

I Bet You Are Wrong: Gambling Adversarial Networks for Structured Semantic Segmentation

Laurens Samson^{1,2}, Nanne van Noord¹, Olaf Booi², Michael Hofmann², Efstratios Gavves¹, and Mohsen Ghafoorian²

¹University of Amsterdam, Amsterdam, Netherlands
Email: {n.j.e.vannoord, e.gavves}@uva.nl

²TomTom, Amsterdam, Netherlands
Email: firstname.lastname@tomtom.com

Abstract

Adversarial training has been recently employed for realizing structured semantic segmentation, in which the aim is to preserve higher-level scene structural consistencies in dense predictions. However, as we show, value-based discrimination between the predictions from the segmentation network and ground-truth annotations can hinder the training process from learning to improve structural qualities as well as disabling the network from properly expressing uncertainties.

In this paper, we rethink adversarial training for semantic segmentation and propose to formulate the fake/real discrimination framework with a correct/incorrect training objective. More specifically, we replace the discriminator with a “gambler” network that learns to spot and distribute its budget in areas where the predictions are clearly wrong, while the segmenter network tries to leave no clear clues for the gambler where to bet. Empirical evaluation on two road-scene semantic segmentation tasks shows that not only does the proposed method re-enable expressing uncertainties, it also improves pixel-wise and structure-based metrics.

1. Introduction

In the past years, deep neural networks have obtained substantial success in various visual recognition tasks including semantic segmentation [12, 15]. Despite the success of the frequently used (fully) convolutional neural networks [31] on semantic segmentation, they lack a built-in mechanism to enforce global structural qualities. For instance, if the task is to detect a single longest line among several linear structures in the image, then a CNN is prob-

ably not able to properly handle such global consistency and will likely give responses on other candidate structures. This stems from the fact that even though close-by pixels share a fair amount of receptive field, there is no designated mechanism to explicitly condition the prediction at a specific location on the predictions made at other related (close-by or far) locations, when training with a pixel-level loss. To better preserve structural quality in semantic segmentation, several methods incorporate graphical models such as conditional random fields (CRF) [26, 49, 42], or use specific topology targeted engineered loss terms [1, 38]. More recently, adversarial training [16] schemes are being explored [32, 21, 13], where a discriminator network learns to distinguish the distributions of the network-provided dense predictions (fake) and ground-truth labels (real), which directly encourages better inter-pixel consistencies in a learnable fashion. However, as we will show, the visual clues that the discriminator uses to distinguish the fake and real distributions are not always high-level geometrical properties. For instance, a discriminator might be able to leverage the prediction values to contrast the fuzzy fake predictions with the crisp zero/one real prediction values to achieve an almost perfect discrimination accuracy.

Such value-based discrimination results in two undesirable consequences: 1) The segmentation network (“segmenter”) is forced to push its predictions toward zeros and ones and pretend to be confident to mimic such a low-level property of real annotations. This prevents the network from expressing uncertainties. 2) In practice, the softmax probability vectors can not get to exact zeros/ones that requires infinitely large logits. This leaves a permanent possibility for the discriminator to scrutinize the small - but still remaining- value gap between the two distributions, making it needless to learn the more complicated geometrical

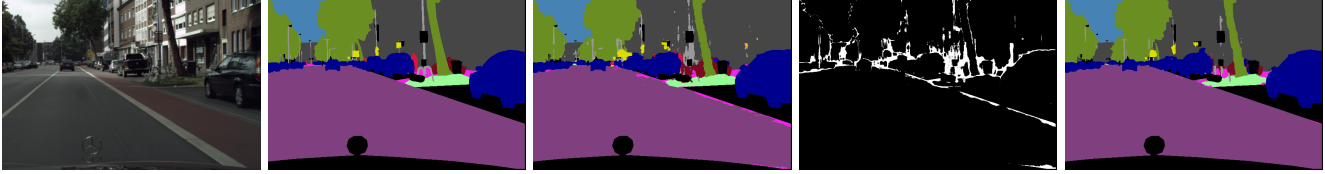


Figure 1: From left to right: sample image from Cityscapes [6], corresponding ground-truth image, predictions from U-Net trained with cross-entropy loss, betting map from the gambler network, predictions from the gambling adversarial nets. Notice e.g. spotted and resolved artefact in predictions from the cross-entropy trained U-Net in bottom right and right side of the road. Best visible zoomed-in on a screen.

discrepancies. This hinders such adversarial training procedures to reach their full potential in learning the scene structure.

The value-based discrimination inherently stems from the fake/real discrimination scheme employed in adversarial structured semantic segmentation. Therefore, we aim to study a surrogate adversarial training scheme that still models the higher level prediction consistencies, but is not trained to directly contrast the real and fake distributions. In particular, we replace the discriminator with a “gambler” network, that given the predictions of the segmenter and a limited budget, learns to spot and invest in areas where the predictions of the network are likely wrong. Put another way, we reformulate the fake/real discrimination problem into a correct/incorrect distinction task. This prevents the segmenter network from faking certainty, since a wrong confident prediction caught by the gambler, highly penalizes the segmenter. See Figures 1 and 2 for getting an overview.

Following are the main contributions of the paper:

- We propose gambling adversarial networks as a novel adversarial training scheme for structured semantic segmentation.
- We show that the proposed method resolves the usual adversarial semantic segmentation training issue with faking confidence.
- We demonstrate that this reformulation in the adversarial training improves the semantic segmentation quality over the baselines, both in pixel-wise and structural metrics on two semantic segmentation datasets, namely the Cityscapes [6] and Camvid [2] datasets.

2. Related work

Structure-preserving semantic segmentation. Several methods have been proposed that use specific loss terms which are targeted at preserving topologies [1, 38, 33] or use graphical models [23, 26, 5, 49, 41, 35, 28, 27, 22, 42] such

as CRFs that model unary, pairwise and/or higher-order potentials, either as a post-processing at the inference time or as integrated training refinement steps. Hand-engineering differentiable targeted loss terms for every desirable structural property is not always feasible in practice. On the other hand, using graphical models either confines the consistency improvements to model low-level features in a local context or imposes high computational costs.

Adversarial semantic segmentation. Adversarial training schemes have been extensively employed in the literature to impose structural consistencies for semantic segmentation [21, 32, 18, 8, 19, 25, 34, 46, 29, 40, 37]. Luc et al. [32] incorporate a discriminator network trained to distinguish the real labels and network-produced predictions. Involving the segmenter in a minimax game with the discriminator motivates the network to bridge the gap between the two distributions and consequently having higher-level consistencies in predicted labels.

More recently, it has been shown that the training dynamics of paired image-to-image translation in general [44, 43] and adversarial semantic segmentation specifically [13, 20, 45, 47], can be improved using paired real/fake embedding losses.

Our method is similar to the aforementioned adversarial formulations in the sense that it also employs a critic network that perceives the whole prediction map, consequently enabling it to model inter-pixel dependencies, and is similarly involved in a minimax game with the segmentation network. Similar to the embedding loss adversarial training, we also leverage the pairing between predictions and ground-truth in our adversarial training, with the difference that we incorporate the ground-truth not as an input but as a supervision for our gambler network. More generally, our method differs in the defined minimax game formulation; the gambler is trained to learn to spot the likely incorrect predictions, while the segmenter is trained to leave as little (structural) clues as possible for the gambler to make an easily profitable investment.

Luc et al. [32] also discuss the value-based discrimination issue, which they attempt to alleviate by feeding

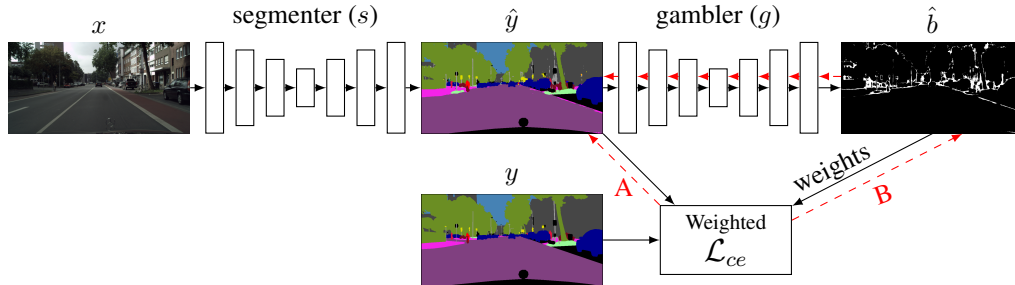


Figure 2: An overview of gambling adversarial networks. The solid black arrows indicate the forward pass. The red dashed arrows represent the two gradient flows of the weighted cross-entropy loss. Gradient flow A provides pixel-level feedback independent of other pixel predictions. Gradient flow B, going through the gambler network, enables feedback reflecting the inter-pixel and structural consistency.

the discriminator with a Cartesian product of the prediction maps and the input image channels. However, their followed strategy resulted in no improvements as reported. This can be attributed to remaining value-based evidence based on values distribution granularity. For instance, a very tiny response to a first-layer edge detector, in this case, can already signify a fake data sample.

Hard-sample mining. Our method is also closely related to the literature on class-imbalance/hard-sample mining. Class-imbalance is another inherent difficulty that needs to be properly tackled when dealing with problems/datasets with imbalanced semantic classes, as is often the case in semantic segmentation. Synthetic minority over-sampling technique (SMOTE) [4] and Mean/median-frequency balancing [10] are common simple strategies that over-sample or scale the loss terms corresponding to the under-represented classes. More recently, focal loss [30] and loss max-pooling [3] improve over the aforementioned by distinguishing between rarity and difficulty; not all samples from a frequent class are easy and not all samples belonging to infrequent classes are difficult. Therefore, focal loss and loss max-pooling address the more generic problem of hard-sample mining. However, the main issue with both is their inherent limitation dealing with label noise and/or ambiguities in the underlying semantics. We can view our gambling adversarial networks as an adversarially learned version of focal loss; the gambler learns to bet on (i.e. up-weight) the samples that it perceives as more difficult and/or more likely to be wrong in predictions from the segmenter. Such a learned approach can alleviate the label noise problem, as a learned network over noisy labels may be able to generalize beyond the noise level in the training dataset [36, 14] or at least soften the erroneous strong weight-increase for noisy samples. Furthermore, focal loss is derived in a pixel-wise manner and therefore cannot provide structural feedback to the segmentation network.

3. Method

In this section, the proposed method, gambling adversarial networks, is described. First, we present the usual adversarial training formulation for structured semantic segmentation and discuss the potential issues with it. Thereafter, we describe gambling adversarial networks and its reformulation of the former.

3.1. Conventional adversarial training

In the usual adversarial training formulation, the discriminator learns to discriminate the ground-truth (real) from the predictions provided by the network (fake). By involving the segmenter in a minimax game, it is challenged to improve its predictions to provide realistic-looking predictions to fool the discriminator [32]. In semantic segmentation, such an adversarial training framework is often employed with the aim to improve the higher-level structural qualities, such as connectivity, inter-pixel (local and non-local) consistencies and smoothness. The minimax game is set-up by forming the following loss terms for the discriminator and segmenter:

$$\mathcal{L}_d(x, y; \theta_s, \theta_d) = \mathcal{L}_{bce}(d(x, s(x; \theta_s); \theta_d), 0) + \mathcal{L}_{bce}(d(x, y; \theta_d), 1), \quad (1)$$

where x and y are the input image and the corresponding label-map, $s(x; \theta_s)$ is the segmenter’s mapping of the input image x to a dense segmentation map parameterized by θ_s , $d(x, y; \theta_d)$ represents the discriminator operating on segmentations y , conditioned on input image x and the binary cross-entropy is defined as $\mathcal{L}_{bce}(\hat{y}, y) = -(y \log \hat{y} + (1 - y) \log(1 - \hat{y}))$, where \hat{y} and y are the prediction and label respectively.

Typically, the loss function for the segmenter is a combination of low-level (pixel-wise) and high-level (adversarial) loss terms [21, 32]:

$$\mathcal{L}_s(x, y; \theta_s, \theta_d) = \mathcal{L}_{ce}(s(x; \theta_s), y) + \lambda \mathcal{L}_{bce}(d(x, s(x; \theta_s); \theta_d), 1), \quad (2)$$

where λ is the importance weighting of the adversarial loss, being the recommended non-saturating reformulation of the original minimax loss term to prevent vanishing gradients [16, 11]. The pixel-level cross-entropy loss \mathcal{L}_{ce} optimizes all the pixels independently of each other by minimizing $\mathcal{L}_{ce}(\hat{y}, y) = -\frac{1}{wh} \sum_{i,j}^{w,h} \sum_k^c y_{i,j,k} \log \hat{y}_{i,j,k}$, where w and h are the width and the height of image x and c is the number of classes in the dataset.

Recently, the usual adversarial training for structured semantic segmentation was suggested to be modified [13, 45, 43] by replacing the binary cross-entropy loss as the adversarial loss term for the segmenter, with a fake/real paired embedding difference loss, where the embeddings are extracted from the adversarially trained discriminator. To be more specific, the adversarial loss term in Equation (2) is replaced by the following embedding loss:

$$\mathcal{L}_{emb}(x, \hat{y}, y; \theta_d) = \|d_e(x, \hat{y}; \theta_d) - d_e(x, y; \theta_d)\|_2, \quad (3)$$

where the function $d_e(x, y; \theta)$ represents the extracted features from a particular layer in the discriminator. As shown in the EL-GAN method, this could significantly stabilize training [13].

Ideally, the discriminator’s decisions are purely based on the structural differences between the real and the fake predictions. However, in semantic segmentation, it is often possible for the discriminator to perfectly distinguish the labels from the predictions based on the values. The output of the segmenter is a softmax vector per pixel, which assigns a probability to every class that ranges between zero and one. In contrast, the values in the ground-truth are either zeros or ones due to the one-hot encoding. Such value-based discrepancy can yield unsatisfactory gradient feedback, since the segmenter might be forced to mimic the one-hot encoding of the ground-truth instead of the global structures. Additionally, the value-based discrimination is a never-ending problem since realizing exact ones and zeros requires infinite large logits, however, in practise, the segmenter always leaves a small value-based gap that can be exploited by the discriminator. Another undesired outcome is the loss of ability to express uncertainties, since all the predictions will converge towards a one-hot representation to bridge the value-based gap between the one-hot labels and probabilistic predictions.

3.2. Gambling Adversarial Networks

To prevent the adversarial network from utilizing the value-based discrepancy, we propose gambling adversar-

ial networks, which focuses solely on improving the structural inconsistencies. Instead of the usual real/fake adversarial training task, we propose to modify the task to learn to distinguish incorrect predictions given the whole prediction map. Different from a discriminator, the critic network (gambler) does not observe the ground-truth labels, but solely the RGB-image in combination with the prediction of the segmentation network (segmenter). Given a limited investment budget, the gambler predicts an image-sized betting map, where high bets indicate pixels that are likely incorrectly classified, given the contextual prediction clues around it. Since the gambler receives the entire prediction, structurally ill-formed predictions, such as non-smoothness, disconnectivities and shape-anomalies are clear visual clues for profitable investments for the gambler. An overview of gambling adversarial networks is provided in Figure 2.

Similar to conventional adversarial training, the gambler and segmenter play a minimax game; The gambler maximizes the expected weighted pixel-wise cross-entropy where the weights are determined by its betting map, while the segmenter attempts to improve its predictions such that the gambler does not have clues where to bet:

$$\mathcal{L}_g(x, y; \theta_s, \theta_g) = -\frac{1}{wh} \sum_{i,j}^{w,h} g(x, s(x; \theta_s); \theta_g)_{i,j} \mathcal{L}_{ce}(s(x; \theta_s)_{i,j}, y_{i,j}), \quad (4)$$

where $g(x, s(x; \theta_s); \theta_g)_{i,j}$ is the amount of budget the gambler invests on position (i, j) .

The segmenter network minimizes the opposite:

$$\mathcal{L}_s(x, y; \theta_s, \theta_g) = \mathcal{L}_{ce}(s(x; \theta_s), y) - \mathcal{L}_g(x, y; \theta_s, \theta_g). \quad (5)$$

Similar to conventional adversarial training, the segmentation network optimizes a combination of loss terms: a per-pixel cross-entropy loss and an inter-pixel adversarial loss. It should be noted that the gambler can easily maximize this loss by betting infinite amounts on all the pixels. Therefore, it is necessary to limit the budget the gambler can spend. We accommodate this by turning the betting map into a smoothed probability distribution:

$$g(x, \hat{y}; \theta_g)_{i,j} = \frac{g_\sigma(x, \hat{y}; \theta_g)_{i,j} + \beta}{\sum_{k,l}^{w,h} g_\sigma(x, \hat{y}; \theta_g)_{k,l} + \beta}, \quad (6)$$

where β is a smoothing factor and $g_\sigma(x, \hat{y}; \theta_g)_{i,j}$ represents the sigmoid output of the gambler network for pixel with the indices i, j . Smoothing the betting map regularizes the model to spread its bets over multiple pixels instead of focusing on a single location.

The adversarial loss causes two different gradient streams for the segmentation network, as shown in Figure 2, where the solid black and dashed red arrows indicate the

forward pass and backward gradient flows respectively. In the backward pass, the gradient flow A pointing directly towards the prediction provides pixel-wise feedback independent of the other pixel predictions. Meanwhile, the gradient flow B, going through the gambler network, provides feedback reflecting inter-pixel and structural consistencies.

4. Experimental results

In this section, we discuss the datasets and metrics for the evaluation of gambling adversarial networks. Thereafter, we describe the different network architectures for the segmenter and gambler networks and provide details for training. Finally, we report the results of our experiments.

4.1. Experimental setup

Datasets. We conduct experiments on two different urban road-scene semantic segmentation datasets, but hypothesize that the method is generic and can be applied to any segmentation dataset.

Cityscapes. The Cityscapes [6] dataset contains 2975 training images, 500 validation images and 1525 test images with a resolution of 2048×1024 consisting of 19 different classes, such as cars, persons and road signs. For pre-processing of the data, we down-scale the images to 1024×512 , perform random flipping and take random crops of 512×512 for training. Furthermore, we perform intensity jittering on the RGB-images.

Camvid. The urban scene Camvid [2] dataset consists of 429 training images, 101 validation images and 171 test images with a resolution of 960×720 . We apply the same data augmentations as described above, except that we do not perform any down-scaling.

Metrics. In addition to the mean intersection over union (IoU), we also quantify the structural consistency of the segmentation maps. Firstly, we compute the BF-score [7], which measures whether the contours of objects in the predictions match with the contours of the ground-truth. A point on the contour line is a match if the distance between the ground-truth and prediction lies within a toleration distance τ , which we set to 0.75 % of the image diagonal as suggested in [7]. Furthermore, we utilize a modified Hausdorff distance to quantitatively measure the structural correctness [9]. We slightly modify the original Hausdorff distance, to prevent it from being overwhelmed by outliers:

$$d_H(X, Y) = \frac{1}{2} \sum \left\{ \frac{1}{|X|} \sum_{x \in X} \inf_{y \in Y} d(x, y), \frac{1}{|Y|} \sum_{y \in Y} \inf_{x \in X} d(x, y) \right\}, \quad (7)$$

where X and Y are the contours of the predictions and labels from a particular class and $d(x, y)$ is the Euclidean distance. We average the score over all the classes that are present in the prediction and the ground-truth.

Network architectures. For comparison, we experiment with two well-known baseline segmentation network architectures. Firstly, a U-Net [39] based architecture as implemented in Pix2Pix [21], which is an encoder-decoder structure with skip connections. The encoder consists of nine down-sampling blocks containing a convolutional layer with batch normalization and ReLu. The decoder blocks are the same, except that the convolutions are replaced by transposed convolutions. Furthermore, we conduct experiments with PSPNet [48], which utilizes a pyramid pooling module to capture more contextual information. Similar to [48], we utilize an ImageNet pre-trained ResNet-101 [17] as backbone.

For the gambler network, we utilize the same networks as the segmentation network. When training with the U-Net based architecture, the gambler network is identical except that it contains only six down-sampling blocks. For the PSPNet, the architecture of the gambler and segmenter are identical. For the baseline adversarial methods, we utilize the PatchGAN discriminator from Pix2Pix [21].

Training. For training the models, we utilize the Adam optimizer [24] with a linearly decaying learning rate over time. Similar to the conventional adversarial training, the gambler and segmenter are trained in an alternating fashion where the gambler is frozen when updating the segmenter and vice versa. Furthermore, we learned that as opposed to conventional adversarial training, our network does not require separate pre-training and in general, we observe that the training is less sensitive to hyperparameters. Details of the hyperparameters can be found in the supplementary material.

4.2. Results

Confidence expression. As discussed before, value-based discrimination encourages the segmentation network to mimic the one-hot vectors of the ground-truth, resulting in loss of ability to express uncertainty. We hypothesize that reformulating the fake/real discrimination in the adversarial training to a correct/incorrect distinction scheme will mitigate the issue. To verify this, the mean and standard deviation of the maximum class-likelihood value in every softmax vector for each pixel is tracked on the validation set over different training epochs and the results are depicted in Figure 3. We conducted this experiment with the U-Net based architecture on Cityscapes, but we observed the same phenomena with the other segmentation network and on the other dataset. One can observe that for both the standard

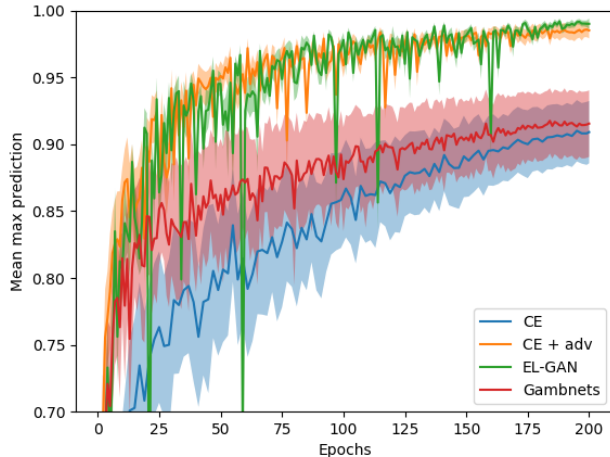


Figure 3: Mean maximum class-likelihoods (mean confidence) over time on the Cityscapes [6] validation set. Solid central curves and the surrounding shaded area represent the mean and standard deviation respectively.

Method	Mean max
Cross-entropy	$90.7 \pm 2.3 \%$
Cross-entropy + adversarial	$98.4 \pm 0.5 \%$
EL-GAN	$98.9 \pm 0.2 \%$
Gambling nets	$91.4 \pm 2.4 \%$

Table 1: Mean maximum value in every softmax vector on the Cityscapes [7] validation set averaged over the last 10 epochs.

adversarial training and EL-GAN that discriminate the real from the fake predictions, the predictions are converging towards one, with barely any standard deviation. For the gambling adversarial networks, the uncertainty of the predictions is well-preserved. In Table 1, the average mean maximum over the last 10 epochs is shown, which confirms that the gambling adversarial networks maintain the ability to express the uncertainty similar to the cross-entropy model, while the existing adversarial methods attempt to converge to a one-hot vector.

U-Net based segmenter. First, we compare the baselines with the gambling adversarial networks on the Cityscapes [6] validation set with the U-Net based architecture. The results in Table 2 show the gambling adversarial networks perform better on the pixel-wise metric (IoU), but also on the structural metrics. In Table 3, the IoU per class is provided for the same experiments. The gambling adversarial networks perform better on most of the classes. Moreover, performance particularly improves on the classes with finer structures, such as traffic light and person. In Table 4, we report the BF-score per class, where the gambling adversarial networks outperform the other methods on al-

Method	Mean IoU	BF-score	Hausdorff
CE	52.7	49.0	36.8
Focal loss [30]	56.2	55.3	30.2
CE + adv [32]	56.3	57.3	31.3
EL-GAN [13]	55.4	54.2	31.6
Gambling nets	57.9	58.5	27.6

Table 2: Results on Cityscapes [6] with U-Net based architecture [21] as segmentation network.

most all classes. Moreover, similar to the IoU, we observe the most significant improvements on the more fine-grained classes, such as rider and pole.

In Figure 4, one qualitative sample is depicted, in the supplementary material more samples are provided. The adversarial methods resolve some of the artifacts, such as the odd pattern in the car on the left. Moreover, the boundaries of the pedestrians on the sidewalk become more precise. We also provide an example betting map predicted by the gambler, given the predictions from the baseline model trained with cross-entropy in combination with the RGB-image. Note that the gambler bets on the badly shaped building in the prediction and responds to the artifacts in the car.

PSPNet segmenter. We conduct experiments with the PSPNet [48] segmenter on the Camvid [2] and Cityscapes [6] datasets. In Table 5, the results are shown on the Cityscapes validation set. Again, the gambling adversarial networks perform better than the existing methods, on both of the structure-based metrics as well as the mean IoU. In Figure 5, a qualitative sample is shown, more can be found in the supplementary material. The gambling adversarial networks provides more details to the traffic lights. Also, the structure of the sidewalk shows significant improvements over the predictions from the model trained with standard segmentation loss.

The quantitative results on the Camvid [2] test set are shown in Table 6. The gambling adversarial networks achieve the highest score on the structure-based metrics, but the standard adversarial training [32] performs best on the IoU. In the supplementary material, we provide qualitative results for the Camvid [2] test set and extra images for the aforementioned experiments.

5. Discussion

Correct/incorrect versus real/fake discrimination. We reformulated the adversarial real/fake discrimination task into training a critic network to learn to spot the likely incorrect predictions. As shown in Figure 3, the discrimination of real and fake causes undesired gradient feedback, since all the softmax vectors converge to a one-hot vector. We

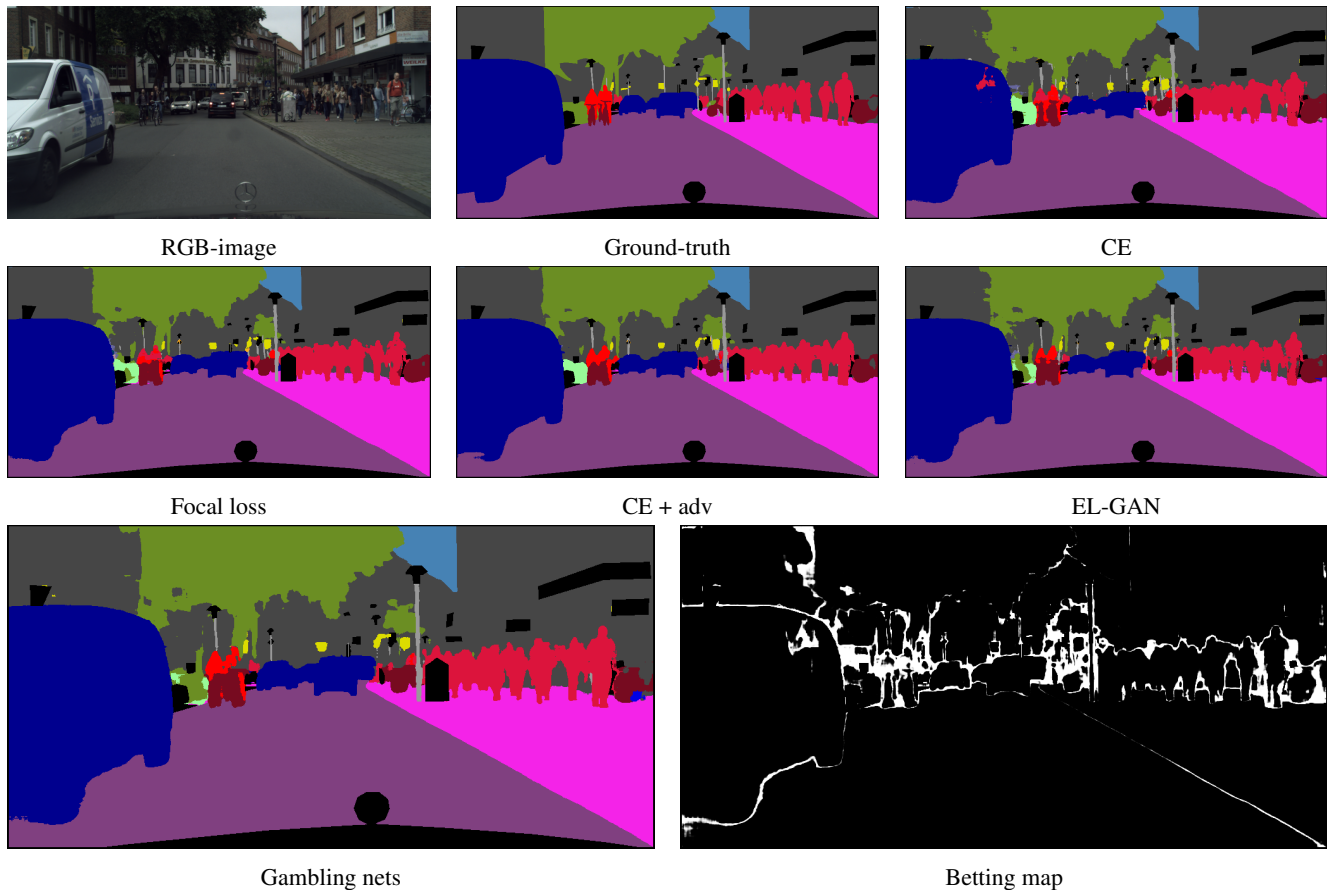


Figure 4: Qualitative results on Cityscapes [6] with the U-Net based architecture [21].

Method	road	swalk	build	wall	fence	pole	tlight	sign	veg.	ter.	sky	pers	rider	car	truck	bus	train	mbike	bike	mean
CE	95.2	68.4	84.4	26.0	30.9	43.0	38.9	51.3	87.2	50.3	91.5	59.0	32.6	85.5	22.8	43.2	19.2	15.4	57.4	57.2
Focal loss [30]	96.0	71.3	87.1	32.2	34.9	48.6	47.6	57.8	88.9	54.2	92.7	62.9	33.5	87.2	28.5	47.5	18.3	19.3	60.0	56.2
CE + adv [32]	95.9	72.7	83.5	28.9	35.2	49.8	47.8	59.3	89.0	54.8	92.3	66.4	38.4	87.2	27.8	41.4	15.3	20.3	62.5	56.3
EL-GAN [13]	96.1	71.1	86.8	33.5	37.0	48.7	46.6	57.3	88.9	53.6	92.9	62.6	34.4	87.1	26.0	38.3	16.3	17.8	58.9	55.4
Gambling	96.3	73.0	87.6	33.4	39.1	52.9	51.3	61.9	89.7	55.8	93.1	68.1	38.9	88.7	30.3	40.2	11.5	24.8	63.2	57.9

Table 3: IoU per class on the validation set of Cityscapes [6] with U-Net based architecture [21] as segmentation network

Method	road	swalk	build	wall	fence	pole	tlight	sign	veg.	ter.	sky	pers	rider	car	truck	bus	train	mbike	bike
CE	84.8	69.0	77.3	15.6	13.7	66.4	31.3	53.7	82.3	28.7	82.0	47.5	29.2	76.0	8.3	12.2	2.6	8.9	44.1
Focal loss [30]	87.2	72.4	80.7	19.7	16.0	71.0	40.1	62.3	86.1	35.8	84.8	51.6	32.0	79.4	9.2	18.0	4.3	12.1	50.5
CE + adv [32]	82.6	72.3	79.8	16.2	16.2	72.1	43.6	65.7	86.2	34.5	83.3	54.8	34.4	78.8	8.7	17.5	4.4	14.0	52.0
EL-GAN [13]	86.9	72.3	79.9	19.3	16.4	70.7	38.2	63.4	85.5	32.7	84.0	51.2	32.7	78.1	9.5	16.8	4.8	8.8	47.0
Gambling	87.4	74.3	81.3	20.7	18.6	74.0	45.7	67.8	87.2	35.4	85.4	57.0	38.8	80.0	11.2	19.3	4.4	15.6	52.9

Table 4: BF-score [7] per class on the validation set of Cityscapes [6] with U-Net based architecture [21] as segmentation network

empirically showed that this behavior is caused by a value-based discrimination of the adversarial network. Moreover, modifying the adversarial task to correct/incorrect discrimination solves several problems. First of all, the reason to apply adversarial training to semantic segmentation is to improve on the high-level structures. However, the value-

based discriminator is not only providing feedback based on the visual difference between the predictions and the labels, but also an undesirable value-based feedback. Moreover, updating the weights in a network with the constraint that the output must be a one-hot vector complicates training unnecessarily. Finally, the value-based discriminator

Method	Mean IoU	BF-score	Hausdorff
CE	72.4	69.0	19.4
Focal loss [30]	71.5	67.4	21.2
CE + adv [32]	68.0	67.0	20.9
EL-GAN [13]	71.3	67.0	21.2
Gambling nets	73.1	70.1	18.7

Table 5: Results on Cityscapes [6] with PSPNet [48] as segmentation network.

Method	Mean IoU	BF-score	Hausdorff
CE	72.5	71.8	17.9
Focal loss [30]	70.8	71.4	17.7
CE + adv [32]	72.7	72.7	17.1
EL-GAN [13]	70.1	69.6	19.1
Gambling nets	72.1	73.8	16.0

Table 6: Results on Camvid [2] with PSPNet [48] as segmentation network.

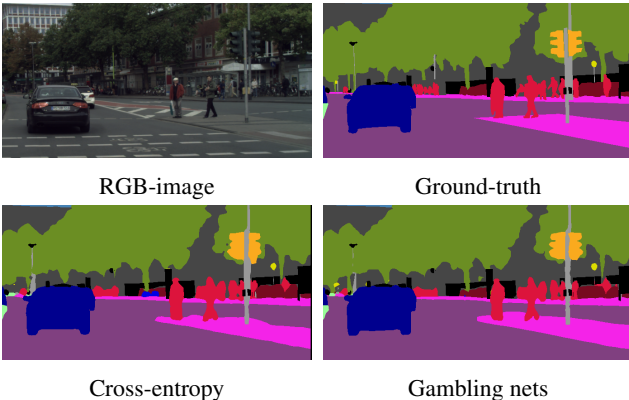


Figure 5: Qualitative results on Cityscapes [6] with PSPNet [48] as segmentation network

hinders the network from properly disclosing uncertainty. Both the structured prediction and expressing uncertainty can be of great value for semantic segmentation, e.g. in autonomous driving and medical imaging applications. However, changing the adversarial task to discriminating the correct from the incorrect predictions resolves the aforementioned issues. The segmentation network is not forced to imitate the one-hot vectors, which preserves the uncertainty in the predictions and simplifies the training. Although we still notice that the gambler sometimes utilizes the prediction values by betting on pixels where the segmenter is uncertain, we also obtain improvements on the structure-based metrics compared to the existing adversarial methods.

Gambling adversarial networks vs. focal loss The adversarial loss in gambling adversarial networks resembles

the focal loss [30], since both methods up-weight the harder samples that contain more useful information for the update. The focal loss is defined as: $\mathcal{L}_{foc}(y, \hat{y}, p_t) = -(1 - p_t)^\gamma y \log \hat{y}$, where p_t is the probability of the correct class and γ is a focusing factor, which indicates how much the easier samples are down-weighted. The advantage of the focal loss is that the ground-truth is exploited to choose the weights, however, the downside is that the focal loss might be over-pronouncing the ambiguous or incorrectly labeled pixels. The adversarial loss in gambling adversarial networks learns the weighting map, which can mitigate the noise effect. Moreover, the adversarial loss generates an extra flow of gradients (flow B), as observable in Figure 2. Gradient stream A provides information to the segmentation network independent of other pixel predictions similar to the focal loss, whereas gradient stream B provides gradients reflecting structural qualities, which is lacking in case of the focal loss.

Insights into betting maps Inspecting the betting maps (see for instance Figure 4), we observe that some of the bets correspond to the class borders, especially the ones that seemingly do not match the visual evidence in the underlying image or the expected shape of the object. We should note that even though there are chances that the ground-truth labels on the borders are different from the predictions, blindly betting on all the borders is not even close to the optimal policy. The clear bad structures in the predictions, e.g. the weird prediction of rider inside the car or the badly formed wall on the left side, are still more rewarding investments that are also being spotted by the gambler.

6. Conclusion

In this paper, we studied a novel reformulation of adversarial training for semantic segmentation, in which we replace the discriminator with a gambler network that learns to use the inter-pixel consistency clues to spot the wrong predictions. We showed that involving the segmenter in a minimax game with such a gambler results in notable improvements in structural and pixel-wise metrics, as measured on two road-scene semantic segmentation datasets.

References

- [1] A. BenTaieb and G. Hamarneh. Topology aware fully convolutional networks for histology gland segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 460–468. Springer, 2016. 1, 2
- [2] G. J. Brostow, J. Fauqueur, and R. Cipolla. Semantic object classes in video: A high-definition ground truth database. *Pattern Recognition Letters*, 30(2):88–97, 2009. 2, 5, 6, 8, 11

- [3] S. R. Buló, G. Neuhold, and P. Kotschieder. Loss max-pooling for semantic image segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 7082–7091. IEEE, 2017. 3
- [4] N. V. Chawla, K. W. Bowyer, L. O. Hall, and W. P. Kegelmeyer. Smote: synthetic minority over-sampling technique. *Journal of artificial intelligence research*, 16:321–357, 2002. 3
- [5] 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*, 40(4):834–848, 2017. 2
- [6] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, and B. Schiele. The cityscapes dataset for semantic urban scene understanding. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 3213–3223, 2016. 2, 5, 6, 7, 8, 11, 12, 13
- [7] G. Csurka, D. Larlus, F. Perronnin, and F. Meylan. What is a good evaluation measure for semantic segmentation?. In *BMVC*, volume 27, page 2013. Citeseer, 2013. 5, 6, 7
- [8] W. Dai, J. Doyle, X. Liang, H. Zhang, N. Dong, Y. Li, and E. P. Xing. Scan: Structure correcting adversarial network for chest x-rays organ segmentation. *arXiv preprint arXiv:1703.08770*, 1, 2017. 2
- [9] M.-P. Dubuisson and A. K. Jain. A modified hausdorff distance for object matching. In *Proceedings of 12th international conference on pattern recognition*, volume 1, pages 566–568. IEEE, 1994. 5
- [10] D. Eigen and R. Fergus. Predicting depth, surface normals and semantic labels with a common multi-scale convolutional architecture. In *Proceedings of the IEEE international conference on computer vision*, pages 2650–2658, 2015. 3
- [11] W. Fedus, M. Rosca, B. Lakshminarayanan, A. M. Dai, S. Mohamed, and I. Goodfellow. Many paths to equilibrium: Gans do not need to decrease a divergence at every step. *arXiv preprint arXiv:1710.08446*, 2017. 4
- [12] A. Garcia-Garcia, S. Orts-Escolano, S. Oprea, V. Villena-Martinez, and J. Garcia-Rodriguez. A review on deep learning techniques applied to semantic segmentation. *arXiv preprint arXiv:1704.06857*, 2017. 1
- [13] M. Ghafoorian, C. Nugteren, N. Baka, O. Booij, and M. Hofmann. El-gan: embedding loss driven generative adversarial networks for lane detection. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 0–0, 2018. 1, 2, 4, 6, 7, 8
- [14] M. Ghafoorian, J. Teuwen, R. Manniesing, F.-E. de Leeuw, B. van Ginneken, N. Karssemeijer, and B. Platel. Student beats the teacher: deep neural networks for lateral ventricles segmentation in brain mr. In *Medical Imaging 2018: Image Processing*, volume 10574, page 105742U. International Society for Optics and Photonics, 2018. 3
- [15] S. Ghosh, N. Das, I. Das, and U. Maulik. Understanding deep learning techniques for image segmentation. *arXiv preprint arXiv:1907.06119*, 2019. 1
- [16] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio. Generative adversarial nets. In *Advances in neural information processing systems*, pages 2672–2680, 2014. 1, 4
- [17] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. 5
- [18] W.-C. Hung, Y.-H. Tsai, Y.-T. Liou, Y.-Y. Lin, and M.-H. Yang. Adversarial learning for semi-supervised semantic segmentation. *arXiv preprint arXiv:1802.07934*, 2018. 2
- [19] Y. Huo, Z. Xu, S. Bao, C. Bermudez, A. J. Plassard, J. Liu, Y. Yao, A. Assad, R. G. Abramson, and B. A. Landman. Splenomegaly Segmentation using Global Convolutional Kernels and Conditional Generative Adversarial Networks. *Proceedings of SPIE*, 10574:10574 – 10574 – 7, 2018. 2
- [20] J.-J. Hwang, T.-W. Ke, J. Shi, and S. X. Yu. Adversarial structure matching loss for image segmentation. *arXiv preprint arXiv:1805.07457*, 2018. 2
- [21] P. Isola, J.-Y. Zhu, T. Zhou, and A. A. Efros. Image-to-image translation with conditional adversarial networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 1125–1134, 2017. 1, 2, 3, 5, 6, 7, 11, 12
- [22] M. Jaderberg, K. Simonyan, A. Vedaldi, and A. Zisserman. Deep structured output learning for unconstrained text recognition. *arXiv preprint arXiv:1412.5903*, 2014. 2
- [23] T. Joy, A. Desmaison, T. Ajanthan, R. Bunel, M. Salzmann, P. Kohli, P. H. Torr, and M. P. Kumar. Efficient relaxations for dense crfs with sparse higher-order potentials. *SIAM Journal on Imaging Sciences*, 12(1):287–318, 2019. 2
- [24] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 5, 11
- [25] S. Kohl, D. Bonekamp, H.-P. Schlemmer, K. Yaqubi, M. Hohenfellner, B. Hadaschik, J.-P. Radtke, and K. Maier-Hein. Adversarial Networks for the Detection of Aggressive Prostate Cancer. *arXiv preprint arXiv:1702.08014*, 2017. 2
- [26] P. Krähenbühl and V. Koltun. Efficient inference in fully connected crfs with gaussian edge potentials. In *Advances in neural information processing systems*, pages 109–117, 2011. 1, 2
- [27] M. Larsson, A. Arnab, F. Kahl, S. Zheng, and P. Torr. Learning arbitrary pairwise potentials in crfs for semantic segmentation. *terrain*, 4:4. 2
- [28] M. Larsson, A. Arnab, S. Zheng, P. Torr, and F. Kahl. Revisiting deep structured models for pixel-level labeling with gradient-based inference. *SIAM Journal on Imaging Sciences*, 11(4):2610–2628, 2018. 2
- [29] Z. Li, Y. Wang, and J. Yu. Brain Tumor Segmentation Using an Adversarial Network. In *International MICCAI Brainlesion Workshop*, pages 123–132. Springer, 2017. 2
- [30] T.-Y. Lin, P. Goyal, R. Girshick, K. He, and P. Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017. 3, 6, 7, 8

- [31] 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*, pages 3431–3440, 2015. 1
- [32] P. Luc, C. Couprie, S. Chintala, and J. Verbeek. Semantic segmentation using adversarial networks. *arXiv preprint arXiv:1611.08408*, 2016. 1, 2, 3, 6, 7, 8
- [33] Z. Mirikharaji and G. Hamarneh. Star shape prior in fully convolutional networks for skin lesion segmentation. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 737–745. Springer, 2018. 2
- [34] P. Moeskops, M. Veta, M. W. Lafarge, K. A. Eppenhof, and J. P. Pluim. Adversarial Training and Dilated Convolutions for Brain MRI Segmentation. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support*, pages 56–64. Springer, 2017. 2
- [35] S. T. Namin, M. Najafi, M. Salzmann, and L. Petersson. A multi-modal graphical model for scene analysis. In *2015 IEEE Winter Conference on Applications of Computer Vision*, pages 1006–1013. IEEE, 2015. 2
- [36] N. Natarajan, I. S. Dhillon, P. K. Ravikumar, and A. Tewari. Learning with noisy labels. In *Advances in neural information processing systems*, pages 1196–1204, 2013. 3
- [37] V. Nguyen, T. F. Y. Vicente, M. Zhao, M. Hoai, and D. Samaras. Shadow detection with conditional generative adversarial networks. In *ICCV: IEEE International Conference on Computer Vision*, pages 4520–4528. IEEE, 2017. 2
- [38] O. Oktay, E. Ferrante, K. Kamnitsas, M. Heinrich, W. Bai, J. Caballero, S. A. Cook, A. De Marvao, T. Dawes, D. P. O'Regan, et al. Anatomically constrained neural networks (acnns): application to cardiac image enhancement and segmentation. *IEEE transactions on medical imaging*, 37(2):384–395, 2017. 1, 2
- [39] O. Ronneberger, P. Fischer, and T. Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 5
- [40] S. K. Sadanandan, J. Karlsson, and C. Whlby. Spheroid segmentation using multiscale deep adversarial networks. In *ICCVW: IEEE International Conference on Computer Vision Workshops*, pages 36–41, Oct 2017. 2
- [41] A. G. Schwing and R. Urtasun. Fully connected deep structured networks. *arXiv preprint arXiv:1503.02351*, 2015. 2
- [42] S. Tsogkas, I. Kokkinos, G. Papandreou, and A. Vedaldi. Deep learning for semantic part segmentation with high-level guidance. *arXiv preprint arXiv:1505.02438*, 2015. 1, 2
- [43] C. Wang, C. Xu, C. Wang, and D. Tao. Perceptual adversarial networks for image-to-image transformation. *IEEE Transactions on Image Processing*, 27(8):4066–4079, 2018. 2, 4
- [44] T.-C. Wang, M.-Y. Liu, J.-Y. Zhu, A. Tao, J. Kautz, and B. Catanzaro. High-resolution image synthesis and semantic manipulation with conditional gans. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8798–8807, 2018. 2
- [45] Y. Xue, T. Xu, H. Zhang, L. R. Long, and X. Huang. Segan: Adversarial network with multi-scale l1 loss for medical image segmentation. *Neuroinformatics*, 16(3-4):383–392, 2018. 2, 4
- [46] D. Yang, D. Xu, S. K. Zhou, B. Georgescu, M. Chen, S. Grbic, D. Metaxas, and D. Comaniciu. Automatic liver segmentation using an adversarial image-to-image network. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*, pages 507–515. Springer, 2017. 2
- [47] F. G. Zanjani, D. A. Moin, B. Verheij, F. Claessen, T. Chericci, T. Tan, et al. Deep learning approach to semantic segmentation in 3d point cloud intra-oral scans of teeth. In *International Conference on Medical Imaging with Deep Learning*, pages 557–571, 2019. 2
- [48] 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*, pages 2881–2890, 2017. 5, 6, 8, 11, 12, 13
- [49] S. Zheng, S. Jayasumana, B. Romera-Paredes, V. Vineet, Z. Su, D. Du, C. Huang, and P. H. Torr. Conditional random fields as recurrent neural networks. In *Proceedings of the IEEE international conference on computer vision*, pages 1529–1537, 2015. 1, 2

Supplementary Material

In this section, we provide the hyperparameters for gambling adversarial networks and extra qualitative results for the different experiments.

Hyperparameters

In the following paragraphs, the hyperparameters for the experiments in the results section are described.

U-Net based architecture on Cityscapes. Training details for the experiment on Cityscapes [6] with U-Net based architecture [21]. We trained the segmenter and gambler in alternating fashion of 200 and 400 iterations respectively over 300 epochs with a batch size of 4. The betting maps are calculated with a smoothing factor β of 0.02. Details for the segmenter and gambler are as following:

- **Segmenter:** optimizer: Adam [24], learning rate: 1e-4, beta1: 0.5, beta2: 0.99, adversarial coefficient λ : 1.0, weight decay: 5e-4.
- **Gambler:** optimizer: Adam [24], learning rate 1e-4, beta1: 0.5, beta2: 0.99, weight decay: 5e-4.

PSPnet on Cityscapes. Training details for the experiment on Cityscapes [6] with PSPNet [48]. We trained the segmenter and gambler in alternating fashion of 800 and 800 iterations respectively over 200 epochs with a batch size of 3. The betting maps are calculated with a smoothing factor β of 0.02. Details for the segmenter and gambler are as following:

- **Segmenter:** optimizer: Adam [24], learning rate: 2.5e-5, beta1: 0.5, beta2: 0.99, adversarial coefficient λ : 1.0, weight decay: 5e-4.
- **Gambler:** optimizer: Adam [24], learning rate 2.5e-5, beta1: 0.5, beta2: 0.99, weight decay: 5e-4.

PSPNet on Camvid. Training details for the experiment on Camvid [2] with PSPNet [48]. We trained the segmenter and gambler in alternating fashion of 100 and 200 iterations respectively over 100 epochs with a batch size of 2. The betting maps are calculated with a smoothing factor β of 0.02. Details for the segmenter and gambler are as following:

- **Segmenter:** optimizer: Adam [24], learning rate: 5e-5, beta1: 0.5, beta2: 0.99, adversarial coefficient λ : 0.5, weight decay: 5e-4.
- **Gambler:** optimizer: Adam [24], learning rate 5e-5, beta1: 0.5, beta2: 0.99, weight decay: 5e-4.

Qualitative results

In Figures 6 and 7, extra qualitative results are depicted for the experiments on the Cityscapes [6] validation set for the U-Net based architecture [21] and PSPNet [48] as segmentation network. In Figure 8, some samples are shown on the test set of Camvid [2] with PSPNet as segmentation network.

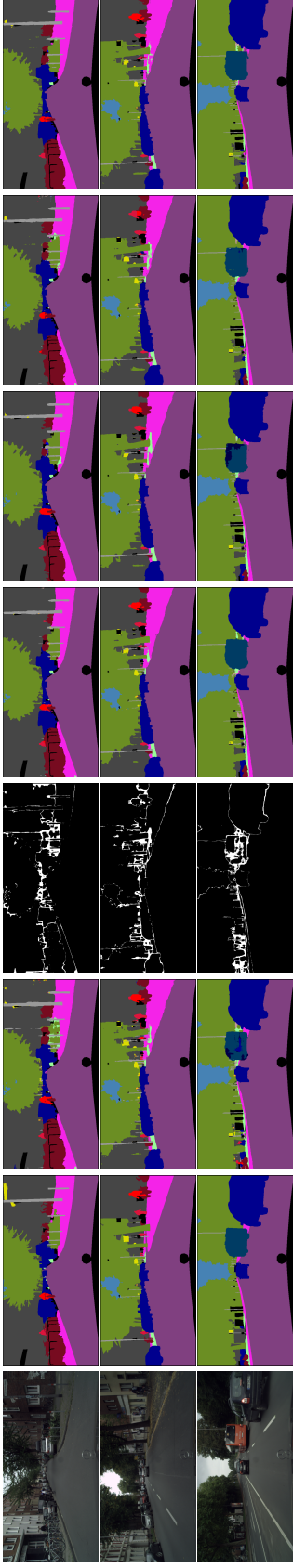


Figure 6: From left to right: RGB-image, ground-truth, CE, betting-map, focal loss, CE + adv, EL-GAN, gambling nets. The betting map is a prediction with as input the RGB image and the CE prediction. Results are for the Cityscapes [6] validation set with the U-Net based architecture [21]. Best visible zoomed-in on a screen.

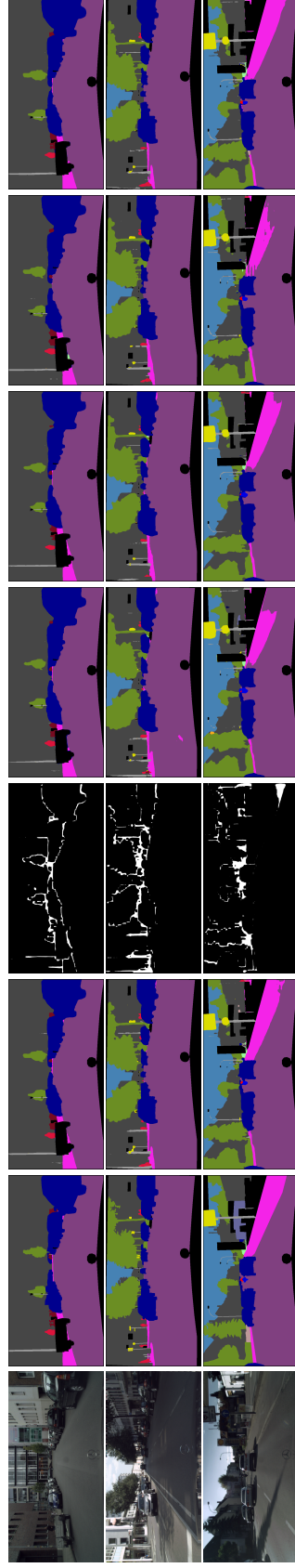


Figure 7: From left to right: RGB-image, ground-truth, CE, betting-map, focal loss, CE + adv, EL-GAN, gambling nets. The betting map is a prediction with as input the RGB image and the CE prediction. Results are for the Cityscapes [6] validation set with PSPNet [48]. Best visible zoomed-in on a screen.

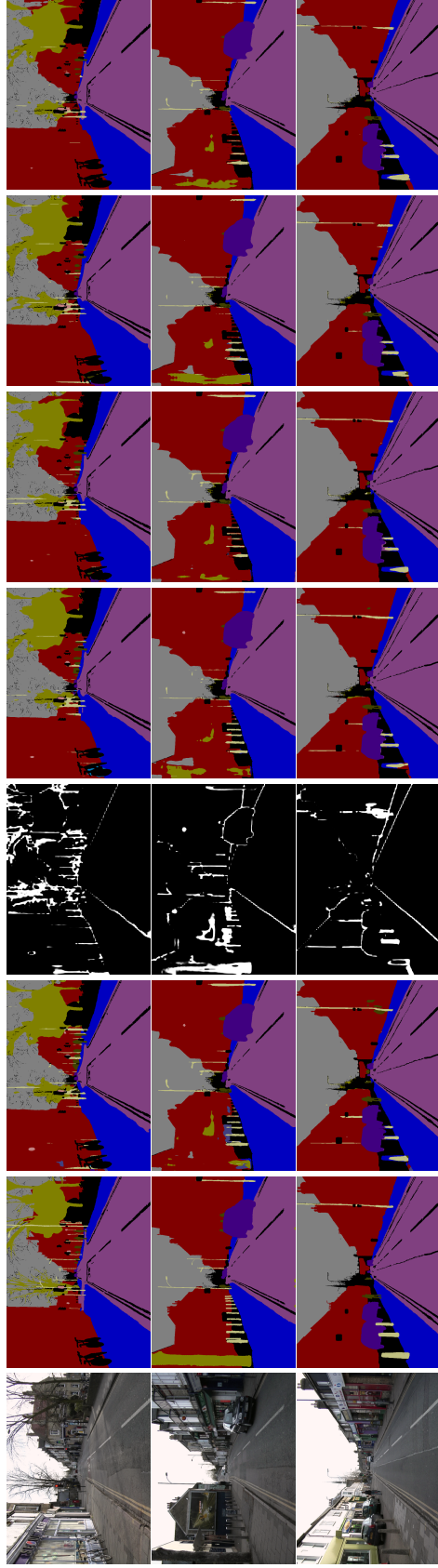


Figure 8: From left to right: RGB-image, ground-truth, CE, betting-map, CE, betting-map, gambling nets, EL-GAN, gambling nets. The betting map is a prediction with as input the RGB image and the CE prediction. Results are for the Camvid [6] test set with PSPNet [48]. Best visible zoomed-in on a screen.