# Grid Saliency for
# Context Explanations of Semantic Segmentation

**Lukas Hoyer**
Bosch Center for Artificial Intelligence
lukas.hoyer@outlook.com

**Mauricio Munoz**
Bosch Center for Artificial Intelligence
AndresMauricio.MunozDelgado@bosch.com

**Prateek Katiyar**
Bosch Center for Artificial Intelligence
Prateek.Katiyar@bosch.com

**Anna Khoreva**
Bosch Center for Artificial Intelligence
Anna.Khoreva@bosch.com

**Volker Fischer**
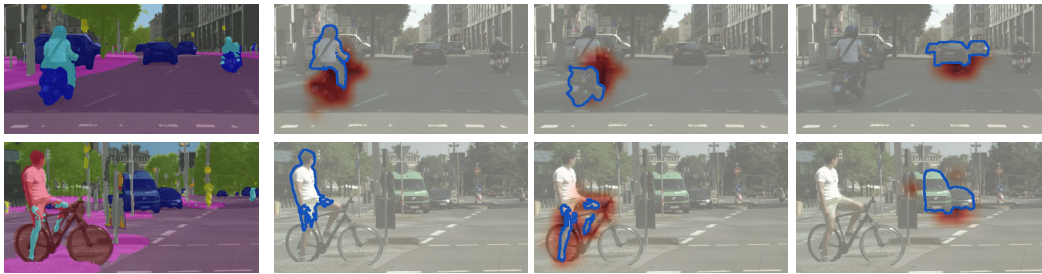Bosch Center for Artificial Intelligence
Volker.Fischer@bosch.com

## Abstract

Recently, there has been a growing interest in developing saliency methods that provide visual explanations of network predictions. Still, the usability of existing methods is limited to image classification models. To overcome this limitation, we extend the existing approaches to generate *grid saliencies*, which provide spatially coherent visual explanations for (pixel-level) dense prediction networks. As the proposed grid saliency allows to spatially disentangle the object and its context, we specifically explore its potential to produce context explanations for semantic segmentation networks, discovering which context most influences the class predictions inside a target object area. We investigate the effectiveness of grid saliency on a synthetic dataset with an artificially induced bias between objects and their context as well as on the real-world Cityscapes dataset using state-of-the-art segmentation networks. Our results show that grid saliency can be successfully used to provide easily interpretable context explanations and, moreover, can be employed for detecting and localizing contextual biases present in the data.

## 1 Introduction

In many real-world scenarios, the presence of an object, its location and appearance are highly correlated with the contextual information surrounding this object, such as the presence of other nearby objects or more global scene semantics. For example, in the case of an urban street scene, a cyclist is more likely to co-occur on a bicycle and a car to appear on the road below sky and buildings (cf. objects and their context explanations in Fig. 1). These semantic correlations are inherently present in real-world data. A data-driven model, such as a deep neural network, is prone to exploit these statistical biases during training in order to improve its prediction performance. These biases picked up by the network during training can lead to erroneous predictions and impair network generalization (cf. the effect of context on misclassification in Fig. 2). An effective and safe utilization of deep learning models for real-world applications, e.g. autonomous driving, requires a good understanding of these contextual biases inherent in the data and the extent to which a learned model incorporated them into its decision making process.

Saliency methods [4, 5, 6, 7, 8, 9] have become a popular tool to explain predictions of a trained model by highlighting parts of the input that presumably have a high relevance for its predictions.

(a) Semantic segment.          (b) Context explanations (red) for different segments (outlined in blue)

Figure 1: Context explanations by grid saliency for semantic segmentation [1, 2] on Cityscapes [3]. Grid saliency not only can contextually explain correct predictions: in the first row the network looks at the motorbike to correctly predict the class rider (light blue); but can also explain erroneous predictions: in the second row the upper body of the rider is incorrectly predicted as person, but for this prediction the bicycle is not salient in contrast to the correctly predicted legs of the rider.

However, to the best of our knowledge the existing saliency methods are mostly focused on image classification networks and thus are not able to spatially differentiate between prediction explanations. In this work, we propose a way to extend existing saliency methods designed for image classification towards (pixel-level) dense prediction tasks, which allows to generate spatially coherent explanations by exploiting spatial information in dense predictions. We call our approach *grid saliency*, which is a perturbation-based saliency method, formulated as an optimization problem of identifying the minimum unperturbed area of the image needed to retain the network predictions inside a target object region. As our grid saliency allows to differentiate between objects and their associated context areas in the saliency map, we specifically explore its potential to produce *context explanations* for semantic segmentation networks. The contextual information is known to be one of the essential recognition cues [10, 11, 12], thus we aim to investigate which local and global context is the most relevant for the network class predictions inside a target object area (see Fig. 1 and Fig. 2 for examples).

In real-world scenarios, context biases are inherently present in the data. To evaluate whether the proposed grid saliency is sensitive to context biases and has the ability to detect them, we introduce a synthetic toy dataset for semantic segmentation, generated by combining MNIST digits [13] with different fore- and background textures, for which we artificially induce a context bias between the digit and the background texture. Besides detecting the mere presence of the context bias, we also analyze the ability of our method to localize it in the image (see Sec. 4). We employ this dataset to compare our approach with different baselines, i.e., introduced extensions of gradient-based saliency methods of [4, 14, 15] to produce context explanations. We show that the proposed dataset can serve as a valid benchmark for assessing the quality of saliency methods to detect and localize context biases. We find that gradient-based techniques, in contrast to our grid saliency, are ill-suited for the context bias detection. By design they tend to produce noisy saliency maps which are not faithful to the context bias present in the data, whereas our method has higher sensitivity for context bias and thus can also precisely localize it in the image. We further evaluate grid saliency performance to produce context explanations for semantic segmentation on the real-world Cityscapes dataset [3] and experimentally show that the produced context explanations faithfully reflect spatial and semantic correlations present in the data, which were thus picked up by the segmentation network [1, 2].

To the best of our knowledge, we are the first to extend saliency towards dense-prediction tasks and use it to produce context explanations for semantic segmentation. In summary, our contributions are the following: (1) We propose an extension of saliency methods designed for classification towards dense prediction tasks. (2) We exploit the proposed grid saliency to produce context explanations for semantic segmentation and show its ability to detect and localize context biases. (3) We create a synthetic dataset to benchmark the quality of produced explanations as well as their effectiveness for context bias detection/localization. (4) We investigate the faithfulness of context explanations for semantic segmentation produced by the grid saliency on real-world data.

## 2   Related Work

**Explanations.**   Many methods attempt to interpret the network decision making process by producing explanations via bounding boxes [16, 17] or attributes [18], providing textual justifications [19, 20, 21] or generating low-level visual explanations [4, 22, 5, 6]. Our work builds on top of

2

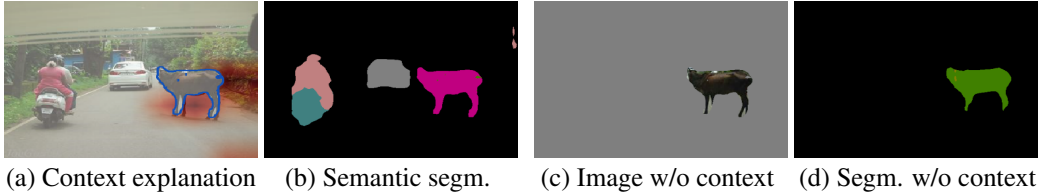| (a) Context explanation | (b) Semantic segm. | (c) Image w/o context | (d) Segm. w/o context |

Figure 2: Effect of context on semantic segmentation and the context explanation provided by grid saliency for an erroneous prediction, the image is taken from MS COCO [23]. The grid saliency (a) shows the responsible context for misclassifying the cow (green) as horse (purple) in the semantic segmentation (b). It shows the training bias that horses are more likely on road than cows. Removing the context (c) yields a correctly classified cow (d).

the latter approaches, also known as *saliency methods*, which try to identify the image pixels that contribute the most to the network prediction. These methods mostly focus on the task of image classification and can be divided into two categories: gradient-based and perturbation-based methods.

Gradient-based methods [24, 4, 25] compute a saliency map that visualizes the sensitivity of each image pixel to a specific class prediction, which is obtained by backpropagating the gradient for the prediction with respect to the image and estimating how moving along the gradient influences the class output. To circumvent noise and visual diffusion in saliency maps, [14] proposes to sum up the gradients over the intensity-scaled input versions, while [15] averages over many noisy samples of the input. Other methods [5, 26, 27] explore integrating network activations into their saliency maps. Gradient-based methods mostly rely on heuristics for backpropagation and as has been shown by [7] may provide explanations which are not faithful to the model or data.

Perturbation-based methods [6, 28, 9, 8] evaluate the class prediction change with respect to a perturbed image, e.g. in which specific regions of the image are either replaced with the mean image values or removed by applying blur or Gaussian noise. The approach in [6] formulates this output change as an optimization problem of the original and perturbed image, while [28] estimates the perturbation mask by training an auxiliary network. To account for image discontinuities, [9, 8] compute the saliency of the masked region by marginalizing it out either over neighboring image regions or by conditioning the trained generative model on the non-masked region and then estimating the classification change. Perturbation approaches might be vulnerable to network artifacts resulting in arbitrary saliency regions [8]. To overcome this, [6, 28] resort to perturbe large image regions.

The above methods are limited to explanations of image classification. In this work, we propose a way to extend them to produce grid saliency maps for dense prediction networks as well (see Sec. 3). We showcase the usability of grid saliency for context explanations of semantic segmentation (see Sec. 4 and 5). The closest related work to produce context explanations is [29], which proposes a network that jointly learns to predict and contextually explain its decisions. In contrast to [29], we focus on the post hoc model explanations and are not limited to the image classification task.

**Semantic segmentation.** CNNs have become a default technique for semantic segmentation [30, 31, 32]. [33] first showcased the use of CNNs for segmentation. Since then, multiple techniques have been proposed, from utilizing dilated convolutions [34, 35] and post-processing smoothing operations [36, 37, 38, 39] to employing spatial pyramid pooling [40, 30, 41, 2] and encoder/decoder architectures [42, 43, 44, 45]. As the accuracy of these methods comes at a high computational cost, there has been an increasing interest in developing real-time semantic segmentation networks with low memory needs [46, 1, 47, 48, 49]. In this work, we aim to produce context explanations for semantic segmentation. To investigate the effectiveness of our approach and show its generalization across architectures, we employ DeepLabv3+ [2] and U-Net [44] with different backbones [1, 50, 51].

## 3 Method

In Sec. 3.1 we introduce the grid saliency method, which allows to produce spatially coherent explanations for dense predictions, and present a way to use it for context explanations of semantic segmentation. Next, in Sec. 3.2 we extend the popular gradient-based saliency methods [4, 14, 15] to produce spatial explanations as well, which we later compare with grid saliency in Sec. 4.

## 3.1 Grid Saliency via Perturbation

Let $f : I \to O$ denote the prediction function, e.g. a deep neural network, which maps a grid input space $I = \mathbb{R}^{H_I \times W_I \times C_I}$ to a grid output space $O = \mathbb{R}^{H_O \times W_O \times C_O}$, where $W$ and $H$ are the respective width and height of the input and output, and $C_I$ and $C_O$ are the number of input channels (e.g. 3 or 1 for images) and output prediction channels (e.g. number of classes for semantic segmentation). To keep the discussion concrete, we consider only images as input, $x \in I$, and per-pixel dense predictions of the network $f(x) \in O$ as output. The goal is to find the smallest saliency map $M \in [0, 1]^{H_I \times W_I}$ that must retain in the image $x$ in order to preserve the network prediction in the request mask area $R \in \{0, 1\}^{H_O \times W_O}$ for class (channel) $c \in \{1, ..., C_O\}$. Further on, for simplicity we assume that the input and output spatial dimensions are the same.

Our method builds on top of perturbation saliency methods [6, 28, 9, 8] designed for image classification. They aim to find salient image regions most responsible for a classifier decision by replacing parts of the image with uninformative pixel values, i.e. perturbing the image, and evaluating the corresponding class prediction change. We follow the same image perturbation strategy as [6]. Let $p$ denote a perturbation function that removes information from an image $x$ outside of the saliency $M$. For example, such perturbation function can be the interpolation between $x$ and $a \in I$, where $a$ can be a constant color image, gaussian blur, or random noise. In this case, $p(x, M) = x \circ M + a \circ (1 - M)$. Note, that in practice $M$ operates on a lower resolution to avoid adversarial artifacts [6, 28] and is later upsampled to the original image resolution. In addition, the pixel values of the perturbed image $p(x, M)$ are clipped to preserve the range of the original image space.

With this notation in hand, we can formulate the problem of finding the saliency map $M$ for the prediction of class $c$ as the following optimization problem:

$$M^*(x, c) = \underset{M}{\operatorname{argmin}} \, \lambda \cdot \|M\|_1 + \| \max(f_c(x) - f_c(p(x, M)), 0)\|_1, \tag{1}$$

where $\| \cdot \|_1$ denotes the $l_1$ norm and $f_c(x)$ is the network prediction for class $c$. The first term can be considered as a mask loss that minimizes the salient image area and perturbs the original image as much as possible. The second term serves as a preservation loss which ensures that the network prediction $f_c(p(x, M))$ for class $c$ on the perturbed image $p(x, M)$ reaches at least the confidence of the network prediction $f_c(x)$ on the original unperturbed image. Thus, the second loss term can be considered as a penalty for not meeting the constraint $f_c(p(x, M)) \geqslant f_c(x)$, hence the use of $\max(\cdot, 0)$ in Eq. 1. The parameter $\lambda$ controls the sparsity of $M$.

We then can spatially disentangle explanations given in the saliency map $M$ for the network predictions in the requested area of interest $R$ from the explanations for the other predictions, by restricting the preservation loss to the request mask $R$ in Eq. (1):

$$M^*_{\text{grid}}(x, R, c) = \underset{M}{\operatorname{argmin}} \, \lambda \cdot \|M\|_1 + \frac{\|R \circ \max(f_c(x) - f_c(p(x, M)), 0)\|_1}{\|R\|_1}. \tag{2}$$

Further on, we will refer to $M^*_{\text{grid}}$ in Eq. (2) as a grid saliency map.

**Context explanations for semantic segmentation.** We now adapt the grid saliency formulation from Eq. (2) to specifically provide *context explanations* for the requested area of interest $R$. Context explanations are of particular interest for semantic segmentation, as context often serves as one of the main cues for semantic segmentation networks (see Fig. 2). Thus, here we focus on context explanations for semantic labelling predictions and assume that $R$ is the area covering the object of interest in the image $x$. To optimize for salient parts of the object context, we integrate the object request mask $R$ in the perturbation function. For the request mask $R$, the perturbed image $p(x, R) \in I$ contains only the object information inside $R$ and all the context information outside $R$ is removed (with a constant color image $a$). For optimization, we will now use this new perturbed image $p(x, R)$ instead of the maximally perturbed image $p(x, M = 0) = a$ and denote the context perturbation function as $p_{\text{context}}(x, R, M) = x \circ M + p(x, R) \circ (1 - M)$.

The context saliency map for class $c$ and request object $R$ can be obtained via optimization of

$$M^*_{\text{context}}(x, R, c) = \underset{M}{\operatorname{argmin}} \, \lambda \cdot \|M\|_1 + \frac{\|R \circ \max(f_c(x) - f_c(p_{\text{context}}(x, R, M)), 0)\|_1}{\|R\|_1}, \tag{3}$$

where the saliency map is optimized to select the minimal context necessary to at least yield the original prediction for class $c$ inside the request mask $R$. Note that, $M^*_{\text{context}}$ can be an empty mask if no context information is needed to recover the original prediction inside $R$.

## 3.2 Gradient-Based Variants

Another way to produce spatially coherent saliency maps is to make use of the popular gradient-based saliency methods. Thus, we additionally consider the Vanilla Gradient (VG) [24], Integrated Gradient (IG) [14], and SmoothGrad (SG) [15] saliency methods.

Let $G(x, c) = \partial g_c(x)/\partial x \in \mathbb{R}^{H_I \times W_I \times C_I}$ denote the gradient of the classification network prediction $g_c(x) \in \mathbb{R}$ for class $c$ with respect to the input image $x \in I$. For the classification task, the saliency maps of VG, IG and SG are computed as:

$$M^{\text{VG}}(x, c) = \sum_{c \in C_I} \left| G(x, c) \right|, \qquad M^{\text{SG}}(x, c) = \sum_{c \in C_I} \left| \frac{1}{n} \sum_{k=1}^{n} G\left(x + \mathcal{N}(0, \sigma^2), c\right) \right|,$$
$$M^{\text{IG}}(x, c) = \sum_{c \in C_I} \left| \frac{1}{n} \sum_{k=1}^{n} G\left(\frac{k}{n}x, c\right) \right|, \tag{4}$$

where $n$ is the number of approximation steps for IG or the number of samples for SG, and $\mathcal{N}(0, \sigma^2)$ represents Gaussian noise with standard deviation $\sigma$.

Following Sec. 3.1, we next extend the above approaches to produce the saliency $M$ for dense predictions $f_c(x) \in O$ and to spatially disentangle explanations given in the saliency $M$ for the network predictions in the request area $R$ from other predictions. For a given input $x$ and a binary request mask $R$, we denote the normalized network prediction score for class $c$ in the request area $R$ as $S(x, R, c) = \|R \circ f_c(x)\|_1 / \|R\|_1$, $S(x, R, c) \in \mathbb{R}$. Similarly to $G(x, c)$, we define $G_{\text{grid}}(x, R, c) \coloneqq \partial S(x, R, c)/\partial x \in \mathbb{R}^{H_I \times W_I \times C_I}$ which directly yields $M_{\text{grid}}^{\text{VG/SG/IG}}(x, R, c)$ by replacing $G(x, c)$ in Eq. 4 with $G_{\text{grid}}(x, R, c)$. For the gradient-based context saliency, as in Sec. 3.1 only salient pixels outside of the object area are considered, i.e.

$$M_{\text{context}}^{\text{VG/IG/SG}}(x, R, c) \coloneqq (1 - R) \circ M_{\text{grid}}^{\text{VG/IG/SG}}(x, R, c). \tag{5}$$

Gradient-based saliency maps are prone to be noisy. Thus, to circumvent this and also make them more comparable to the lower resolution perturbation-based grid saliency, in our experiments we apply the spatial mean filter on top of the saliency map with a $(W_I/W_S) \times (H_I/H_S)$ kernel and stride, where $W_S \times H_S$ is the resolution of the perturbation-based saliency map.

## 4 Context Bias Detection on Synthetic Data

Although context biases are inherently present in real-world data, in practice it is hard to measure and alter these correlations and thus to employ this data for benchmarking context bias detection. In order to evaluate whether the grid saliency is sensitive to context biases and has the ability to detect them, in Sec. 4.1 we introduce a synthetic dataset for semantic segmentation with artificially induced context biases.[1] We use this dataset to compare the grid saliency methods proposed in Sec. 3 and show that this dataset can serve as a valid benchmark for assessing the quality of saliency methods for context bias detection and localization, see Sec. 4.2.

### 4.1 Benchmark for Context Bias Detection

**Dataset.** The proposed synthetic toy dataset consists of gray scale images of size $64 \times 64$ pixels, generated by combining upscaled digits from MNIST [13] with foreground and background textures from [52, 53], as can be seen in Fig. 3 (a). In order to introduce different context information for the digit, two background textures are used for the upper and lower half of the image. For each synthetic image a corresponding segmentation ground truth is generated, where the MNIST digit defines the mask and the semantic class (including a background class).

To evaluate the ability of saliency methods to explain context biases, we propose to generate biased and unbiased versions of the dataset. For the unbiased version (DS$^{\text{no-bias}}$), all fore- and background textures appear with equal probability for all digits. For the biased version, a single digit class is coupled with a specific background texture, located randomly either in the upper or lower half of the background. We consider two variants of it, with a weakly DS$^{\text{w-bias}}$ and strongly induced bias

---

[1]The code for generating the proposed synthetic dataset with induced context biases can be found here: https://github.com/boschresearch/GridSaliency-ToyDatasetGen.
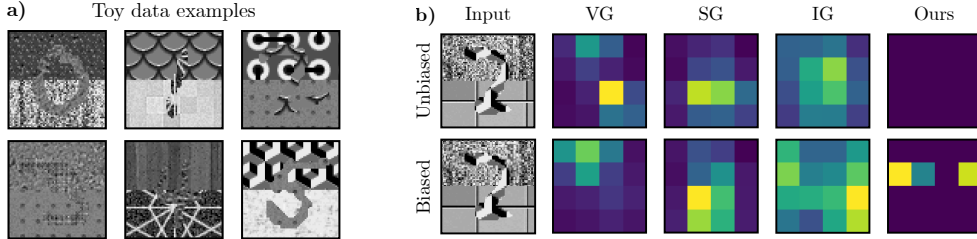
Figure 3: **a)** Synthetic dataset images, see Sec. 4.1 for details. **b)** Context saliency maps of different methods. For the image on the left, the different context saliency maps $M_{\text{context}}$ are shown for the networks [1] trained on the unbiased (top row) and biased (bottom row) dataset versions. For the biased version, the digit *2* is biased with the background texture in the top half of the image. Our perturbation-based grid saliency is able to precisely detect this bias, in contrast to other methods.

$DS^{\texttt{s-bias}}$. For the dataset with a strongly induced bias $DS^{\texttt{s-bias}}$, a specific texture appears if and only if a certain digit class is present. For a weakly induced bias $DS^{\texttt{w-bias}}$, a specific texture always appears along with the biased digit but also uniformly appears with other digits. From the pool of 25 textures, 5 textures are chosen randomly to induce context bias for one of 10 digits. For all 50 texture/digit combinations, a weakly and strongly biased dataset variant with train/test splits is generated.

**Evaluation metrics.** To evaluate the extent to which a network is able to pick up a context bias present in the training data, we propose to measure on the unbiased $DS^{\texttt{no-bias}}$ test set the performance of the network trained on the biased $DS^{\texttt{w/s-bias}}$ training set, using the standard Intersection over Union (IoU, see [33]) metric for semantic segmentation. If the network has picked up the bias, we expect to see a significant drop in IoU for the biased digit segmentation.

To benchmark how well different saliency methods (perturbation- or gradient-based) can detect the context bias of semantic segmentation networks, we propose to evaluate to which extent a context saliency map $M_{\text{context}}(x, R, c)$ (cf. Eq. 3 and 5 in Sec. 3) for the request object mask $R$ and its corresponding class $c$ is concentrated on the ground truth context area $C = 1 - R_{GT}$, where $R_{GT}$ is the ground truth mask of the object $R$, by using a context bias detection metric (CBD):

$$\text{CBD}(x, R, c) = \|C \circ M_{\text{context}}(x, R, c)\|_1 / \|C\|_1. \tag{6}$$

To benchmark the ability of different saliency methods to localize a context bias, we propose to measure how much of the context saliency $M_{\text{context}}(x, R, c)$ falls into the ground truth biased context area $C_{bias}$, the upper or lower half of the image $x$ by the design of $DS^{\texttt{w/s-bias}}$, where $C_{\text{bias}}$ is a binary mask of the biased context area. We refer to this metric as a context bias localization (CBL):

$$\text{CBL}(x, R, c) = \|C_{\text{bias}} \circ M_{\text{context}}(x, R, c)\|_1 / \|C \circ M_{\text{context}}(x, R, c)\|_1. \tag{7}$$

In our experiments we report mIoU and the mean CBD and CBL measures (mCBD, mCBL) per biased digit class $c$, averaging the results across all images in the $DS^{\texttt{w/s-bias}}$ test set, 5 randomly selected bias textures and 5 different initial random seeds for the training set generation.

## 4.2 Experimental Results

**Implementation details.** We use the U-Net [44] architecture with a VGG16 [51] backbone. As request mask, the segmentation prediction of the target digit is used. The saliency maps with a size of $4 \times 4$ are optimized using SGD with momentum of 0.5 and a learning rate of 0.2 for 100 steps starting with a 0.5 initialized mask. A weighting factor of $\lambda = 0.05$ is used (see Eq. 3). A constant color image is used for perturbation. Further implementation details are provided in the supp. material.

**Bias in the trained network.** We first investigate if the networks trained on $DS^{\text{w-bias}}$ and $DS^{\text{s-bias}}$ have picked up the induced weak and strong context biases in the training data. For this purpose, we evaluate their performance in terms of mIoU on the unbiased dataset $DS^{\text{no-bias}}$ and report the results in Fig. 4 (a), which visualizes the digit-wise mIoU with respect to the biased digit. The first row of the heat map (labeled as N) in Fig. 4 (a) shows the performance of the networks trained on $DS^{\text{no-bias}}$. We observe a clear drop in performance for biased digits (diagonal elements) in comparison to the first row. As expected, the performance drop is higher for the stronger bias. Moreover, the mIoU of the unbiased digits (non-diagonal elements) is also affected by the introduced bias. For example,
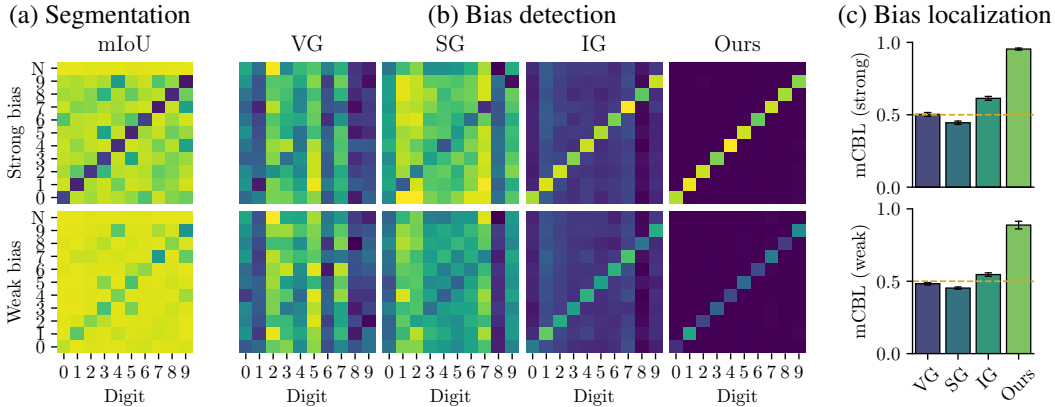
6

Figure 4: Context bias detection/localization of different saliency methods. **(a)** shows the segmentation mIoU and **(b)** & **(c)** report the context bias detection and localization results, for the strongly and weakly biased datasets, respectively. In **(a)**-**(b)**, the $y$ axis indicates the network bias towards the digit (0-9), with $N$ denoting the unbiased setting. The $x$ axis indicates the corresponding digit result. In contrast to VG [24], IG [14], and SG [15], our grid saliency is able to accurately detect (see diagonal elements in **(b)**) and localize **(c)** the induced context bias, see Sec. 4.2 and suppl. material for details.

inducing a bias for the digit nine leads to a decreased performance for the digit four (see second row in Fig. 4 (a)). We observe that this effect mostly occurs for the similar looking digits and, most likely, is caused by the fact that on the unbiased dataset the bias textures also occur with the unbiased digits, resulting in the confusion of similar looking digits for the network. From the observed mIoU drop for biased digits we can conclude that the networks have picked up the introduced bias. However, in real world it is often impossible to collect fully unbiased data. Thus, we next evaluate the ability of our grid saliency to detect context bias only using the biased data.

**Context bias detection with grid saliency.** In the last column of Fig. 4 (b), we report the context bias detection results for our perturbation-based grid saliency method, described in Sec. 3.1, using the CBD metric (see Eq. 6). The mCBD values are visualized with respect to networks trained on data biased to different digits $DS^{s/w\text{-bias}}$ (y-axis) and for the different digit classes (x-axis) in the biased test set. The only exception is the first row (labeled as N), where for comparison we show the results with no bias, for the network trained on $DS^{no\text{-bias}}$. We observe that our grid saliency shows substantial evidence of context bias for digits with induced bias (diagonal elements), both strong and weak. Even for the weak bias in Fig. 4 (b) the grid saliency still clearly differentiates between biased and unbiased digits (diagonal vs. non-diagonal elements). Note that the bias detection using the grid saliency does not require an unbiased test set. In the suppl. material, we also study the influence of the bias texture as well as the choice of hyperparameters.

**Comparison across different saliency methods.** In Fig. 4 (b) and (c) we compare our perturbation-based grid saliency with the context saliency extensions of gradient-based methods, i.e. VG, SG, and IG (see Eq. 4), using mCBD and mCBL metrics proposed in Sec. 4.1. From Fig. 4 (b) we see that VG and SG are not able to reliably detect the context bias, while IG achieves a comparable performance to our grid saliency. However, in contrast to the perturbation saliency, IG has also high mCBD values for unbiased digits, complicating its use in practice, as one would need to tune a reliable detection threshold, which is particularly challenging for weakly biased data.

**Context bias localization.** In addition to the detection performance, we evaluate how well the saliency methods localize the context bias using the mCBL measure. Fig. 4 (c) shows that VG, IG, and SG have low localization performance, comparable to random guessing ($\sim 0.5$ mCBL), while our grid saliency is the only method which is able to accurately localize the context bias (mCBL above $0.9$) on both strongly and weakly biased data. We evaluate the ability of our perturbation grid saliency to detect and localize context bias across different semantic segmentation networks. For U-Net [44] with the VGG16 [51], ResNet18 [50] and MobileNetv2 [1] backbones, we observe similar bias detection and localization performance (see also supp. material).
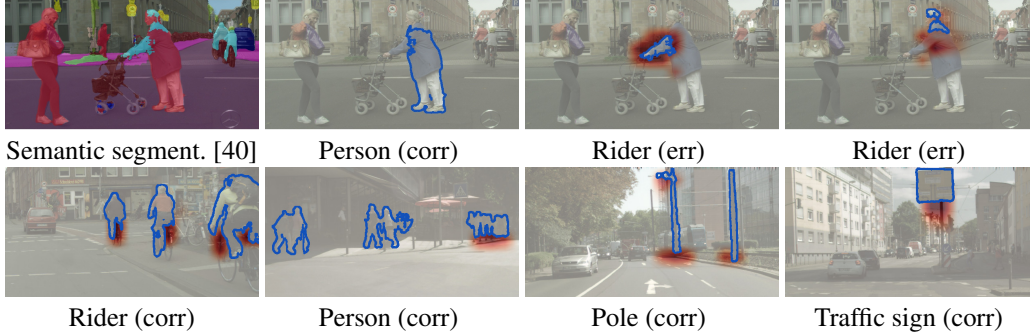
Figure 5: Context explanations by grid saliency for erroneous (err) and correct (corr) semantic segmentation predictions on Cityscapes [3]. Salient red regions visualize the most important context $M_{\text{context}}$ for the requested segment prediction $R$ outlined in blue. See Sec. 5 for discussion.

## 5 Cityscapes Experiments

Motivated by the success of our perturbation-based grid saliency at detecting context biases on the synthetic dataset described in Sec. 4, we next apply it to Cityscapes [3] in order to produce and analyze context explanations for semantic segmentation in real-world scenes.

**Experimental setup.** We use 500 finely annotated images of the Cityscapes validation set, considering only a subset of classes for the analysis. For our experiments we use the Deeplabv3+ [2] network with a Mobilenetv2 backbone [1]. Our optimization setup largely carries over from Sec. 4.2, with the exception that we optimize a coarse 16 by 32 pixel mask using SGD with a learning rate of 1 for 80 steps and use $\lambda = 0.01$. Additional implementation details are provided in the supp. material.

**Analysis of context explanations.** In order to gain a global understanding of the model context bias, we aggregate statistics from the produced grid saliency maps for several requested classes across all validation images. These statistics are summarized in Fig. 6. For each image and class, the context saliency is computed for all sufficiently large instances (at least 10k pixels), as can be seen in the second row of Fig. 5. Next, for each image and semantic class we compute the weighted label distribution of context salient pixels. The weight of each pixel is given by the computed saliency intensity at said pixel and its context class label is taken from the ground truth segmentation. For each requested object class, we group by label across all images and sum the saliency values within each group. This yields an accumulated saliency value for each label. We normalize the results by the sum of all saliency values across all classes and images to yield a probability distribution across labels.

To show that context explanations produced by our grid saliency are meaningful, we additionally compare its accumulated context class statistics with their baseline class distributions. These baselines are computed in the same manner as above, but instead of relying on the optimized saliency map we use a fixed dilation of the contour of all predicted object segments considered in the image, thereby capturing the immediate context around each object in a uniform set of directions.

Fig. 6 compares these accumulated statistics of context explanations for the requested object classes, given on the $y$ axis. We observe that the context explanations are focused across reasonable and somewhat expected class subsets, which vary per class. For instance, in comparison to the baseline, context saliency for the rider class shifts attention from its spatially co-occurring classes such as road, vegetation and building, to mainly the bicycle class. Car context saliency attention is decreased on building and vegetation and mainly focused on road, sidewalk and other cars. Bicycle context saliency mostly attends to sidewalk rather than the road class. Overall, our grid saliency is able to provide sensible and coherent explanations for network decision making, which reflect semantic dependencies present in street scenes. These quantitative findings are validated by the qualitative results in Fig. 5. The second row of Fig. 5 gives several representative examples of context saliencies. Please refer to the supp. material for additional examples and a high-resolution version of Fig. 6 as well as a qualitative model comparison on MS COCO.

**Context explanations of erroneous predictions.** One of the motivations for this work is also to explain unexpected model behavior. Specifically, in case of erroneous predictions (i.e., Fig. 1 and 2), we wish to understand to which extent context bias contributes to the failure. Table 1 shows class statistics of context explanations for correct and erroneous predictions, where the intersection of the
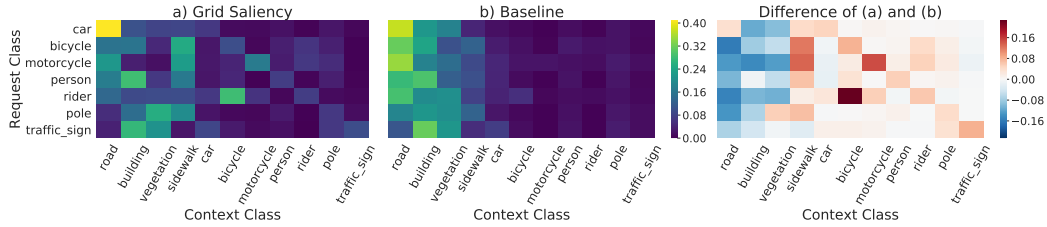
Figure 6: Class statistics of context explanations on Cityscapes, see Sec. 5 for discussion.

Table 1: Class statistics of context explanations of correct and erroneous predictions on Cityscapes.

| GT | Prediction | Context class | | | | | | | | | |
| | | road | bicycle | motorcycle | car | vegetation | building | sidewalk | rider | person | pole |
|---|---|---|---|---|---|---|---|---|---|---|---|
| rider | rider | 0.17 | **0.22** | 0.03 | 0.07 | 0.09 | 0.11 | 0.10 | 0.07 | 0.02 | 0.02 |
| | person | 0.13 | 0.10 | 0.01 | 0.05 | 0.10 | **0.24** | 0.05 | 0.08 | 0.11 | 0.04 |
| person | person | 0.12 | 0.03 | 0.00 | 0.09 | 0.15 | **0.22** | 0.13 | 0.01 | 0.05 | 0.05 |
| | rider | 0.09 | **0.18** | 0.02 | 0.04 | 0.10 | 0.13 | 0.07 | 0.08 | 0.15 | 0.03 |
| | car | 0.19 | 0.01 | 0.00 | **0.32** | 0.05 | 0.11 | 0.03 | 0.01 | 0.19 | 0.01 |
| bicycle | bicycle | 0.18 | 0.06 | 0.01 | 0.05 | 0.09 | 0.15 | **0.19** | 0.05 | 0.04 | 0.04 |
| | motorcycle | 0.19 | **0.27** | 0.13 | 0.01 | 0.02 | 0.11 | 0.13 | 0.03 | 0.01 | 0.04 |
| car | car | **0.39** | 0.02 | 0.00 | 0.06 | 0.11 | 0.10 | 0.08 | 0.01 | 0.04 | 0.04 |
| | motorcycle | 0.07 | 0.06 | **0.22** | 0.17 | 0.14 | 0.07 | 0.02 | 0.01 | 0.00 | 0.04 |
| | person | 0.14 | 0.01 | 0.00 | **0.28** | 0.03 | 0.14 | 0.07 | 0.01 | 0.22 | 0.02 |
| motorcycle | motorcycle | **0.18** | 0.04 | 0.10 | 0.04 | 0.07 | 0.14 | **0.18** | 0.05 | 0.03 | 0.06 |
| | bicycle | 0.13 | 0.10 | 0.07 | 0.10 | 0.10 | **0.14** | 0.11 | 0.01 | 0.05 | 0.03 |
| | car | 0.16 | 0.01 | **0.18** | 0.10 | 0.10 | 0.10 | 0.11 | 0.05 | 0.04 | 0.06 |
| pole | pole | 0.04 | 0.02 | 0.00 | 0.05 | **0.28** | 0.23 | 0.10 | 0.00 | 0.03 | 0.04 |
| | building | 0.04 | 0.01 | 0.00 | 0.04 | 0.07 | **0.57** | 0.02 | 0.00 | 0.01 | 0.05 |
| | vegetation | 0.03 | 0.00 | 0.00 | 0.03 | **0.63** | 0.08 | 0.03 | 0.00 | 0.00 | 0.07 |

ground truth mask and the correct or misclassified prediction region are used as the request mask $R$, and similarly to Fig. 6 semantic class statistics are computed inside the salient context for all sufficiently large request masks (at least 625 pixels). We observe that the saliency severely changes for the error cases in comparison to the correct predictions. E.g. for the correctly classified rider class, context saliency mostly focuses on bicycle (22%), but is much less present (10%) when rider is mistaken as person and in this case the context saliency shifts from bicycle to building (24%). The opposite effect is observed when person is misclassified as rider. Sidewalk saliency is decreased when bicycle is misclassified as motorcycle (19% vs. 13%) and when pole is misclassified as building (10% vs. 2%). In general, grid saliency is able to provide reasonable explanations of erroneous predictions.

Fig. 1 (second row) and Fig. 5 (first row) showcase how applying grid saliency on the same input image but different prediction outputs is useful for isolating failures caused by context bias. Both of them show examples of context explanations for erroneous segmentations of a single object. In Fig. 5, the arms, head and shoulders of a pedestrian are classified as rider, while the rest of the body is correctly identified as person. The context saliency for the former body parts activates highly on the arms and stroller handles, whereas activations for person does not highlight this support. Thus, a reasonable conclusion is that the misclassification may be attributed to the arm pose and potentially also to the context bias given by the similar appearance of the stroller and bicycle handles.

## 6 Conclusions

We proposed spatial grid saliency, a general framework to produce spatially coherent explanations for (pixel-level) dense prediction networks, which to the best of our knowledge is the first method to extend saliency techniques beyond classification models. We investigated the ability of grid saliency to provide context explanations for semantic segmentation, showing its effectiveness to detect and localize context bias on the synthetic toy dataset specifically designed with an artificially induced bias to benchmark this task. Our results on the real-world data indicated that grid saliency can be successfully employed to produce easily interpretable and faithful context explanations for semantic segmentation, helping to discover spatial and semantic correlations in the data picked up by the network. We hope the proposed grid saliency and the insights of this work can contribute to a better understanding of semantic segmentation networks or other models for dense prediction, elucidating some aspects of the problem that have not been well explored so far.

Besides enabling visual explanations for dense-prediction tasks, we see potential utility of grid saliency for the following applications: *1) Architecture comparison:* Context explanations produced by grid saliency can be used to compare architectures with respect to their capacity to either learn or to be invariant towards context. *2) Network generalization via active learning:* Existing context biases might impair network generalization. E.g., cows might mostly appear on grass during training. A network that was trained and evaluated on this data and picked up that bias will perform poorly in real-world cases, where the cow, for example, appears on road and gets misclassified as the horse class, see Fig. 2. Actions can be taken, such as targeted extra data collection, to improve network generalization. *3) Adversarial detection:* Grid saliency can be used to detect and localize adversarial patches outside object boundaries (e.g. [54]). Cases for which the salient regions lie largely outside an object would strongly indicate the presence of an adversary or misguided prediction. We consider the above-mentioned utilities of grid saliency interesting and promising future research directions.

# References

[1] Mark B. Sandler, Andrew G. Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4510–4520, 2018.

[2] Liang-Chieh Chen, Yukun Zhu, George Papandreou, Florian Schroff, and Hartwig Adam. Encoder-decoder with atrous separable convolution for semantic image segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[3] Marius Cordts, Mohamed Omran, Sebastian Ramos, Timo Rehfeld, Markus Enzweiler, Rodrigo Benenson, Uwe Franke, Stefan Roth, and Bernt Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[4] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *International Conference on Learning Representations (ICLR)*, 2013.

[5] Ramprasaath R. Selvaraju, Abhishek Das, Ramakrishna Vedantam, Michael Cogswell, Devi Parikh, and Dhruv Batra. Grad-cam: Why did you say that? `arXiv:1611.07450`, 2016.

[6] Ruth Fong and Andrea Vedaldi. Interpretable explanations of black boxes by meaningful perturbation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, pages 3449–3457, 2017.

[7] Julius Adebayo, Justin Gilmer, Michael Muelly, Ian J. Goodfellow, Moritz Hardt, and Been Kim. Sanity checks for saliency maps. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2018.

[8] Chun-Hao Chang, Elliot Creager, Anna Goldenberg, and David Duvenaud. Explaining image classifiers by counterfactual generation. In *International Conference on Learning Representations (ICLR)*, 2019.

[9] Luisa M. Zintgraf, Taco Cohen, Tameem Adel, and Max Welling. Visualizing deep neural network decisions: Prediction difference analysis. In *International Conference on Learning Representations (ICLR)*, 2017.

[10] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Loddon Yuille. The role of context for object detection and semantic segmentation in the wild. 2014.

[11] Jasper RR Uijlings, Arnold WM Smeulders, and Remko JH Scha. The visual extent of an object. *International journal of computer vision*, 96(1):46–63, 2012.

[12] Aymen Azaza, Joost van de Weijer, Ali Douik, and Marc Masana. Context proposals for saliency detection. *Computer Vision and Image Understanding*, 174:1–11, 2018.

[13] Yann LeCun, Corinna Cortes, and Christopher J.C. Burges. The MNIST database of handwritten digits. Technical report, 1998.

[14] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. Axiomatic attribution for deep networks. In *International Conference on Machine Learning (ICML)*, 2017.

[15] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda B. Viégas, and Martin Wattenberg. Smoothgrad: removing noise by adding noise. `arXiv:1706.03825`, 2017.

[16] Andrej Karpathy and Li Fei-Fei. Deep visual-semantic alignments for generating image descriptions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[17] M. Oramas JoséOramas and Tinne Tuytelaars. Modeling visual compatibility through hierarchical mid-level elements. `arXiv:1604.00036`, 2016.

[18] Sadaf Gulshad, Jan Hendrik Metzen, Arnold W. M. Smeulders, and Zeynep Akata. Interpreting adversarial examples with attributes. `arXiv:1904.08279`, 2019.

[19] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.

[20] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. Generating visual explanations. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.

[21] Yinpeng Dong, Hang Su, Jun Zhu, and Bo Zhang. Improving interpretability of deep neural networks with semantic information. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

[22] Bolei Zhou, Aditya Khosla, Àgata Lapedriza, Aude Oliva, and Antonio Torralba. Learning deep features for discriminative localization. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[23] Tsung-Yi Lin, Michael Maire, Serge J. Belongie, Lubomir D. Bourdev, Ross B. Girshick, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C. Lawrence Zitnick. Microsoft coco: Common objects in context. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014.

[24] Matthew D. Zeiler and Rob Fergus. Visualizing and understanding convolutional networks. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2014.

[25] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin A. Riedmiller. Striving for simplicity: The all convolutional net. In *International Conference on Learning Representations (ICLR)*, 2014.

[26] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, Wojciech Samek, and Oscar Deniz Suarez. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. In *PloS one*, 2015.

[27] Jianming Zhang, Sarah Adel Bargal, Zhe L. Lin, Jonathan Brandt, Xiaohui Shen, and Stan Sclaroff. Top-down neural attention by excitation backprop. *International Journal of Computer Vision (IJCV)*, 2017.

[28] Piotr Dabkowski and Yarin Gal. Real time image saliency for black box classifiers. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2017.

[29] Maruan Al-Shedivat, Kumar Avinava Dubey, and Eric P. Xing. Contextual explanation networks. `arXiv:1705.10301`, 2018.

[30] Hengshuang Zhao, Jianping Shi, Xiaojuan Qi, Xiaogang Wang, and Jiaya Jia. Pyramid scene parsing network. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 6230–6239, 2017.

[31] Tobias Pohlen, Alexander Hermans, Markus Mathias, and Bastian Leibe. Full-resolution residual networks for semantic segmentation in street scenes. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3309–3318, 2017.

[32] Guosheng Lin, Anton Milan, Chunhua Shen, and Ian D. Reid. Refinenet: Multi-path refinement networks for high-resolution semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5168–5177, 2017.

[33] Evan Shelhamer, Jonathan Long, and Trevor Darrell. Fully convolutional networks for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.

[34] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan Loddon Yuille. Semantic image segmentation with deep convolutional nets and fully connected crfs. In *International Conference on Learning Representations (ICLR)*, 2015.

[35] Fisher Yu and Vladlen Koltun. Multi-scale context aggregation by dilated convolutions. In *International Conference on Learning Representations (ICLR)*, 2016.

[36] Guosheng Lin, Chunhua Shen, Ian D. Reid, and Anton van den Hengel. Efficient piecewise training of deep structured models for semantic segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[37] Shuai Zheng, Sadeep Jayasumana, Bernardino Romera-Paredes, Vibhav Vineet, Zhizhong Su, Dalong Du, Chang Huang, and Philip H. S. Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.

[38] Siddhartha Chandra and Iasonas Kokkinos. Fast, exact and multi-scale inference for semantic image segmentation with deep gaussian crfs. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.

[39] Varun Jampani, Martin Kiefel, and Peter V. Gehler. Learning sparse high dimensional filters: Image filtering, dense crfs and bilateral neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.

[40] Liang-Chieh Chen, George Papandreou, Iasonas Kokkinos, Kevin Murphy, and Alan Loddon Yuille. Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 40:834–848, 2018.

[41] Golnaz Ghiasi and Charless C. Fowlkes. Laplacian pyramid reconstruction and refinement for semantic segmentation. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.

[42] Seunghoon Hong, Hyeonwoo Noh, and Bohyung Han. Decoupled deep neural network for semi-supervised semantic segmentation. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2015.

[43] Hyeonwoo Noh, Seunghoon Hong, and Bohyung Han. Learning deconvolution network for semantic segmentation. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*, 2015.

[44] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.

[45] Vijay Badrinarayanan, Alex Kendall, and Roberto Cipolla. Segnet: A deep convolutional encoder-decoder architecture for image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence (PAMI)*, 39:2481–2495, 2016.

[46] Rudra P. K. Poudel, Stephan Liwicki, and Roberto Cipolla. Fast-scnn: Fast semantic segmentation network. `arXiv:1902.04502`, 2019.

[47] Adam Paszke, Abhishek Chaurasia, Sangpil Kim, and Eugenio Culurciello. Enet: A deep neural network architecture for real-time semantic segmentation. `arXiv:1606.02147`, 2017.

[48] Davide Mazzini. Guided upsampling network for real-time semantic segmentation. 2018.

[49] Hengshuang Zhao, Xiaojuan Qi, Xiaoyong Shen, Jianping Shi, and Jiaya Jia. Icnet for real-time semantic segmentation on high-resolution images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.

[50] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778, 2016.

[51] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *International Conference on Learning Representations (ICLR)*, 2015.

[52] SiteOrigin. Background image generator. `http://bg.siteorigin.com/`.

[53] Subtle patterns. `https://www.toptal.com/designers/subtlepatterns/`.

[54] Mark Lee and Zico Kolter. On physical adversarial patches for object detection. *arXiv preprint arXiv:1906.11897*, 2019.

[55] Tijmen Tieleman and Geoffrey Hinton. Lecture 6.5-rmsprop: Divide the gradient by a running average of its recent magnitude. *COURSERA: Neural Networks for Machine Learning 4*, 2012.

[56] François Chollet. Xception: Deep learning with depthwise separable convolutions. *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2017.

# 7 Further Details on Context Bias Detection on Synthetic Data

## 7.1 Synthetic Dataset

The proposed synthetic toy dataset consists of grayscale images of size $64 \times 64$ pixels, generated by combining digits from MNIST [13] with foreground and background textures from [52, 53], as can be seen in Fig. 7. Each MNIST image is bilinearly upsampled to $64 \times 64$ and used as a binary mask for the digit shape (with threshold $= 0.5$), which interior and exterior regions are filled with fore- and background textures, respectively. In total, a pool of 25 textures are considered for fore- and background generation. To induce the texture variance among instances and make the segmentation task more challenging, we crop the textures randomly from the texture images of at least size $400 \times 400$ generated via [52]. One crop for each texture is visualized in Fig. 8. The texture variation between different random crops is shown in Fig. 9. For each synthetic image a corresponding segmentation ground truth is generated, where the MNIST digit defines the mask and the semantic class. Overall, 11 semantic classes are considered, 10 digits plus the background class.

For all dataset variants, a training and a test set are generated from the original MNIST training and validation sets respectively, using all 25 textures. The training set contains 50k images and the test set consist of 1k images. For the training set, the seed for setting up the random generator can be varied in order to obtain different dataset variants. It specifies the digit order as well as the texture selection. Code to generate the toy dataset is provided under `https://github.com/boschresearch/GridSaliency-ToyDatasetGen`.
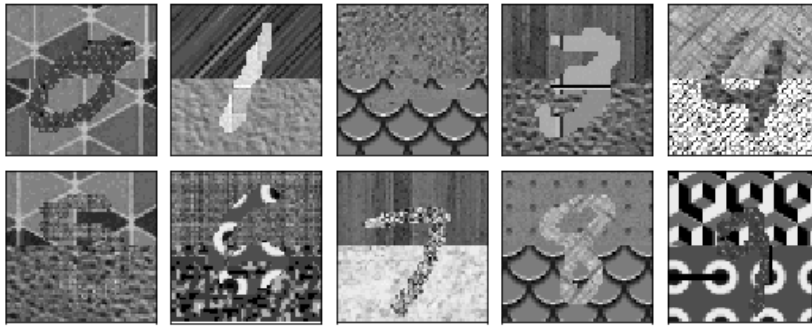


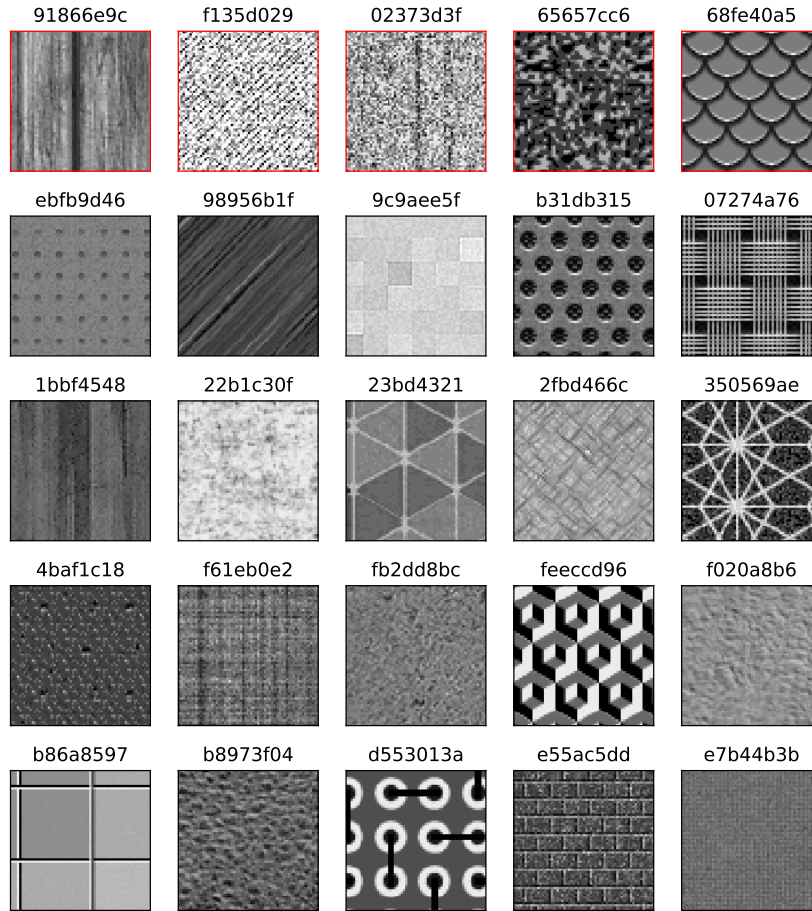Figure 7: Examples from the toy dataset.

Figure 8: Random crop of each texture in the texture pool. The textures used as bias texture are shown in the first row (red framed).
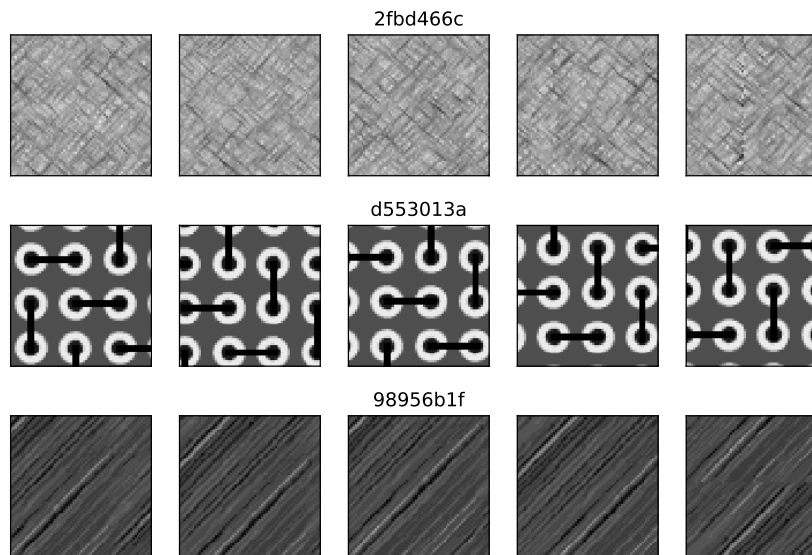


Figure 9: Variance of textures across different random crops.

## 7.2 Experimental Setup

We utilize the U-Net [44] architecture with VGG16 [51], Mobilenetv2 [1], and ResNet18 [50] backbones, with a pixel-wise softmax layer (11 classes) as the final layer. For training, we use the multi-class cross entropy loss with the RMSPROP [55] optimizer. The networks are trained for 50 epochs, with the batch size 100; the learning rate is set to $10^{-3}$ and $10^{-4}$ for the first 40 and last 10 epochs, respectively.

To avoid adversarial artifacts in the estimation of the perturbation-based context saliency, we use the saliency maps of a lower resolution for optimization, i.e. $4 \times 4$ for $M_{context}$, and then upsample them to the original size of $64 \times 64$ using nearest neighbor upsampling for evaluation. We use a set of different constant perturbation grayscale values $\{0, 0.25, 0.5, 0.75, 1.0\}$ and for each image choose the one with the lowest loss value for $M = 0$, to avoid hurting the segmentation by having a low contrast between the object and the perturbation image. In order to avoid optimizing for the border artifacts of the network predictions, for the preservation loss computation, the request mask is eroded with a $3 \times 3$ kernel size.

The saliency maps are optimized using SGD with a momentum of $0.5$ and a learning rate of $0.2$ for 100 steps, starting with the initial map $M = 0.5$ and clipping saliency map values below $0.2$ to $0$ at each step in order to avoid artifacts. A weighting factor of $\lambda = 0.05$ is used (see Eq. 3). After the optimization, we check if an empty mask achieves a lower loss value than the optimized one and use the mask with the lowest loss as the final saliency map.

For the IG saliency method, the number of interpolation steps $n$ is set to 25 and for the SG saliency, we use 25 Gaussian noise samples with $\mu = 0, \sigma = 0.15(\max(x) - \min(x))$.

All computations were done on a Nvidia Titan Xp GPU with 12 GB memory.

## 7.3 Bias Detection

In this and the following section, we provide a more detailed analysis of the different saliency methods benchmarked on the toy dataset described in the main paper in Sec. 4. Along with our approach we additionally consider the Vanilla Gradient (VG) [24], Integrated Gradient (IG) [14], and SmoothGrad (SG) [15] saliency methods. Instead of aggregating the IoU, context bias detection (CBD) and context bias localization (CBL) metrics over multiple bias textures, we show them separately, to gain some insight how the bias texture may affect the performance of the context saliencies. The results are averaged over 5 different random seeds used for the training data generation.

We evaluate different networks trained on data biased towards a specific digit with a specific bias texture. When testing on the unbiased dataset, there is a drop in the segmentation IoU for the biased digit, which shows that the network has picked up the bias. As can be seen in Fig. 10 a) and 11 a) the extent of the drop is stable across different bias textures and is mostly affected by the bias digit.

For the different saliency methods, we have checked if this bias can be detected only using the biased dataset. For VG and SG (see Fig. 10 b/c) and 11 b/c)), the CBD highly deviates between the biased digits (diagonal), with the amount and direction heavily dependent on the bias texture. In practice, these methods are not applicable as there is no control over the bias texture. By design these gradient-based methods are more sensitive to high frequency patterns and thus lead to unfaithful explanations [7]. For IG and our context saliency (see Fig. 10 d/e) and 11 d/e)), we can observe a significantly smaller dependency on the bias texture allowing a bias detection to be independent from the bias texture. Moreover, for the weakly biased dataset, our grid saliency also roughly reflects the extent of the bias, which we observed in the segmentation IoU drop on the unbiased dataset (compare Fig. 11 a) and Fig. 11 e)).

Note that the mIoU drop for unbiased digits which look similar to the biased ones (e.g., four and nine) is not reflected in the generated saliency maps. By design, our method only looks for positive evidence (context that is present in the image and supports the classification) without taking into account negative evidence for the prediction. By applying an biased network to an unbiased dataset it is exposed to those negative biases as the bias texture acts as negative evidence for unbiased digits.
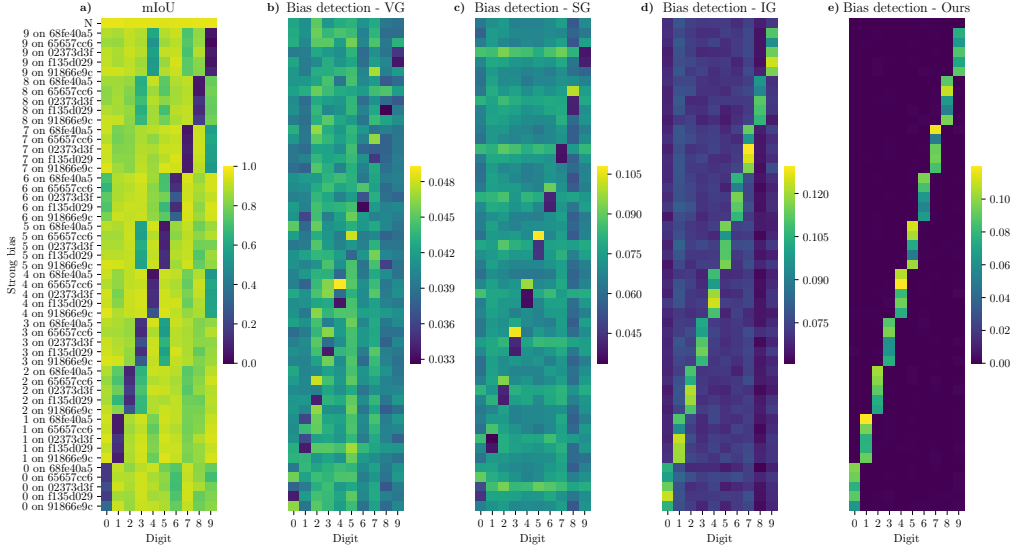
Figure 10: Context bias detection comparison of different saliency methods for strongly biased datasets. Instead of averaging over five different bias textures as done in Fig. 3, each bias texture is listed separately to show how the bias texture affects the different saliency methods. The bias detection performance is measured using the CBD metric (see Eq. 6).
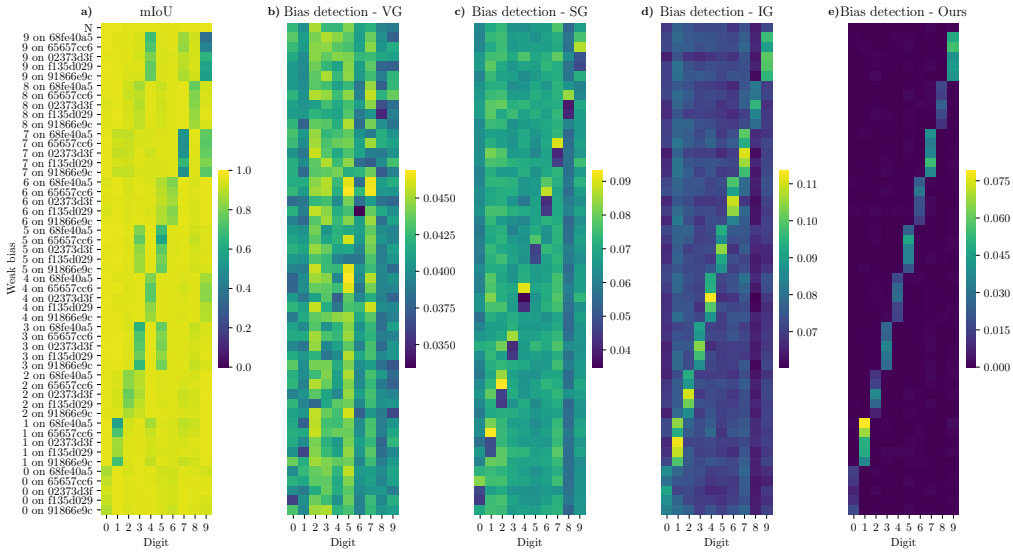


Figure 11: Context bias detection comparison of different saliency methods for weakly biased datasets. Instead of averaging over five different bias textures as done in Fig. 3, each bias texture is listed separately to show how the bias texture affects the different saliency methods. The bias detection performance is measured using the CBD metric (see Eq. 6).

## 7.4 Bias Localization

A similar effect can be seen for bias localization in Fig. 12 and 13. While VG and SG (column a) and b)) highly depend on the bias texture causing the localization even to focus more on the unbiased half, IG and our method (column c) and d)) depend significantly less on the bias texture. However, IG only achieves a bias localization slightly above random guessing while our method is able to localize both strong and weak biases very well.
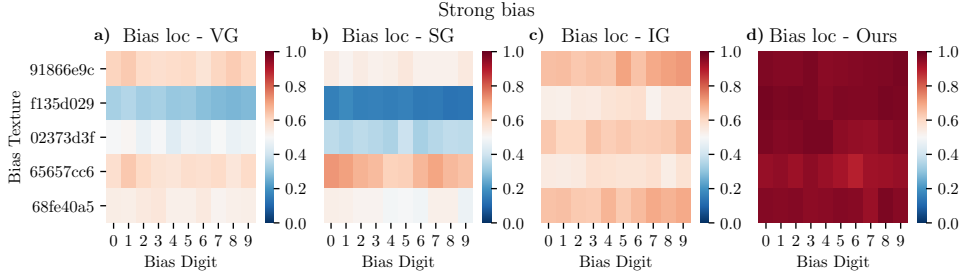
Figure 12: Bias localization across different grid saliency methods for strongly biased datasets. The CBL metric (see Eq. 7), averaged over five training dataset generation seeds, is shown with respect to bias texture and bias digit.
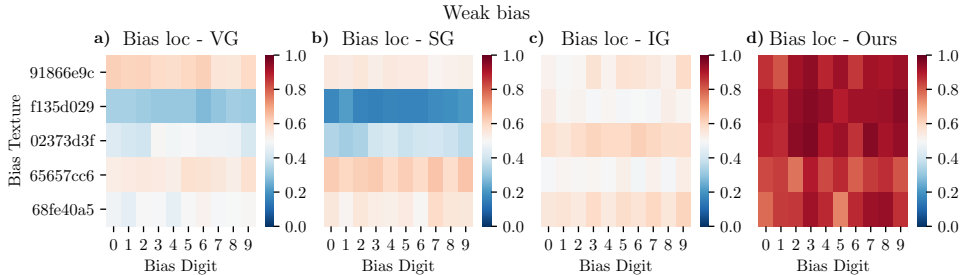


Figure 13: Bias localization across different grid saliency methods for weakly biased datasets. The CBL metric (see Eq. 7), averaged over five training dataset generation seeds, is shown with respect to bias texture and bias digit.

## 7.5 Network Comparison

In order to show that our grid saliency can be applied to several network architectures, we have repeated the experiments from Sec. 4.2 of the main paper with different backbones for Unet [44]. For that purpose, we have chosen VGG16 [51], ResNet18 (RN) [50], and MobileNetv2 (MN) [1] due to their different structure. All values are aggregated over five different bias textures and five training set generation seeds showing both the mean and the standard deviation. For the architecture comparison, only the weakly biased dataset variants are used.

**Segmentation mIoU on the unbiased test set.** First, we have checked if the different architectures have picked up the bias from the dataset by applying network instances trained on biased datasets to the unbiased test set and calculated the segmentation mIoU. As can be see in Fig. 15 a), all architectures have a clear drop in mIoU for the biased digits (diagonal) in comparison to the baseline mIoU of the unbiased network (top row N). Therefore, we can conclude that all networks have picked up the bias. VGG achieves the best base segmentation performance and has also a smaller drop in mIoU for biased digits, meaning that it is less susceptible to a context bias. For that reason, we have chosen this configuration for the main paper.

**Bias Detection and Localization.** Next, we have applied the context saliency methods to the different versions of the biased datasets to check if they are able to show the presence of a bias in the biased dataset itself.

For VG and SG (see Fig. 14), the mean over different bias textures is not susceptible to the bias, however, there is a high standard deviation for the biased digits (diagonal) caused by the dependence on the bias texture (see Sec. 7.3). IG and our method are both able to detect the context bias independent from the chosen network architecture. For VG, SG, and IG, we also observe different levels of focus on the context for different backbones, independent of the digit and if it is biased or not. For the bias localization, our method stably outperforms VG, SG, and IG across different backbones (see Fig. 16). For grid saliency we also observe a slight localization improvement for VGG over MN and RN, in contrast to other methods.
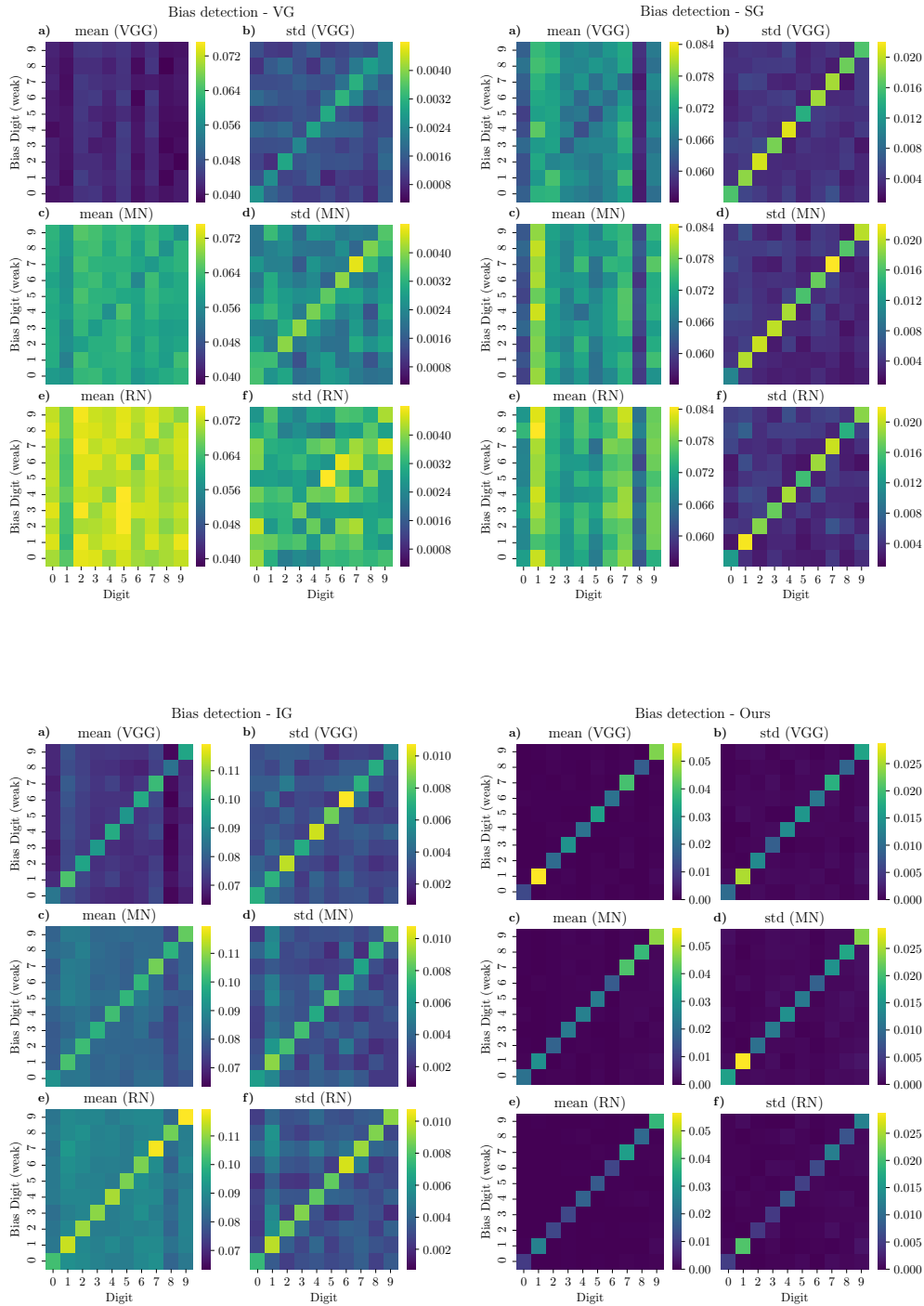
Figure 14: Comparison of the different grid saliency methods VG (top-left), SG (top-right), IG (bottom-left), and our perturbation saliency (bottom-right) across different network architectures (from top to bottom: VGG16 (VGG), Mobilenetv2 (MN), and ResNet18 (RN)). For each combination of grid saliency method and architecture, the mean and standard deviation (std) plots for context bias detection (CBD metric) are shown. All values are aggregated over five different bias textures and training set generation seeds.
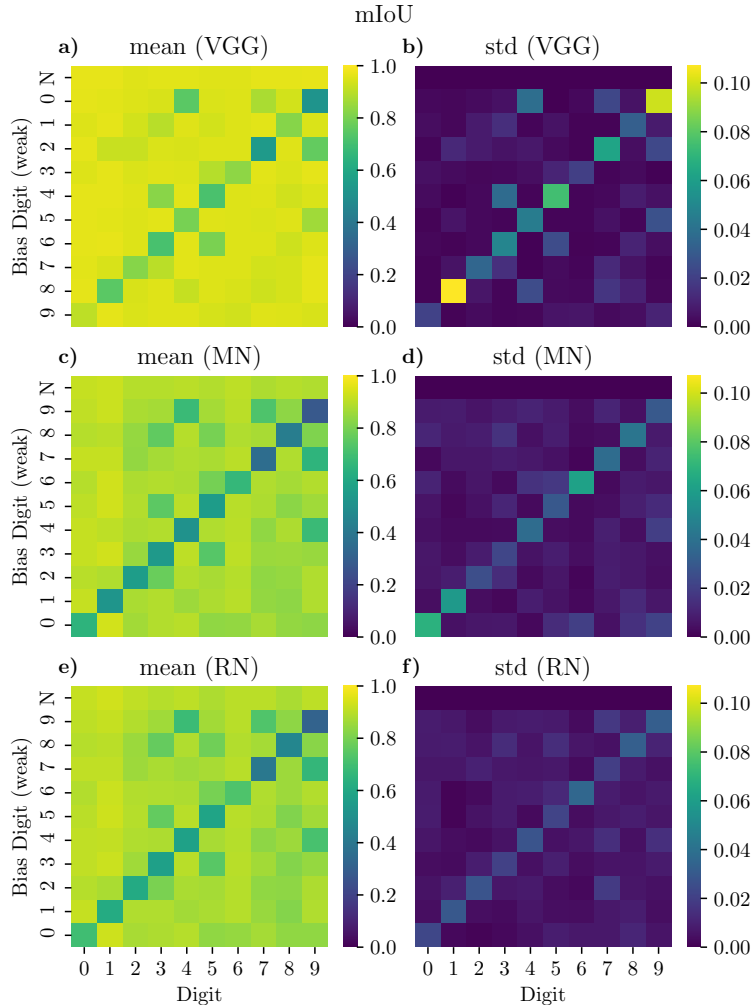
18

Figure 15: Comparison of the segmentation mIoU on the unbiased test set for different network architectures: VGG16 (VGG), Mobilenetv2 (MN), and ResNet18 (RN). All values are aggregated over five different bias textures and training set generation seeds. Both the mean and the standard deviation (std) are visualized.
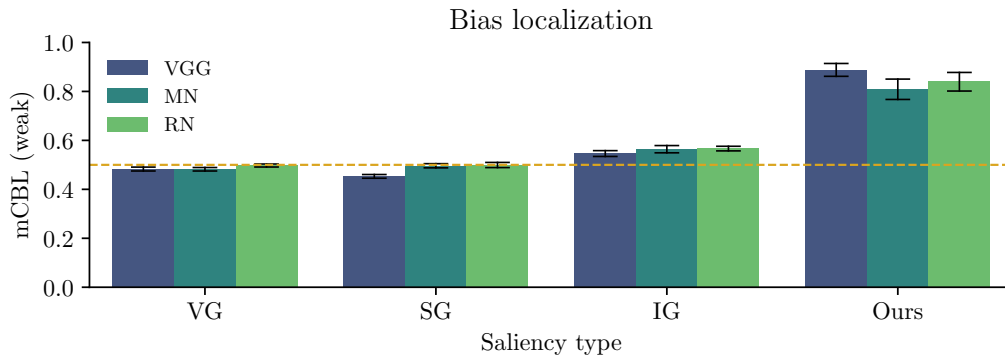


Figure 16: Comparison of localization capability for different network architectures ($x$ axis) using the CBL metric. All values are aggregated over ten bias digits, five different bias textures and five training set generation seeds. The error bars indicate the standard deviation over the bias digits.

19

## 7.6    Effect of Optimization Parameters

In Fig. 17 we report the effect of optimization parameters for grid saliency on context biased detection and localization performance (CBD, CBL). Influence of different optimization parameters on the quality of the obtained context explanations was evaluated on the basis of parameter sweeps around the chosen setting for the loss weighting $\lambda$, the learning rate and momentum, as well as initialization of the saliency mask (red points in Fig. 17). In our experiments in the main paper, all optimization parameters were set up by jointly looking at the two loss term values in Eq. 2, Sec. 3.1 and visual inspection of saliencies over a small image subset.

We notice that grid saliency shows comparable performance over a broad space of parameter settings (observing mostly smooth degradation with suboptimal parameter choices), with $\lambda$ clearly controlling the trade off between bias detection and localization quality (higher $\lambda$ value leads to a smaller salient region, see Sec. 3.1 of the main paper for method details). One sees, that the chosen parameter setting (red points in Fig. 17) balances well CBD and CBL performances (each high).
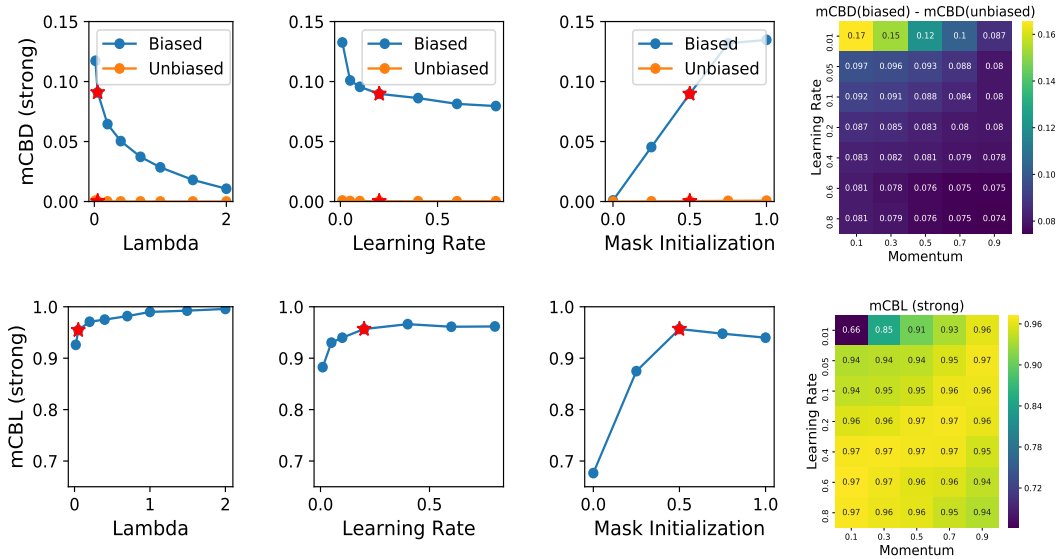


Figure 17: Effect of the grid saliency optimization parameters on the quality of context bias detection (mCBD) and localization (mCBL). mCBD (top row) and mCBL (bottom row) where computed for different settings of optimization parameters (from left to right): weighting lambda, learning rate, grey value initialization. All values are averaged over five randomly selected bias digits, five bias textures, and three training set generation seeds. The optimization parameters used in our experiments are depicted as red points. One sees, that the chosen parameter setting balances well mCBD and mCBL performances (each high). Note that for mCBD a high difference between biased and unbiased is beneficial.

# 8 Cityscapes Experiments

## 8.1 Experiment Setup

**Cityscapes Dataset.** A central motivation of this study is to better understand the behavior of segmentation networks. Cityscapes images describe rich scenes that make it arguably easy for segmentation models to learn clear context biases. We use the 500 (finely annotated) validation images available for Cityscapes as a starting point. We choose to analyze the validation data instead of training data in order to highlight behaviors that are likely to occur in a real world deployment scenario.

For the semantic segmentation, we have used the state-of-the-art network Deeplabv3+ [2] with a Mobilenetv2 backbone [1]. The weights were obtained from the original Deeplab repository[2].

**Implementation details.** The computation pipeline used to derive Fig. 6 in the main paper and its detailed version in Fig. 20 requires a selection of hyperparameters and design choices, which we explain here in more detail.

The saliency mask $M^*_{grid}$ is obtained by optimizing Eq. 2 with SGD, thus there is no guarantee of convergence to a global optimum. The loss function used in the optimization procedure (see in Eq.2) seeks to find a balance (partially controlled by $\lambda$) between penalizing the size of the produced saliency region and the preservation loss, which measures how well the softmax scores inside the request mask $R$ were restored to (at least) their initial values prior to perturbation, i.e. removing the image background. Typically, this preservation loss can be interpreted as a percentage relating to how much of the original softmax score was restored. So, for a loss of 0.1 or smaller, 90% or more of the original softmax activation scores must be restored. Samples that do not converge to a preservation loss of 0.1 or smaller are ignored in the computation. Similarly to the synthetic dataset, bilinear upsampling is used to upsample the optimized mask to the input image size.

The preservation loss in Eq. 2 is by definition normalized by the size of $R$, thus the size of $R$ does not directly influence the optimization convergence. In Fig. 18 we show the effect of the size of $R$, where context saliencies for each request mask $R$ were obtained with the same optimization parameters. Independent of the request mask $R$ size, for all riders salient context always falls on bikes.



Figure 18: Influence of the size of the request mask $R$ on context explanations provided by grid saliency. Independent of the request mask $R$ size, for all riders salient context always falls on bikes. Note that the context grid saliency explanations were calculated for each request region separately and are only combined together for the visualization.

Fig. 19 shows some intermediate results for the analysis of a single frame, which illustrates the method and the optimization in a slightly more detailed way. The components include (a) the input

---

[2]`http://download.tensorflow.org/models/deeplabv3_mnv2_cityscapes_train_2018_02_05.tar.gz`

a) Input frame  b) Softmax output on input frame

c) Perturbed context  d) Softmax output on peturbed context

e) Optimized image  f) Softmax output on optimized image

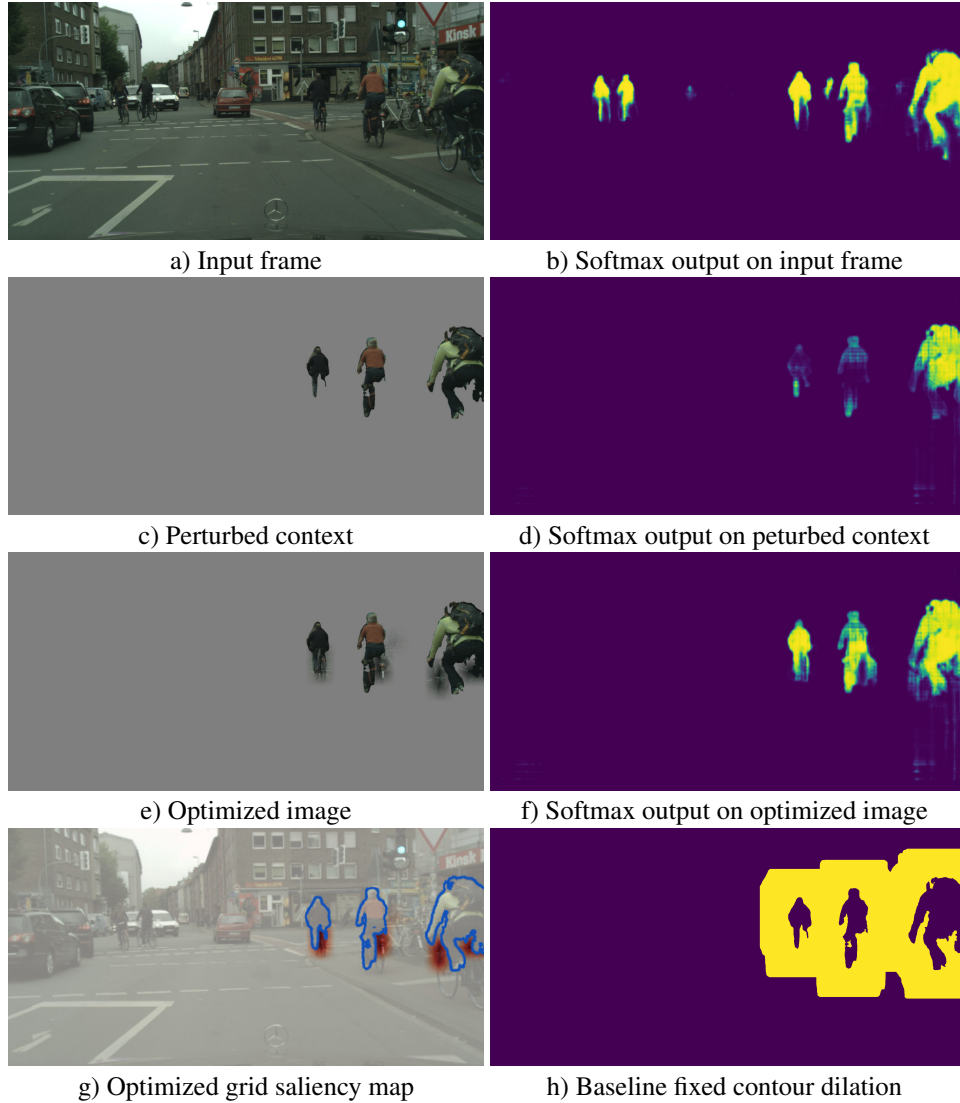g) Optimized grid saliency map  h) Baseline fixed contour dilation

Figure 19: Intermediate visualizations for the rider class request mask on an input frame.

frame, (b) the softmax scores for the "rider" class on the input frame, (c) the image with perturbed context given a request mask for the "rider" class, (d) the softmax output score of (c), (e) the optimized image (request mask plus learned background context), (f) the output softmax scores on (e) as an input, (g) the learned grid saliency, and (h) the dilation around the prediction used for the comparison to baseline.

Generally, penalizing the size of the saliency region is not sufficient to avoid spurious activations. We find that in general, the saliency mask at convergence represents a superset of the important pixels, i.e. it is often possible to slim down the saliency further without severely harming the preservation loss. Rather than adding regularizers directly in the loss function as [6], we simply clip all mask activation values smaller than or equal to 0.2 back to 0.0 for each step of the optimization. We find this leads to considerably less noisy activations and more focused and spatially coherent saliencies. Moreover, we manage potential border artifacts of the network predictions by eroding the request mask of each instance with a $3 \times 3$ erosion kernel.

Optimizing a low resolution saliency mask makes it hard to deal with small object instances in images. Specifically, a single pixel of the coarse mask already corresponds to a relatively wide spatial context for small objects in the background of a scene. Because we count class labels in the activated saliency, this adds significant noise to the resulting statistics. In order to counteract this, we remove

22

any connected components of the predicted object instances that have a pixel count smaller than $10k$ (or $625$ for context explanations of erroneous predictions) from our request masks. This filters out a considerable number of images from the heatmap computation. Table 2 and Table 3 count the number of samples available for each class that pass the selection criteria used for the experiments.

The baseline metric dilates the contours of the predicted instances in a given image that pass the minimum size requirement using a total of $400k$ pixels for the dilated region. If multiple instances are present, these dilation pixels are shared amongst them. The number was chosen so as to be able to capture a meaningful portion of an object's spatial context, even if said object is large and in the foreground. An example baseline mask with the fixed dilation can be seen in Fig. 19 (h).

Table 2: Number of images in the Cityscapes validation dataset used to compute context grid saliency of the model predictions for each request class. For each request class, an image is included if and only if its computed grid saliency is not empty or invalid. This is the case only if the following conditions are met: 1. The request mask for the image contains at least one connected component larger than $10k$ pixels after applying the border erosion kernel. 2. The grid saliency computation does not converge to an empty context saliency mask. 3. The saliency optimization procedure converges with a preservation loss of $0.1$ or less.

| Request (prediction) class | Number of samples |
|---|---|
| rider | 39 |
| car | 249 |
| person | 45 |
| bicycle | 105 |
| motorcycle | 16 |
| pole | 162 |
| traffic sign | 81 |

Table 3: Number of images in the Cityscapes validation dataset used to compute context grid saliency of the correct and erroneous model predictions. For each class, an image is included if and only if its computed grid saliency is not empty or invalid. This is the case only if the following conditions are met: 1. The request mask for the image contains at least one connected component larger than $625$ pixels after applying the border erosion kernel. Here the request mask consists of the intersection between the ground truth specified in the first column and the model prediction specified in the second column. As these segments are typically much smaller and more disconnected, the above threshold was reduced to from $10k$ to $625$ pixels. 2. The grid saliency computation does not converge to an empty context saliency mask. 3. The saliency optimization procedure converges with a preservation loss of $0.1$ or less.

| GT class | Prediction class | Number of samples |
|---|---|---|
| rider | rider | 160 |
| | person | 67 |
| person | person | 217 |
| | rider | 45 |
| | car | 29 |
| bicycle | bicycle | 262 |
| | motorcycle | 12 |
| car | car | 323 |
| | motorcycle | 8 |
| | person | 33 |
| motorcycle | motorcycle | 44 |
| | bicycle | 18 |
| | car | 9 |
| pole | pole | 410 |
| | building | 324 |
| | vegetation | 145 |

## 8.2 Further Quantitative Examples

Fig. 21 shows additional representative qualitative examples of our method for a variety of request classes. Some interesting effects include saliency activations on road signs marked on the pavement for the car class (Fig. 21 (e)), as well as the behavior of the method for occluded objects. As can be seen in subfigures (b) and (d), the grid saliency for occluded objects tends to activate more strongly along the entire contour of the object. In particular, subfigure (d) contains two objects of similar size, pose and background, where only the second is partially occluded.

In certain cases, the prediction of the image with completely perturbed context is already so good so that no gradients are available for the preservation loss, and the optimization then focuses on minimizing the total loss by simply getting rid of the saliency activation altogether. This typically results in large object instances of clearly identifiable classes such as cars or pedestrians with no saliency activations. Fig. 21 i) illustrates this for the left and center group of pedestrians. We consider this effect in itself to be a valid solution. An empty context saliency means that the requested object is self-containing and context is not necessary for its segmentation.
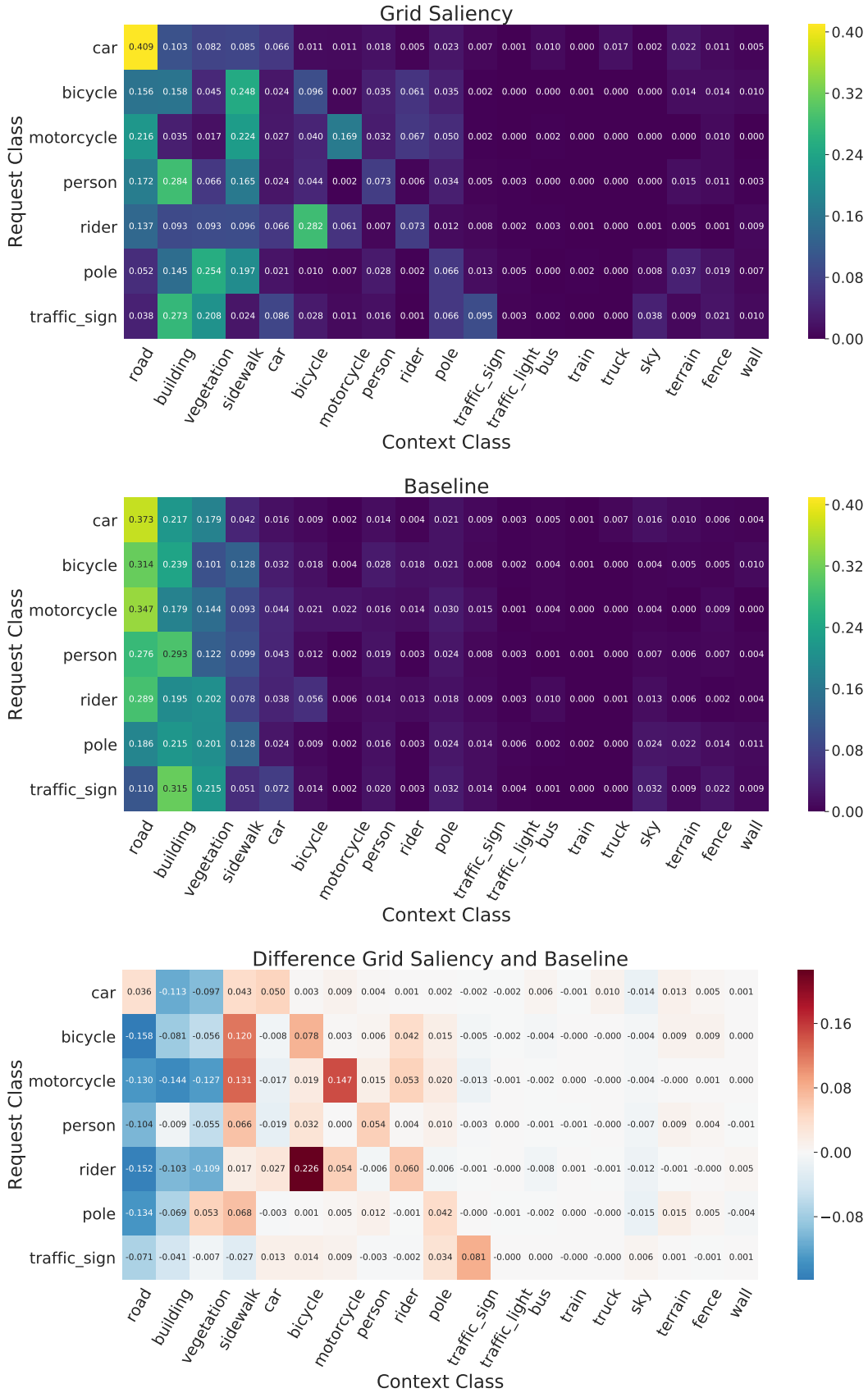
Figure 20: Comparison of saliency distributions for a selected set of classes across the Cityscapes validation data. Note that the rows of the *Grid Saliency* and *Baseline* maps do not sum up to 1, because only a subset of classes is considered. 25
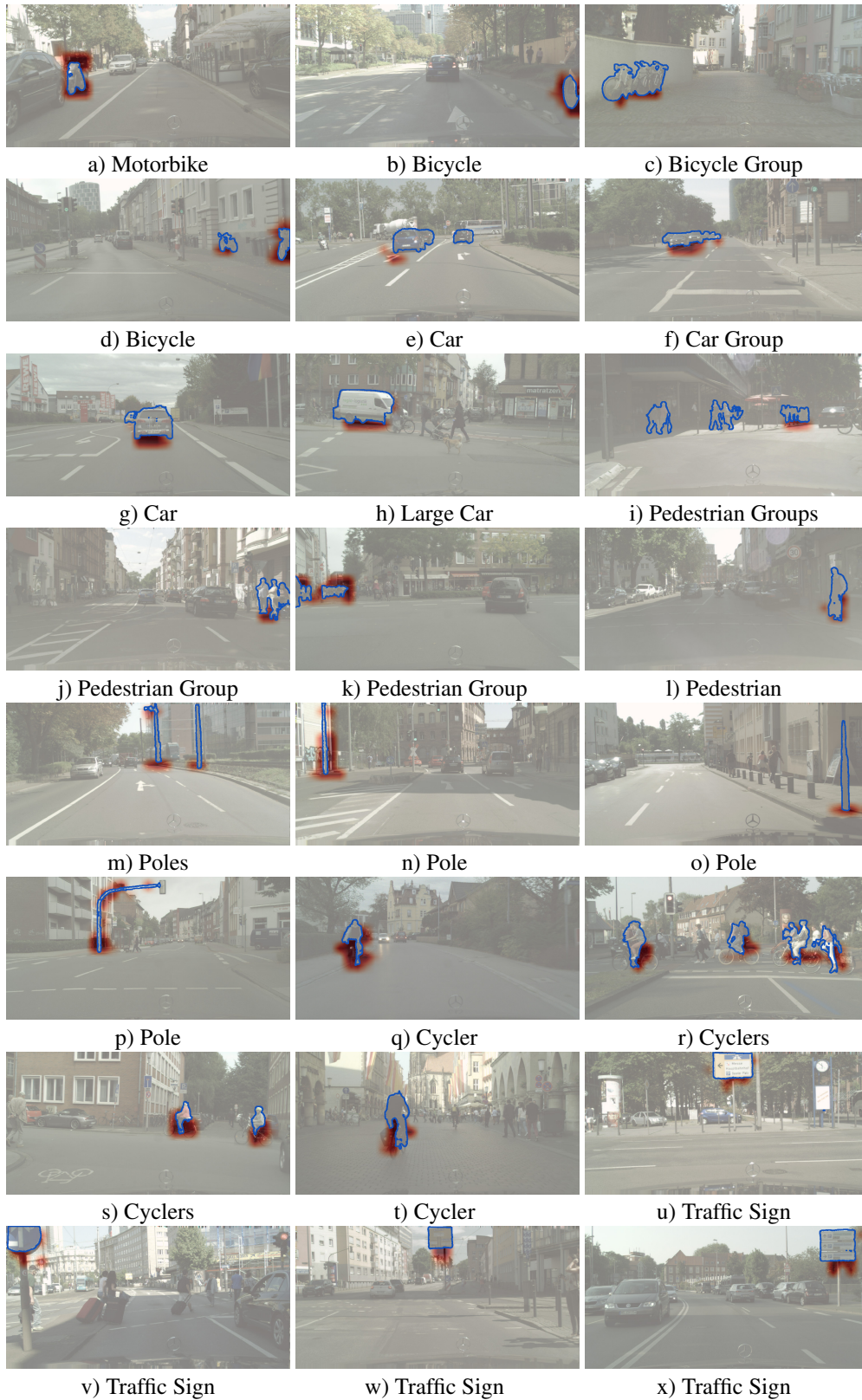
a) Motorbike      b) Bicycle      c) Bicycle Group

d) Bicycle      e) Car      f) Car Group

g) Car      h) Large Car      i) Pedestrian Groups

j) Pedestrian Group      k) Pedestrian Group      l) Pedestrian

m) Poles      n) Pole      o) Pole

p) Pole      q) Cycler      r) Cyclers

s) Cyclers      t) Cycler      u) Traffic Sign

v) Traffic Sign      w) Traffic Sign      x) Traffic Sign

Figure 21: Cityscapes examples supporting the class statistics.

# 9 Additional Results on MS COCO

In order to show that our grid saliency can be applied on different types of real data and to various network architectures, we have repeated the experiments from Sec. 4.2 of the main paper on the MS COCO images [23] with different backbones for the state-of-the-art Deeplabv3+ [2] network. In particular, we have chosen MobileNetv2 [1] and Xception (XC) [56] as backbones due to their different structure. All implementation details and optimization parameter settings are borrowed from the Cityscapes experiments.

In Fig. 22, we show some examples of context explanations on MS COCO obtained with our grid saliency method. We observe that grid saliency provides sensible and coherent explanations for network decision making, which reflect semantic dependencies present in the data (e.g. boat appears on the water). Context explanations produced by grid saliency can be also utilized to compare architectures with respect to their capacity to either learn or to be invariant towards context. E.g., in Fig. 22 the segmentation network with MobileNetv2 (MN) backbone learnt to rely more on context in contrast to its variant with a more powerful Xception (XC) backbone, which, for example, does not look at rails as context to segment the train.
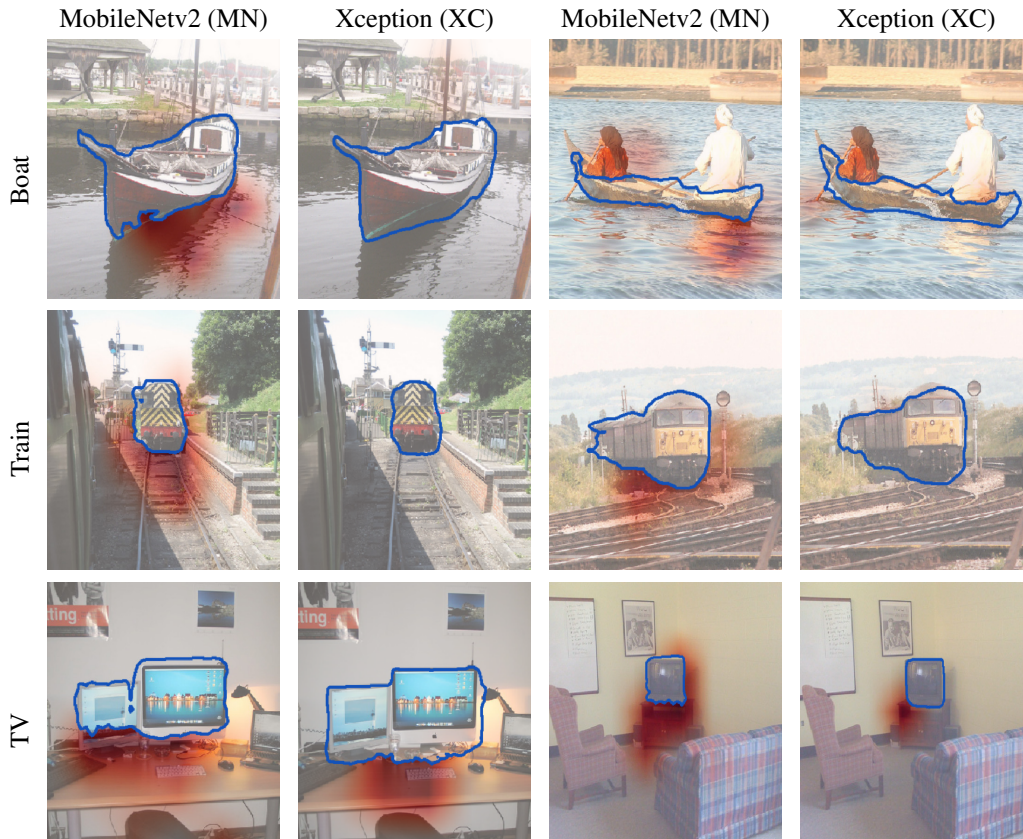


Figure 22: Context explanations produced by grid saliency on MS COCO [23] for the Deeplabv3+ [2] network with MobileNetv2 [1] and Xception (XC) [56] as backbones.