring objects can help resolve ambiguities between visually similar classes and improve performance in various vision tasks including object detection [13, 3] and segmentation [24]. However, objects can also appear in previously unseen context or be absent from a very typical context. For example, we might find a keyboard on a desk without a monitor (object-without-

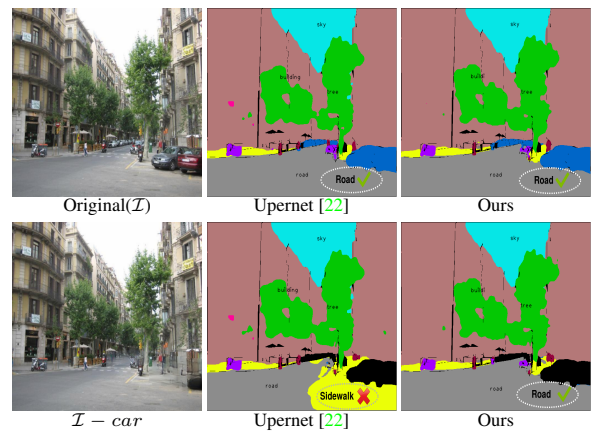


Figure 1: An example of the sensitivity of road and sidewalk segmentation to the context object *car*. Removing *car* from the image (second row) causes segmentation errors in the baseline model which hallucinates a sidewalk (yellow) when there is none. Our model trained with proposed data-augmentation is more robust to these context changes.

context), or find a monitor without a keyboard (context-without-object). While humans can handle both these atypical scenarios gracefully, computer vision models often fail by ignoring the visual evidence for the object in object-without-context case or hallucinating objects which are not actually present in the image in context-without-object case. For example, in our experiments we find that *keyboard* is often not recognized without a nearby *monitor*, and semantic segmentation of roads suffers without *cars* (see Figure 1). While context can be an important cue, this kind of too heavy or even pathological dependency on contextual signals is undesirable, and it is important to systematically identify and ideally fix such cases. In this work, we analyze and quantify the effect contextual information on two tasks, multi-label classification and semantic segmentation.

Context includes a lot of different kinds of information,

including co-occurring objects, scene type and lighting. For our analysis, we limit context to only the set of co-occurring objects in the image. While this might seem restrictive, we find in our analysis that image classification and segmentation models learn many interesting and undesirable dependencies between an object and other co-occurring objects (context) in the image. We use object removal as the main methodology to understand and quantify the role of context in downstream vision models. Specifically, we compare the output of the target models on the original input image and an edited version of this image with one object removed from it. If the model heavily uses the contextual relationship between removed object and the objects present in the image, removal will have an adverse effect on the model output. Measuring this helps us quantify the contextual dependencies learnt by the target model.

Ideally we want models which can utilize contextual cues when available, but are robust to variations in context and can detect and segment objects even when they appear out of context. However, machine learning based vision models widely used today are biased to the data they see frequently in training and tend to perform poorly on less frequent situations, for example the object-without-context and context-without-object scenarios we are interested in. We address this by proposing a data augmentation scheme to expose the image classification and segmentation models to different contexts during training, in-order to improve robustness of the model to context changes. This is done by removing selected objects from images and training the models on the edited images to recognize and segment other objects in the image, even with contextual objects removed. Our experiments show that the classification and segmentation models trained with this data augmentation scheme are less sensitive to context changes and perform better on real out-of-context datasets, while preserving the baseline performance on the regular data splits.

To summarize, the main contributions of this paper are as follows: a) We propose an object removal based method to understand and quantify sensitivity of vision models to context, b) We apply this to analyze image classification and segmentation models and find some interesting and undesirable dependencies learnt by the models between classes and contextual objects and c) We propose a data augmentation scheme based on object removal to make the models more robust to contextual variation and show that it helps improve performance in out-of-context scenarios.

2. Related work

The importance of semantic context in visual recognition is a well established with studies showing context can help humans recognize objects faster e.g. when dealing with difficult low resolution images [15, 1]. In computer vision, incorporating context information has been shown to im-

prove performance in various tasks including object recognition [12, 21, 17] and action recognition [9], object detection [3] and segmentation [24]. Earlier approaches built explicit context models for example by incorporating co-occurrences [17] and spatial location statistics [6]. However, recently, explicit context modeling has been replaced the use of deep convolutional neural network (CNN) encoders which summarize the information in the whole image into compact features. Classification and segmentation models, built on top of these deep feature encoders, can exploit information about object and context to achieve good performance [11, 7, 14]. Approaches to improve the use of context in CNNs have been explored including using spatial pyramids [25], atrous convolutions [4] and learning context encoding with a separate neural network [24]. While this implicit context encoding with deep CNNs gives good performance, it is less interpretable and it makes it hard to know whether the models are basing their decisions on visual or contextual evidence. Recent works propose to address this with methods to interpret neural networks by visualizing salient regions for classification decision [18, 23], quantifying how interpretable individual units in network are [2]. While these works focus on interpreting the internal representations of the network, we look at quantifying the context sensitivity of models from the input data perspective and treat the networks as black boxes. By manipulating the input image to remove objects and observing the network output, we quantify the sensitivity of classification and segmentation models to context and discover some interesting and undesirable dependency between classes. A recent work [19] takes a related approach. By adding few out-of-context objects into images they show that object detection networks are brittle to presence of out-of-context objects. However while the focus in [19] is to explore feature interference caused by out of context objects, we aim to quantify and mitigate contextual dependencies between classes. Recent work [8] proposes data augmentation scheme for object detection by adding objects into new contexts with contextual modeling to improve average case performance. In contrast, our work focuses on improving the contextual robustness of segmentation and classification models.

3. Quantifying the role of context

We quantify the contextual dependence of image classification and segmentation models by applying object removal. We propose metrics to measure the effects of context by measuring the change in the target model output for the original and the images edited to remove context objects. Next, we will discuss our removal model, define the robustness metrics and present the data augmentation strategies to reduce the contextual dependence and improve performance in out-of-context situations.

3.1. Object removal

To create edited images with context objects removed, we need a fully automatic object removal model. For this, we utilize ground-truth object masks to remove the desired object and use an in-painting network to fill in the removed region. We base our in-painting network on the model proposed in [20], since this inpainter is directly optimized for removal, and can better handle irregular masks [20], regularly encountered in object removal. More details about the network architecture can be found in the supplementary material. The above removal method works well for medium sized objects, but struggles for large objects since then the in-painter needs to synthesize most of the image. Hence, we impose size restrictions on the objects we choose to remove to be less than 30% of the image. In the classification scenario on the COCO dataset, we consider all 80 object categories for removal. In the segmentation setting on the ADE20k dataset, we consider only the non-stuff categories (90 categories) for removal and measure the effects of removing these objects on the segmentation of all 140 categories. The stuff categories include objects like road, sky and field which are typically very large and hard to inpaint and hence are excluded from removal. An important point to note here is that the in-painter model is not aware of the downstream models and is not optimized to fool or change their decisions. The effects of the in-painter are local and only affects the region the object is removed from. Qualitative examples in Figures 2 and 3 show that the in-painting works reasonably well in the object removal setting.

3.2. Measuring context dependency

To understand the effect contextual cues have on image-classification and segmentation models, we test them models on edited images where a context object has been removed. Precisely, given an original image I containing a set of objects $C = \{c_1, c_2 \dots c_n\}$, we first create a set of edited images $\mathcal{I}_e = \{I - c_i | c_i \in C \text{ and removable}(c_i)\}$. Next we test classification and segmentation models on I and \mathcal{I}_e and the check their outputs for consistency with the performed removal as described below for each task.

Image-level classification. Given a trained classifier S_{c_i} for class c_i , we will now characterize how robust it is to changes in context of c_i . We first obtain classifier scores for the original image I , edited image $I - c_i$ with object c_i removed and for the edited set $\mathcal{I}_{owc} = \{I - c_j : c_j \in I, j \neq i\}$, all of which contain the object c_i but have one context object removed. Ideally, if the classifier S_{c_i} is robust to context changes it should score all the images in \mathcal{I}_{owc} higher than the image $I - c_i$, since $I - c_i$ does not contain the object c_i and the images in \mathcal{I}_{owc} do. Precisely, a classifier robust to context should satisfy the below in-equality:

$$S_{c_i}(I_{owc}) \geq S_{c_i}(I - c_i), \forall I_{owc} \in \mathcal{I}_{owc} \quad (1)$$

We can count the number of times this condition is violated to quantitatively measure the robustness of the classifier.

$$V^{\min}(c_i) = \frac{\sum_I \mathbb{1}[(\min_{I_{owc}} S_{c_i}(I_{owc})) < S_{c_i}(I - c_i)]}{\sum_I \mathbb{1}[c_i \in I]} \quad (2)$$

$$V^{\text{mean}}(c_i) = \frac{\sum_I \mathbb{1}[\mathbb{E}_{I_{owc}} [S_{c_i}(I_{owc})] < S_{c_i}(I - c_i)]}{\sum_I \mathbb{1}[c_i \in I]} \quad (3)$$

where $\mathbb{1}$ is the indicator variable. $V^{\min}(c_i)$ is a strict metric counting instances classifier scores $I - c_i$ higher than any of the edited images, whereas $V^{\text{mean}}(c_i)$ is a softer metric counting instances where $I - c_i$ is scored higher than the average score assigned to the edited images.

Semantic segmentation. To understand the role context plays in this pixel-level labeling task, we analyze the behaviour of a trained segmentation model by removing one object at a time from the original image. Specifically, we measure how the segmentation correctness of the rest of the image changes (as compared to segmentation of the original image) when we remove an object from the original image. Given a segmentation model P , we compute the intersection-over-union (IoU) for a class c_i (w.r.t. ground-truth) on the original image I and edited image $I - c_j$. If the IoU value changes more than threshold α , we consider the segmentation prediction for class c_i to be affected by removal of c_j . Counting these violations we get,

$$AR(c_i, c_j) = \frac{\sum_I \mathbb{1}[|\Delta \text{IoU}_{c_i c_j}| \geq \alpha]}{\sum_I \mathbb{1}[c_i, c_j \in I]} \quad (4)$$

where $\Delta \text{IoU}_{c_i c_j}$ is the change in IoU of class c_i with removal of object c_j and α is the change threshold. The matrix $AR(c_i, c_j)$, represents the fraction of images where removing the object c_j , affects the segmentation of the object c_i with high values of $AR(c_i, c_j)$ indicating that the segmentation model depends heavily on the presence of the context object c_j to segment c_i .

3.3. Data augmentation with object removal

We now present our data augmentation solution to reduce the sensitivity of classification and segmentation models to context distribution. The main idea is to expose these models to training images of object-without-context and context-without-object scenarios. This will help the models deal with the lack of contextual information and hence become more robust to context changes. For this we perform object removal to create edited images with some objects removed and add these edited images to the training batch. Specific details of how to pick objects for removal and how to use them in training for the two tasks are discussed below.

Classification. For classification we experiment with two strategies to use the edited images in training. In the first approach, we refer to as *Data-aug-rand*, we randomly select one object to remove with uniform probability. Edited



Figure 2: Context violations by image-level classifier. The primary object is marked with blue box and the context object is marked with magenta. The first column shows the original image, middle shows the image with only object and the third with only the context. We see that the baseline classifier depends heavily on the context and always scores the context only images (last column) higher than the image with only the primary object (middle column). The data augmented model does better and gets the ordering right.

image is assigned the same labels as the same as the original image excluding the removed object class. Now we train the classifier with original and data augmented images using simple binary cross-entropy loss. In the second approach referred to as *Data-aug-const* we explicitly optimize for robustness by including the inequality in (1) in the loss function. To do this, for a randomly selected images in the training batch, we create the full edited image set $\{I - c_i : c_i \in I\}$. Then we can incorporate the robust-

ness constraint as a hinge loss with final loss being weighted sum of cross-entropy loss and the hinge loss L_h .

$$L_h(I) = \sum_{c_i \in I} \max \left[0, S_{c_i}(I - c_i) - \min_{c_j, j \neq i} S_{c_i}(I - c_j) \right] \quad (5)$$

Segmentation. We also perform data augmentation on the segmentation task by creating edited images by selectively removing objects. The edited images can be used in training the segmentation model in two ways. First we can simply ignore the removed pixels and train the model to predict the original ground-truth labels on the rest of the image (*Ignore*). This helps the model learn that the labeling of a pixel should not be affected by the removal of a context object. Alternatively, we can explicitly tell the model that the removed object is not present by minimizing the likelihood assigned to the removed class at the edited pixel locations (*Negative loss*). The next question is how to sample the objects to remove. We explore three different strategies to select these objects. The first strategy, *Random*, selects one random object to remove from the objects present in the image with uniform probability. However, sometimes the *Random* strategy can select very large object for removal, which can harm the quality of the edited image. To address this we use the *Sizebased* strategy, which select objects based on their relative sizes in the image, giving higher probability to selecting smaller objects. The probability for picking an object is computed as $p(c_i, I) \propto \frac{\sum_{c_i \in I} a(I, c_i)}{a(I, c_i)}$ where $a(I, c_i)$ is the area of the class c_i in image I . We also explore a hard negative mining based strategy, where we create harder training examples for the segmentation model by removing easy classes. This allows the model to focus on segmenting the harder classes while also becoming robust to context. Concretely, in this *HardNegative* strategy we monitor the average cross-entropy segmentation loss $l_{\text{avg}}(c_i)$ for an object class c_i and calculate the probability of removal of c_i as inversely proportional to $l_{\text{avg}}(c_i)$.

4. Experiments and Results

This section presents the results of our analysis of how much the contextual information influences the performance of image classification and segmentation models. Using the robustness metrics defined in Section 3.2, we discover that the classification predictions on many well-performing classes are sensitive to context, and perform poorly on object-without-context and context-without-object images. Similar results are also found in the segmentation setting with the model depending heavily on context objects to correctly segment classes like *road*, *sidewalk*, *grass*. We also present results from our data-augmentation strategies, which help reduce this context dependence and improve robustness, without sacrificing performance.

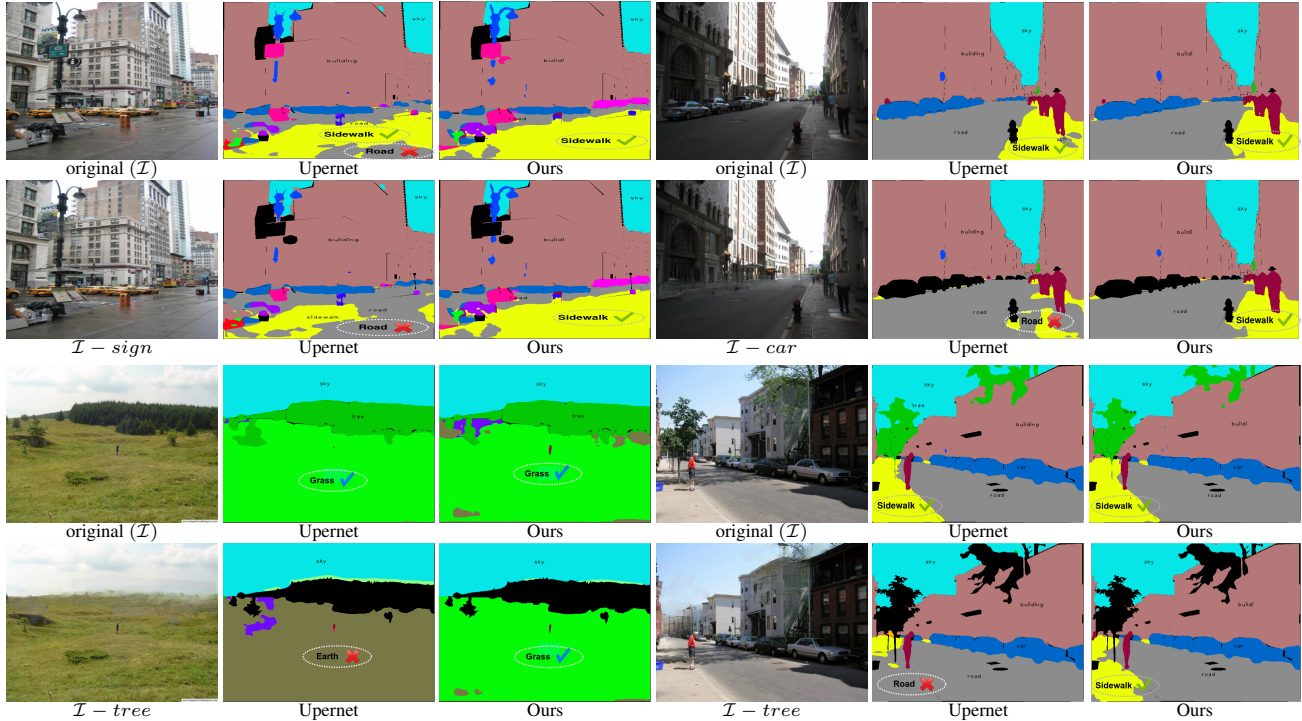


Figure 3: Examples of segmentation failures due to removal of a single context object. We see the segmentation of road, sidewalk and grass affected significantly when context objects like signboard, car and tree is removed (comparing odd and even rows). Model trained with proposed data-augmentation is more robust to these changes.

4.1. Image level classification

4.1.1 Experimental setup for classification

Training data. We run our classification experiments on the COCO dataset [10], which contains 80 labeled object classes in their natural contexts. The dataset also has bounding box and segmentation annotation for each object. We use image-level labels to train the classifiers and use the object segmentation masks to test them with object removal.

Out-of-context testing. Apart from testing the classifier models on regular COCO data we conduct additional experiments to quantify the performance in out-of-context scenarios with natural images. We divide the COCO images into two splits: the first split *Co-occur* with images having at least two objects in them and the second split *Single* with images containing a single object. The *Full* split is all images combining *Co-occur* and *Single*. The idea behind this splitting of the dataset is to separate out images where objects occur in their context (*Co-occur*) and images where object occur alone without the usual co-occurring context objects *Single*. Now we can train our models on the *Co-occur* split and test it on the *Single* split to measure, using only real images, how a classifier trained with only co-occurring objects performs when objects appear without the context seen in training. Additionally we also test our COCO trained models on the *UnRel* dataset [16] which

contains natural images with objects occurring in unusual contexts and relationships. We keep the classes which map to one of the 80 object classes in COCO, leaving 29 classes and 1071 images in the *UnRel* dataset.

Baseline classifier. The image-level classification model we test is based on the architecture proposed in [14]. It consists of a Imagenet [5] pre-trained VGG-19 network for feature extraction network followed by two convolution layers, global max-pooling layer and a linear classification layer with sigmoid activations. The model is trained with binary cross-entropy loss. We train and test the model at single scale at 256x256 resolution, to simplify the analysis. Our classifier achieves similar mAP on real coco data as reported in [14], with our mAP slightly lower (0.600 vs 0.628 in [14]) due to single scale training and testing.

4.1.2 Analyzing classifier robustness to context

To measure the robustness of the trained classifier to context, we test it on real images and edited images and compute the robustness scores V^{\min} and V^{mean} as described in Section 3.2. Table 1 shows the robustness scores averaged over all classes computed on the COCO test along with the standard performance metric mean average precision (mAP) for the baseline classifier (first row). We can see that, despite achieving good mAP (0.6), the baseline clas-

Model	Training Data	COCO test set			Robustness Metrics		UnRel dataset ↑
		Full ↑	Co-occur ↑	Single ↑	V^{\min} ↓	V^{mean} ↓	
Baseline	Full (39k)	0.60	0.57	0.62	34%	24%	0.50
Data-aug-rand	Full (39k)	0.61	0.58	0.65	32%	22%	0.54
Data-aug-const	Full (39k)	0.60	0.58	0.63	25%	14%	0.52
Baseline	Co-occur (30k)	0.56	0.55	0.58	34%	24%	0.46
Data-aug-rand	Co-occur (30k)	0.58	0.57	0.60	31%	21%	0.49
Data-aug-const	Co-occur (30k)	0.58	0.57	0.60	27%	15%	0.51

Table 1: Effect of data augmentation on classification model

Model	all (407 images)		with car (258)		without car (149)	
	Road	Sidewalk	Road	Sidewalk	Road	Sidewalk
Upernet	0.81	0.59	0.86	0.67	0.68	0.40
DataAug	0.82	0.60	0.86	0.65	0.72	0.46

Table 2: Comparing the performance of road and sidewalk segmentation on natural images with and without cars.

sifier trained on full data performs poorly in-terms of robustness metrics. In about 34% of cases the model violates the context consistency requirement of (1). This means in 34% cases, the classifier scores images without the target object higher than an image where object is present but a context object has been removed. Comparing the per-class robustness score, $V^{\min}(c_i)$ and the per-class average precision (AP) (see supplementary for visualization), we see that good performance in AP does not mean the classifier is robust to context. Many classes like mouse, keyboard, sink, tennis racket etc, which are performing well in AP (≥ 0.8), but have poor robustness to changes in context ($V_o^{\min} \geq 50\%$). In extreme case, the *mouse* classifier violates the consistency in more than 90% of cases, despite having very good AP (0.88). This indicates that the classifiers are relying too much on contextual evidence to detect the objects but perform poorly when tested on images where the context distribution is different from training.

We visualize these violations in Figure 2. In the first row we see the case where the keyboard classifier scores the image with the keyboard removed higher (4.67) than the image with the keyboard but with the monitors removed (1.99). Similarly we see the skateboard and the frisbee classifiers relying on person to hallucinate the respective objects. The violations shown in the first three rows of Figure 2 occur in objects with high co-occurrence dependence with other classes. However, such context violations can also be seen in classes like *person* which occur in very different contexts as seen in the last row of Figure 2. Here, the violation occurs in a difficult image where the *person* is small, but a more distinct class with co-occurrence dependence on person is clearly visible (*kite*). The classifier seems uses the *kite* context class to hallucinate that there is a *person*, even when the *person* has been removed.

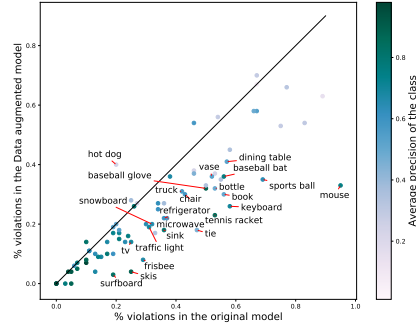


Figure 4: Comparing the % of violations in different classes with and without data augmentation. Points below the diagonal line show improvement with data-augmentation and the ones above degrade. The colors denote the average precision.

4.1.3 Data augmentation to improve robustness

We train two variants of the data-augmented image classification models as described in Section 3.3. The first *Data-aug-rand* learns with standard cross-entropy loss on the edited images with a random object removed and the second *Data-aug-const* which is optimized directly for robustness using a set of edited images and hinge loss.

Quantitative results. We present the evaluation of the data-augmented and the baseline models in Table 1. On models trained with *Full* training data, the data-augmented model *Data-aug-rand* provides a small improvement in overall mAP on the COCO test set (0.61 vs 0.60). However measuring the performance on the two splits *Co-occur* and *Single* reveals that the improvement is significant on the *Single* split (0.65 vs 0.62), indicating that the data augmentation helps the classifier better deal with out of context objects. This is also seen when comparing the performance of the two models on the UnRel dataset, where *data-aug-const* significantly improves over the baseline model (0.54 vs 0.50). This improved robustness of the data augmented classifier to context changes is also measured by our robustness metrics V^{\min} and V^{mean} . *Data-aug-rand* classifier makes overall 2% less violations under both worst-case (V^{\min}) and average-case (V^{mean}) context changes. Directly optimizing the robustness constraints allows the model *Data-aug-const* to significantly improve upon the baseline model in robustness metrics, while still obtaining improvement in the performance metrics. It exhibits much less worst-case (25% vs 34% for baseline) and average-case violations (14% vs 24% for baseline), while improving the performance in the UnRel dataset (0.52 mAP vs 0.50 for baseline). The benefit of optimizing for robustness is clearly seen when we constrain the training data to the *Co-occur* set, where the classifier never sees objects alone. Baseline model trained on the *Co-*

occur set drops in performance on the *Single* (0.58 from 0.62 on when trained on *Full*) and the UnRel test sets (0.46 vs 0.50 with *Full*). However, with data augmentation and enforcing robustness constraints, we can recover some of this performance. On the *Single* test set *Data-aug-const* model trained on *Co-occur* set gets 0.58 mAP compared to 0.60 by baseline model trained on full data and even surpass it on the UnRel test set with 0.51 mAP. This shows that the data augmented model is able to overcome the contextual bias in the training set and perform well in unseen contexts.

When we compare the per-class robustness metrics between regular and data augmented models (data-aug-const), as shown in the Figure 4, we see that data-augmentation significantly reduces the worst case violations (V^{\min}) on well-performing classes. For example, V^{\min} drops from 95% to less 36% for the mouse class and from 58% to 28% for the keyboard class. The effect of this increased robustness is seen in qualitative examples in Figure 2. For example, in the first row the baseline keyboard classifier gives too much weight to evidence from *monitor* and scores the image with only *monitor* higher than the only keyboard. However, the data augmented model correctly orders the images.

4.2. Semantic segmentation

So far, we have seen that multi-label classification models suffer from sensitivity to context, with classifiers often mixing up contextual and visual evidence. Next we will measure the context sensitivity of models in a more local and strongly supervised task of semantic segmentation.

4.2.1 Experimental setup for segmentation

Training and test data. We conduct our semantic segmentation experiments primarily on the ADE20k dataset [26] containing 140 categories of labeled objects, in different settings. Some of the 140 classes are typical background classes like *sky*, *sea* and *wall* and are large and difficult to in-paint and are hence excluded from removal.

Out-of-context testing. Following the process in image-level classification, we also measure the performance of the segmentation models on real out-of-context data. This is done in two ways. First, we train the segmentation model in a restricted setting with only three classes *car*, *road* and *sidewalk*. Now, we can again make two splits of the training and testing images into the *Co-occur* split of images with at-least two objects (3317 images) and the *single* split with only a single object (1693 images). Then we train the segmentation models on *co-occur* split and test on *single* split to see how well it can perform segmentation without context. Additionally we also test the models trained with ADE20k data on the Pascal-context dataset [13] in order to measure the performance under a different context distribution. This is done by manually mapping the 59 labels

in the pascal-context to ADE20k labels and restricting the segmentation model to produce only the mapped labels.

Baseline segmentation model. We use the recent UperNet [22] model, with good results on the ADE20k, as our baseline segmentation model. We train the variant with the Resnet-50 encoder and a Upernet decoder with batch size of 6 images (maximum that fit in GPU) and with the default hyper-parameters suggested by the authors. This model achieves mean intersection-over-union (mIoU) of 0.377 and accuracy of 78.19% with single scale testing.

4.2.2 Context in semantic segmentation

We analyze robustness of the segmentation models to context by removing objects and computing the matrix $AR(c_i, c_j)$ presented in Section 3.2, which measures the % of images where removal of object c_j significantly affects segmentation of object c_i . The matrix $AR(c_i, c_j)$ we obtain for the Upernet model in ADE20k dataset is a sparse matrix with sharp peaks (see supplementary for a visualization). This indicates that the classes depend on specific context objects and are significantly affected by their removal. The sparsity also indicates that the effects on the segmentation are due the class being removed and not in-painting artifacts (otherwise the segmentation would be affected by all removal). Some of dependencies we discover in $AR(c_i, c_j)$ are reasonable and harmless, for example between *pot* and *plant* ($AR = 50\%$). Once you remove the *plant*, *pot* looks more like a *trash can* and the segmentation model often flips the label to *trash can*. However other dependencies are spurious and not desirable. For example, we notice that often the segmentation model uses presence of *car* to differentiate between *road* and *sidewalk*. Removing *car* affects the IoU of the *road* and *sidewalk* in 21% and 22% of cases respectively. This dependence is undesirable, and can be catastrophic in applications like self-driving cars.

We show qualitative examples where removal affects segmentation of Upernet model in Figure 3. The first two rows show the cases where removal of an object negatively impacts the segmentation of other objects. This include cases where removal of *street sign* and *car* severely affects segmentation of *road* and *sidewalk*, and a case where removal of *trees* affects segmentation of *grass*. We can see from these examples that while edit on the image is small and local, the effects of this removal on segmentation prediction is not local. Removal of a small objects can have drastic effects on segmentation in a far-away region.

4.2.3 Data augmentation for segmentation

Next we will look at the results of using data-augmentation for segmentation models. For this purpose we train the Upernet [22] based data-augmented models on the ADE-

Model	Removed Pixels	ADE20k	
		mIoU	Acc
Upernet[22]	-	0.377	78.31
DA (random)	Ignore	0.320	75.2
DA (sizebased)	Ignore	0.379	78.31
DA (hard negative)	Ignore	0.375	77.8
DA (sizebased)	Negative	0.377	78.25
DA (hard negative)	Negative	0.385	78.47

Table 3: Data augmentation results on ADE20k dataset

20k dataset with on three different strategies for selecting the object to remove as discussed in Section 3.3.

Quantitative results. Table 3, shows the results comparing the data-augmented models with the baseline Upernet model. We can see that random sampling strategy, which worked well in image classification, fails here leading to drop in performance. This is because, many object categories in ADE20k dataset are large and difficult to remove like bed, sofa and mountain and random strategy suffers by picking these. Instead when we switch to size-based and hard-negative based sampling, we see that the performance improves and the size-based sampling model achieves the best mIoU of the three models (0.379). Applying negative likelihood loss on the removed object class gets further improvement when combined with hard negative sampling. This model also improves upon the Upernet baseline (achieving 0.385 IoU vs 0.377 by Upernet), despite the fact that the removal based data-augmentation is designed to make the model more robust to contextual variations.

To understand how data-augmentation impacts sensitivity to context, Figure 5 visualizes the maximum sensitivity of a class to removal of other classes, $\max_{c_j} AR(c_i, c_j)$ for different classes with and without data-augmentation. We see that for majority of classes robustness to context improves with data augmentation. For example *pillow* class is only affected 32% of the time with context changes, compared to 53% before data augmentation. Similarly, *road* and *sidewalk* classes are only affected 9% and 14% of the time respectively, compared to 21% and 22% before. This improved robustness translates into better generalization to real out-of-context data. We can see this in Table 2 where the performance of the *road* and *sidewalk* segmentation is measured on the validation set on images with and without cars. On the full set and on the split with cars, we see that the performance of the baseline Upernet and our augmented model (DA hard negative with negative loss) is equivalent. However, when we look at only images without car, the Upernet model performs significantly worse in both road (0.68 vs 0.72 for ours) and sidewalk (0.40 vs 0.46 for ours) segmentation. This quantitatively shows that the baseline model struggles to distinguish between *road* and *sidewalk*

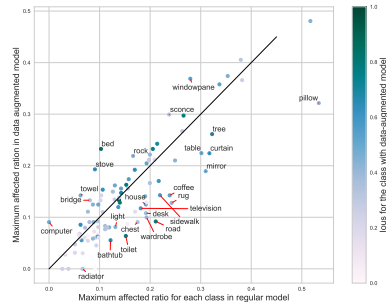


Figure 5: Comparing the context sensitivity of different classes with and without data augmentation with $\max_{c_j} AR(c_i, c_j)$ metric. Points below the diagonal improve with data-augmentation. The color denotes the mIoU.

without *car* in the image, whereas our data augmentation is more robust and performs well even without context (*car*).

We also see the benefit of data augmentation in experiments on restricted *Co-occur* training set and on the Pascal-context dataset. Our data augmented model outperforms the Upernet model (both trained on the ADE20k dataset) when tested on the Pascal-context dataset in both mIoU and pixel accuracy. While the Upernet model achieves mIoU of 0.284 and pixel accuracy of 61.3% our data augmented model achieves 0.293 and 62.10% respectively, indicating that it is able to generalize better when tested on a dataset with different context distribution than one seen during training. Table 4 presents the experiments with the *Co-occur* training set in the three class setting. First we can see that when we switch from training on *Full* training data to *Co-occur* split (containing only images with at least two objects), the performance of the Upernet greatly drops on the *Single* test split (from 0.67 to 0.52). This indicates that the model overfits to the context it sees, and is not able to segment objects when it seeing them out of context. However, with data-augmentation we generate images of objects without context, and can recover most of this performance loss (0.646). Surprisingly, data-augmented model trained on smaller *co-occur* data also outperforms the baseline trained with *Full* data when tested on the *co-occur* split.

Qualitative examples in Figure 3 also show the effect of increased robustness to context. While the baseline Upernet model is affected by context object removal causing drastic changes in predictions of other regions, our data augmented model is more stable. For example the removal of *sign-board*, *car* or *tree* does not effect the segmentation of the *road* or *sidewalk* by our model.

5. Conclusions

We have presented a methodology to analyze and quantify the context sensitivity of image classification and seg-

Model	Training Data	Full	Only Cooccur	Only Single
Upernet	Full (5k)	0.774	0.797	0.670
Data Aug	Full (5k)	0.742	0.754	0.675
Upernet	Co-occur (3.3k)	0.680	0.713	0.520
Data Aug	Co-occur (3.3k)	0.82	0.86	0.646

Table 4: Experiments in three class setting on ADE20k

mentation models, based on editing images to remove objects and measuring the effect on the target model output. Our analysis shows that despite good performance in-terms on mAP, classifiers for certain classes like keyboard, mouse, skateboard are very sensitive to context objects and perform poorly when seen out of context. In semantic segmentation setting, our analysis shows similar dependency between classes. For example we discover that the model depends on the presence of car to segment roads and sidewalk and fails drastically when the car is not present in the image. We present a data augmentation scheme based on object removal to mitigate this and make the classification and segmentation models more robust to context changes. Our experiments show that the proposed data augmentation scheme can help models generalize to out of context scenarios without losing performance in standard setting, indicating that the data augmented models better balance contextual and visual information.

References

- [1] E. Barenholtz. Quantifying the role of context in visual object recognition. *Visual Cognition*, 22(1), 2014. [2](#)
- [2] D. Bau, B. Zhou, A. Khosla, A. Oliva, and A. Torralba. Network dissection: Quantifying interpretability of deep visual representations. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. [2](#)
- [3] S. Bell, C. Lawrence Zitnick, K. Bala, and R. Girshick. Inside-outside net: Detecting objects in context with skip pooling and recurrent neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016. [1, 2](#)
- [4] 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. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 40, 2018. [2](#)
- [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. ImageNet: A Large-Scale Hierarchical Image Database. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009. [5](#)
- [6] C. Desai, D. Ramanan, and C. C. Fowlkes. Discriminative models for multi-class object layout. *International journal of computer vision*, 95(1), 2011. [2](#)
- [7] T. Durand, N. Thome, and M. Cord. Weldon: Weakly supervised learning of deep convolutional neural networks. In *CVPR*, 2016. [2](#)
- [8] N. Dvornik, J. Mairal, and C. Schmid. Modeling visual context is key to augmenting object detection datasets. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018. [2](#)
- [9] M. Jain, J. C. van Gemert, and C. G. Snoek. What do 15,000 object categories tell us about classifying and localizing actions? In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. [2](#)
- [10] 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 *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. [5](#)
- [11] J. Long, E. Shelhamer, and T. Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. [2](#)
- [12] M. Marszalek and C. Schmid. Semantic hierarchies for visual object recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2007. [2](#)
- [13] 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 (CVPR)*, 2014. [1, 7](#)
- [14] M. Oquab, L. Bottou, I. Laptev, and J. Sivic. Is object localization for free?-weakly-supervised learning with convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015. [2, 5](#)
- [15] D. Parikh, C. L. Zitnick, and T. Chen. Exploring tiny images: The roles of appearance and contextual information for machine and human object recognition. *IEEE transactions on pattern analysis and machine intelligence*, 34(10), 2012. [1, 2](#)
- [16] J. Peyre, I. Laptev, C. Schmid, and J. Sivic. Weakly-supervised learning of visual relations. In *ICCV*, 2017. [5](#)
- [17] A. Rabinovich, A. Vedaldi, C. Galleguillos, E. Wiewiora, and S. Belongie. Objects in context. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2007. [2](#)
- [18] M. T. Ribeiro, S. Singh, and C. Guestrin. Why should i trust you?: Explaining the predictions of any classifier. In *Proceedings of the ACM SIGKDD international conference on Knowledge discovery and data mining*, 2016. [2](#)
- [19] A. Rosenfeld, R. Zemel, and J. K. Tsotsos. The elephant in the room. *arXiv preprint arXiv:1808.03305*, 2018. [2](#)
- [20] R. Shetty, M. Fritz, and B. Schiele. Adversarial scene editing: Automatic object removal from weak supervision. In *Advances in Neural Information Processing Systems (NIPS)*, 2018. [3, 10](#)
- [21] A. Torralba, K. P. Murphy, and W. T. Freeman. Using the forest to see the trees: exploiting context for visual object detection and localization. *Communications of the ACM*, 53(3), 2010. [1, 2](#)
- [22] T. Xiao, Y. Liu, B. Zhou, Y. Jiang, and J. Sun. Unified perceptual parsing for scene understanding. In *Proceedings*

of the European Conference on Computer Vision (ECCV), 2018. 1, 7, 8, 12

- [23] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014. 2
- [24] H. Zhang, K. Dana, J. Shi, Z. Zhang, X. Wang, A. Tyagi, and A. Agrawal. Context encoding for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2018. 1, 2
- [25] H. Zhao, J. Shi, X. Qi, X. Wang, and J. Jia. Pyramid scene parsing network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [26] B. Zhou, H. Zhao, X. Puig, S. Fidler, A. Barriuso, and A. Torralba. Scene parsing through ade20k dataset. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017. 7

Appendix A. Object removal model

To remove objects we use the ground truth segmentation masks and dilate them by a small factor (5 in the coco dataset and 7 in the ADE20k dataset). This dilated mask is multiplied with the input image to remove the target object from the image. Then the masked image and the mask is passed to an in-painting network which fills the masked area with a plausible background texture. We use the in-painting architecture and the training procedure proposed in [20]. Table 5 and 6 present the detailed architecture of the in-painting network. We will make the code and pre-trained in-painter models available after the review process.

Appendix B. Analyzing robustness to context

In the main paper, we presented our analysis showing that the classification and segmentation models are sensitive to context and their predictions are significantly affected when presented with edited images with context objects removed. In the following sub-sections we present additional visualizations to support these arguments.

B.1. Image-level Classification

Co-occurrence of objects. An important factor which causes the image-level classification models to use contextual dependencies is the co-occurrence distribution of objects. Many objects in COCO have a strong co-occurrence relation with other objects. We quantify this using the normalized co-occurrence counts for each object with others given by

$$NC(c_i, c_j) = \frac{\text{Count}(c_i \cap c_j)}{\text{Count}(c_i)}$$

. This matrix is visualized in Figure 6. $NC(c_i, c_j)$ takes value between 0 and 1 and represents the fraction of images containing object c_i , which also contains object c_j . We can see that, for classes like skateboard, surfboard, tennis racket and handbag, this ratio is very high ($\geq 90\%$) with the person

Masked Image + mask
Conv 4x4, 64 filters, stride 1
Conv 4x4, 128 filters, stride 2
Conv 4x4, 256 filters, stride 2
Conv 4x4, 512 filters, stride 2
Residual Block, 256 filters
Residual Block, 256 filters
Residual Block, 256 filters
Residual Block, 256 filters
Residual Block, 256 filters
Residual Block, 256 filters
Upsample + Conv 3x3, 256 filters
Upsample + Conv 3x3, 128 filters
Upsample + Conv 3x3, 64 filters
Conv 7x7, 3 filters, stride 1

Table 5: In-painting model architecture starting with input in the first row to the output layer in the last. Each convolutional layer is followed by a Instance Norm layer and a Leaky Relu non-linearity with slope 0.1

Conv 3x3, n filters, stride 1
Instance Norm
Leaky Relu (slope 0.1)
Conv 3x3, n filters, stride 1
Instance Norm
Leaky Relu (slope 0.1)

Table 6: Architecture of the residual block with n filters

class, since these classes often occur with a person holding or riding them. We also see that the matrix is not symmetric. This is because, while the skateboard might occur always with a person, but person class occurs in various contexts without skateboard. However, for some groups of objects like mouse, keyboard and monitor, and spoon, fork, and cup have symmetric co-occurrence relationship.

For many categories, including the cases discussed above, the co-occurrence ratio is very high ($> 60\%$). This causes problems for object classifiers of these categories, as we see in the analysis presented in the main paper. When a small or difficult to detect object class like mouse or skateboard, frequently co-occurs with a more easy to detect object class like monitor or person, the classifiers tend to overuse the contextual relationship for making their classification decisions instead of visual evidence for the object of interest. This leads to failures when the context is different or the object occurs without context.

Relation of performance to robustness. As discussed in section 4.1.2 in the main paper, we find that many well-performing object classes in terms of average precision (AP) perform poorly in terms of robustness. To show this we plot the per-class average precision against the worst-case robustness metric $V^{\min}(c_i)$ in Figure 7. We can see for example that classes like mouse, tennis racket, sports ball, baseball bat and book which have high AP ($g_{eq} 0.6$)

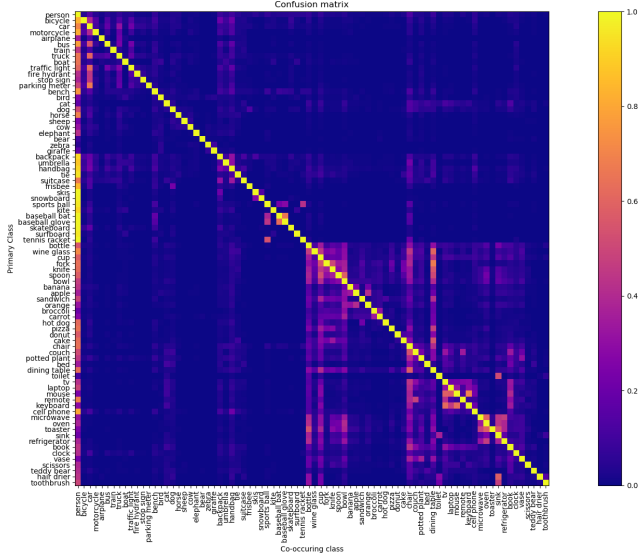


Figure 6: Co-occurrence ratio, $N(c_i, c_j)$ of objects on the COCO dataset

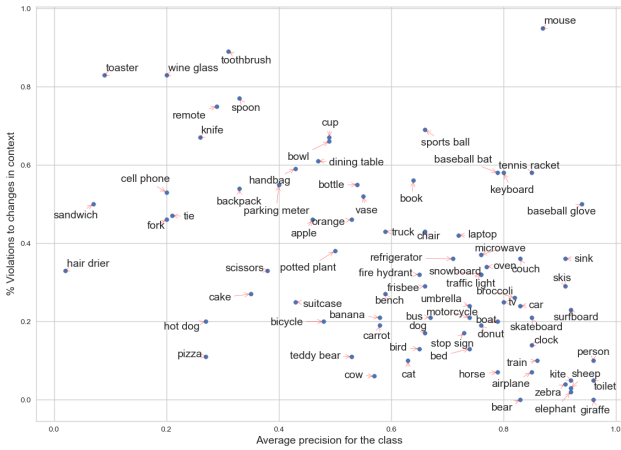


Figure 7: Comparing class-wise average precision to the % of violations in changes to context. Many well-performing categories (high mAP), have high percentage of violations, including mouse, tennis racket, keyboard, book, and sink.

have poor robustness ($V^{\min}(c_i) \geq 0.5$). In all these cases, the classifier seems to predominantly use contextual objects to make their predictions and achieve high average precision. But they fail when presented with object-without-context and context-without-object images, usually scoring the context-without-object images higher. This is also seen in further visual examples presented in Figure 8. Interestingly visually distinct classes like zebra, elephant, giraffe achieve high AP, while also being robust as seen in Figure 7

Object without Context Context without Object

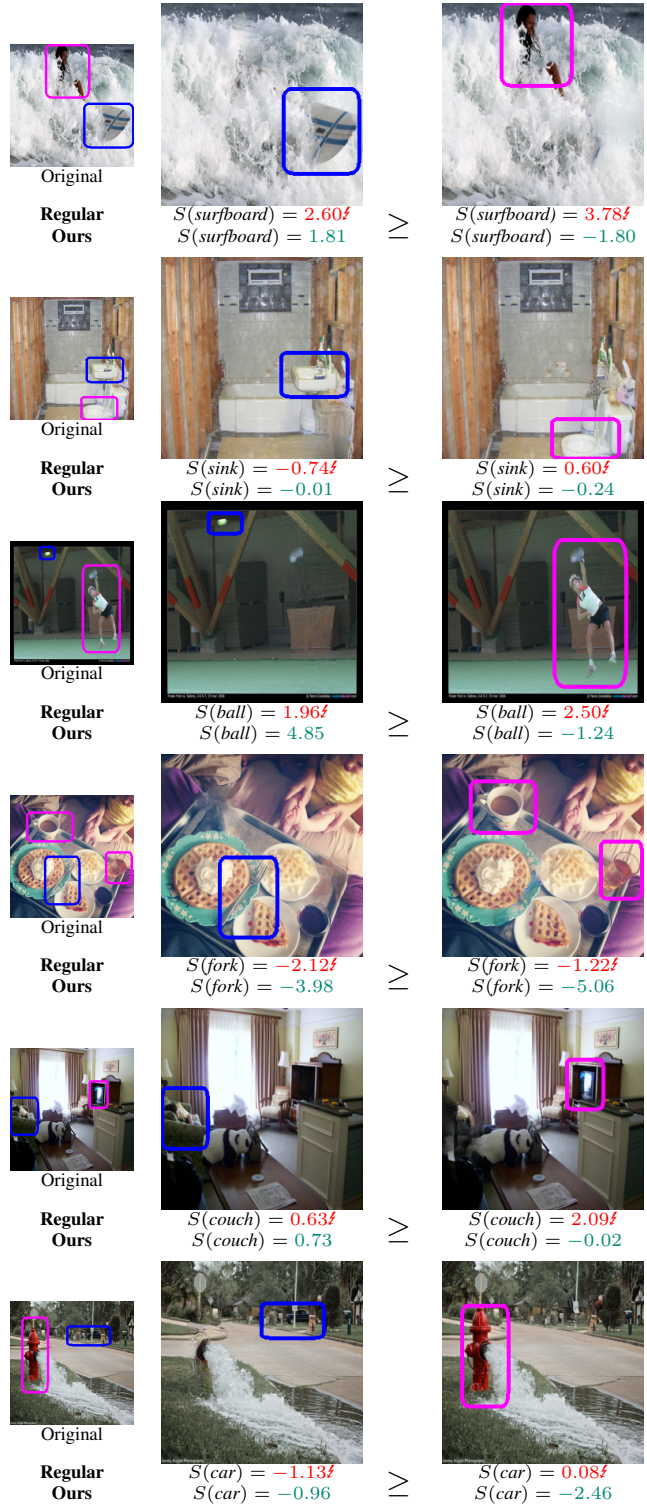


Figure 8: Context violations by image-level classifier. The primary object is marked with blue box and the context object is marked with magenta. The first column shows the original image, middle shows the image with only object and the third with only the context. We see that the baseline classifier depends heavily on the context and always scores the context only images (last column) higher than the image with only the primary object (middle column). The data augmented model does better and gets the ordering right.

B.2. Semantic Segmentation

Visualizing robustness metrics. We compute the robustness metric $AR(c_i, c_j)$, which measures the ratio of instances when segmentation of class c_i is affected by the removal of class c_j , in the ADE20k dataset for the Upernet [22] model. This is visualized in Figure 9. The y-axis is the affected class and the x-axis is the removed object class. We show the rows and the columns which have atleast one entry > 0.1 , for readability. We can see that the $AR(c_i, c_j)$ matrix is very sparse, indicating that the segmentation is not affected by all removal, but of only specific classes. As discussed in section 4.2.2 of the main paper, we can see that classes like road and sidewalk depend on the class car. The sidewalk is also to an extent affected by removal of trees.

To measure the direction of the effect, that is if removal of context harms or improves the segmentation of a class, we visualize the average change in IoU in Figure 10. Surprisingly, we find that not all context removal negatively affects the segmentation. Sometimes removing an object helps the model to resolve ambiguities and fix the segmentation of other objects. We can see in Figure 10 that while majority of change is negative, for a few pairs of objects removal positively affects the IoU. For example removing *lamp* class improves the segmentation of ceiling *light*. Similarly, removing *armchair* improves segmentation of *chair* and *sofa* classes, since the ambiguity is resolved.

Ablation on data augmentation. To understand if removal and in-painting is really needed for data augmentation, we conduct two ablation studies which are presented in Table 7. First we train a version of the baseline model where for each training sample we randomly select an object and set its label to 'ignore'. We use the same sampling strategy as *sizebased* data augmentation model. Hence this mimics exactly the training procedure in DA (sizebased), except without actually removing the object. Comparing the results of this model (No removal (sizebased)) to the data augmented version, we see that removing the object is necessary and simply ignoring the label leads to a severe drop in performance (0.354 vs 0.379). Similarly, in the second experiment we train a model with the object removed but without in-painting. In this case we can see from Table 7 that, having in-painter during augmentation is slightly better than the model without (0.379 vs 0.375).

Further visual examples of the sensitivity of the baseline Upernet [22] model to contextual changes and the robustness provided by data-augmentation is seen in Figure 11.

Confirming the source of sensitivity. Finally we conduct an additional experiment to verify that the volatility we see in the output of the segmentation models on edited images are due to removal of context objects and not due to editing artifacts. To test this we observe the predictions of the model on three version of the input image. First is the orig-

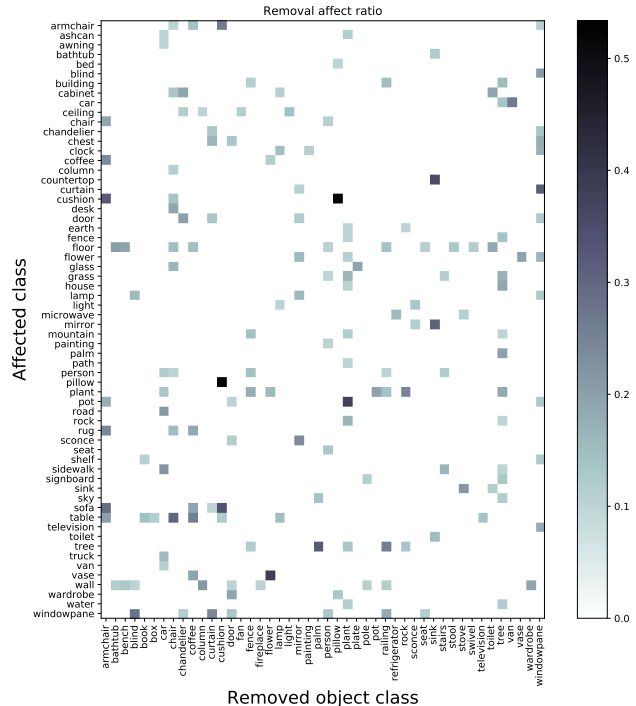


Figure 9: Visualizing frequency with which classes are affected by removal of other objects. Y-axis are the affected objects and the x-axis shows the removed objects

Model	Removed Pixels	ADE20k	
		mIoU	Acc
Upernet[22]	-	0.377	78.31
No removal (sizebased)	Ignore	0.354	77.45
No inpainter (sizebased)	Ignore	0.375	78.25
DA (sizebased)	Ignore	0.379	78.31

Table 7: Data augmentation results on ADE20k dataset

inal image. Second is the image with a context object removed. Finally the last image is the false edited image created by masking and in-painting the input image with the same mask as the context object, except with a horizontal flip. Thus the context object is not removed in the false edited image, but a similar shape and size region is removed and in-painted in a different part of the input image. All three images are fed to the segmentation models and the output is shown in Figure 12. We can see that the Upernet model output is virtually identical on the original image and the false edited image(third row). However the segmentation on the edited image with context object removed is significantly different (second row). This indicates that the segmentation models are not affected by the editing artifacts but by the removal of the context objects as claimed in the paper.

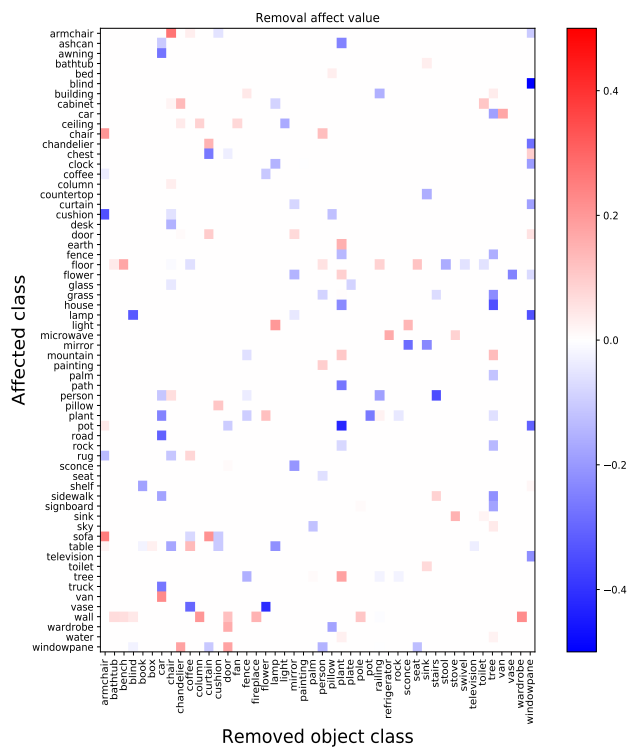


Figure 10: Visualizing the mean change in the IoU of object segmentation with context object removal. Y-axis are the affected objects and the x-axis shows the removed objects

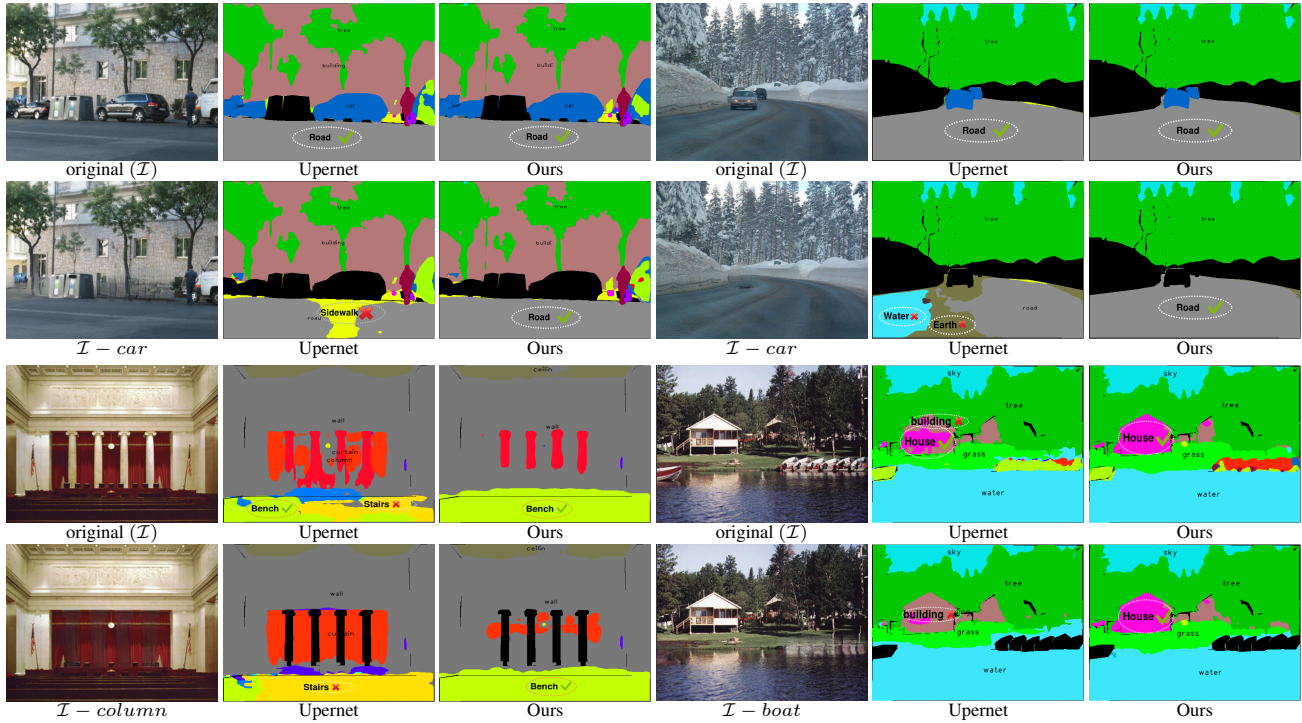


Figure 11: Examples of segmentation failures due to removal of a single context object. We see the segmentation of road, bench and house affected significantly when context objects like car, columns and boat is removed (comparing odd and even rows). Model trained with proposed data-augmentation is more robust to these changes.

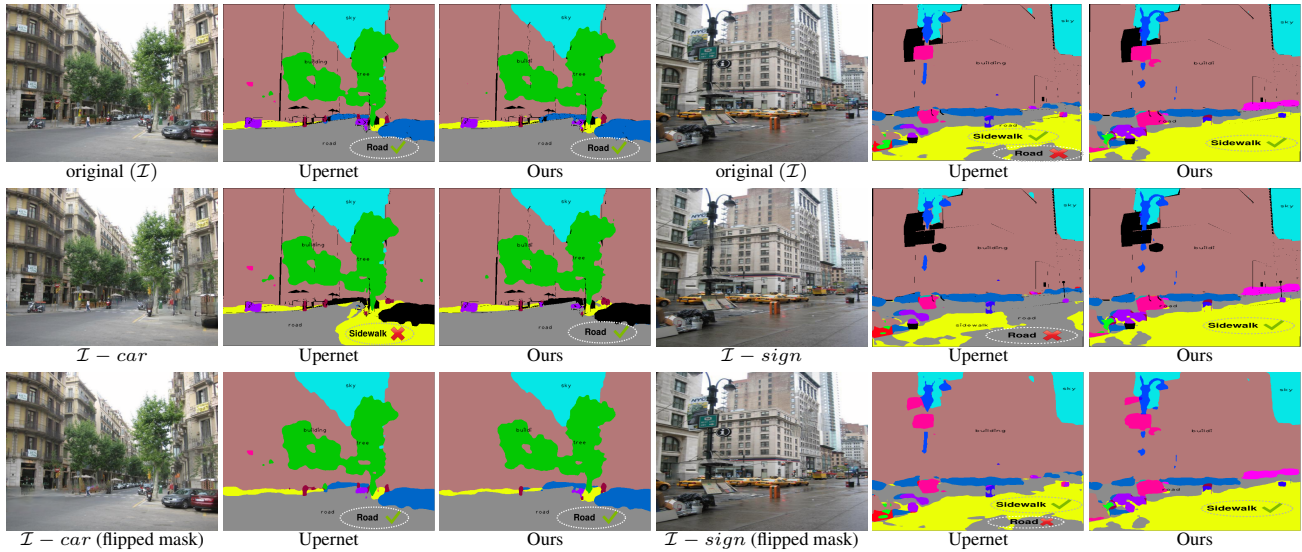


Figure 12: Experiment to verify that the volatility of the segmentation output is due to object removal and not due to editing artifacts. First row shows original images and the segmentations produced for them. Second row shows the edited images with an object removed and the segmentation output for them. Here we can see the segmentation output of Upernet significantly affected by the removal of car and sign. Final, row shows the original image edited with the same object mask as the second row, but horizontally flipped. Thus the object is not removed, but a different part of the image is edited with the same mask. We can see here that the segmentation is not affected at all by this edit and is very similar to the segmentation produced by the original image.