

U-CE: Uncertainty-aware Cross-Entropy for Semantic Segmentation

Steven Landgraf Markus Hillemann Kira Wursthorn Markus Ulrich
 Karlsruhe Institute of Technology (KIT)

(steven.landgraf, markus.hillemann, kira.wursthorn, markus.ulrich)@kit.edu

Abstract

Deep neural networks have shown exceptional performance in various tasks, but their lack of robustness, reliability, and tendency to be overconfident pose challenges for their deployment in safety-critical applications like autonomous driving. In this regard, quantifying the uncertainty inherent to a model’s prediction is a promising endeavour to address these shortcomings. In this work, we present a novel *Uncertainty-aware Cross-Entropy loss (U-CE)* that incorporates dynamic predictive uncertainties into the training process by pixel-wise weighting of the well-known cross-entropy loss (CE). Through extensive experimentation, we demonstrate the superiority of U-CE over regular CE training on two benchmark datasets, Cityscapes and ACDC, using two common backbone architectures, ResNet-18 and ResNet-101. With U-CE, we manage to train models that not only improve their segmentation performance but also provide meaningful uncertainties after training. Consequently, we contribute to the development of more robust and reliable segmentation models, ultimately advancing the state-of-the-art in safety-critical applications and beyond.¹

1. Introduction

Humans often make poor decisions and reach erroneous conclusions while overestimating their abilities, a phenomenon known as the Dunning-Kruger effect [22]. Although deep neural networks are highly effective at solving semantic segmentation problems [34], they also suffer from overconfidence [13]. Additionally, neural networks lack interpretability [12] and struggle to distinguish between in-domain and out-of-domain samples [26]. These flaws are particularly relevant in safety-critical applications, such as autonomous driving [32] and medical imaging [27], as well as in computer vision tasks that have high demands on reliability, like industrial inspection [16, 44] and automation [24, 45], where robust predictions are crucial. Misclassify-

¹Code will be made available once the publication process is complete.

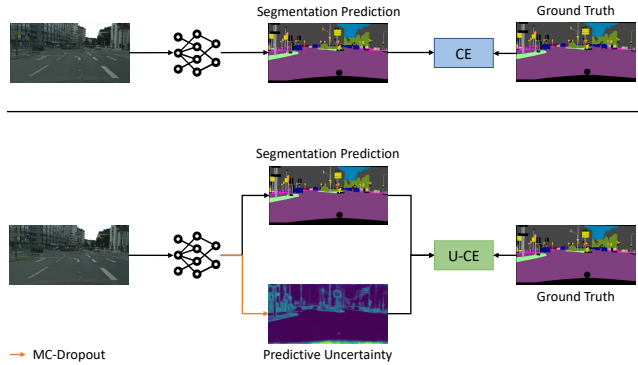


Figure 1. U-CE introduces an uncertainty-aware cross-entropy loss that dynamically incorporates the predictive uncertainties provided by Monte Carlo Dropout (MC-Dropout) into the training process. As a result, we manage to train models that are naturally capable of predicting meaningful uncertainties after training while also improving their segmentation performance.

ing pixels in these contexts can lead to severe consequences, emphasizing the need for robust and trustworthy segmentation models.

Previous work suggests that quantifying the uncertainty inherent to a model’s prediction is a promising endeavour to enhance the safety and reliability of such applications [25–27, 35, 36]. These uncertainties provide additional insights beyond the common softmax probabilities, revealing regions where the model is indecisive and likely to make errors. Surprisingly, the utilization of these uncertainties during the training of segmentation models has not been thoroughly explored.

In this work, we present a novel *Uncertainty-aware Cross-Entropy loss*, referred to as **U-CE**, that addresses this gap by incorporating dynamic uncertainty estimates into the training process as shown in Figure 1. Through pixel-wise uncertainty weighting of the well-known cross-entropy loss (CE), we harness the valuable insights provided by the uncertainties for more effective training. With U-CE, we manage to train models that are naturally capable of predicting meaningful uncertainties after training while simultaneously improving their segmentation performance.

Our contributions can be summarized as follows: Firstly,

we propose the U-CE loss function, which utilizes uncertainty estimates to guide the optimization process, emphasizing regions with high uncertainties. Secondly, we conduct extensive experiments on two benchmark datasets, Cityscapes [8] and ACDC [41], using two common backbones, ResNet-18 and ResNet-101 [15], demonstrating the superiority of U-CE over regular CE training. Lastly, we present additional insights, limitations, and potential improvements for U-CE through multiple ablation studies and a thorough discussion.

2. Related Work

In this section, we briefly review the related work on uncertainty quantification and uncertainty-aware segmentation.

2.1. Uncertainty Quantification

Deep neural networks, with their millions of model parameters and non-linearities, have proven effective in solving complex tasks in natural language processing [38] and computer vision, like semantic segmentation [34]. Unfortunately, due to their complexity, the computation of the exact posterior probability distribution of the network’s output is infeasible [4, 30]. Consequently, approximate uncertainty quantification methods are employed to offer a practical solution to tackle the intractability of the exact posterior distribution. The most prominent methods include Bayesian Neural Networks [31], Monte Carlo Dropout [10], and Deep Ensembles [23]. We will refer to these methods as traditional uncertainty quantification techniques throughout the following.

A mathematically grounded, though computationally complex, approach to uncertainty quantification is provided by Bayesian Neural Networks, which transform a deterministic network into a stochastic one using probabilistic distributions placed over the activations or the weights [18]. For instance, Bayes by Backprob [4] employs variational inference to learn approximate distributions over the weights. These can be used to create an ensemble of models with differently sampled weights to approximate the posterior distribution of the predictions.

Gal and Ghahramani simplify this approximation process by using Monte Carlo Dropout [10]. While dropout is usually applied as a regularization technique [43], Monte Carlo Dropout uses this concept to sample from the posterior distribution of a network’s prediction at test time. In its original form, Monte Carlo Dropout only captures the epistemic uncertainty inherent to the model. To obtain a more comprehensive measure of uncertainty that includes the aleatoric uncertainty, which captures the noise inherent in the observations, Monte Carlo Dropout can be combined with learned uncertainty predictions and assumed density filtering [11, 21, 30].

The current state-of-the-art uncertainty quantification method are Deep Ensembles, which consist of an ensemble of trained models that generate diverse predictions at test time [23]. Due to the introduction of randomness through random weight initialization or different data augmentations across ensemble members [9], Deep Ensembles are well-calibrated [23]. Multiple studies demonstrated that Deep Ensembles generally outperform other uncertainty quantification methods across varying tasks [14, 39, 48]. However, this performance gain is associated with high computational cost.

In addition to the aforementioned approximate uncertainty quantification methods, there has been a growing interest in deterministic single forward-pass approaches, which offer advantages in terms of memory usage and inference time. For example, van Amersfoort *et al.* [46] and Liu *et al.* [29] explore the concept of distance-aware output layers. While these methods demonstrate good performance, they are not competitive with the current state-of-the-art and require significant modifications to the training process [36]. Another approach, proposed by Mukhoti *et al.* [36], simplifies the two previous methods by employing Gaussian Discriminant Analysis for feature-space density estimation after training. Although they perform on par with Deep Ensembles in some settings, their approach still necessitates a more sophisticated training approach. Additionally, fitting the feature-space density estimator is only possible after training, which is not suitable for U-CE where meaningful uncertainties are required during training.

Overall, uncertainty quantification remains an active and evolving field of research, with various approaches offering their own advantages and disadvantages. For our specific case, Monte Carlo Dropout emerges as the preferred option due to its ease of use, minimal impact on the training process, and computational efficiency compared to Deep Ensembles. Through Monte Carlo Dropout sampling, we can compute the predictive uncertainty to apply pixel-wise weighting of the well-known cross-entropy loss. With predictive uncertainties, we refer to the standard deviation of the softmax probabilities of the predicted class provided by Monte Carlo Dropout sampling.

2.2. Uncertainty-aware Segmentation

In the domain of uncertainty-aware segmentation, researchers have explored various techniques to incorporate uncertainty measures into the training process. Surprisingly, traditional uncertainty quantification methods have been largely overlooked or underutilized. We provide an overview of notable works that leverage uncertainty-aware techniques for segmentation tasks in various domains. Additionally, we discuss how U-CE addresses the gap towards full utilization of traditional uncertainty quantification methods during training.

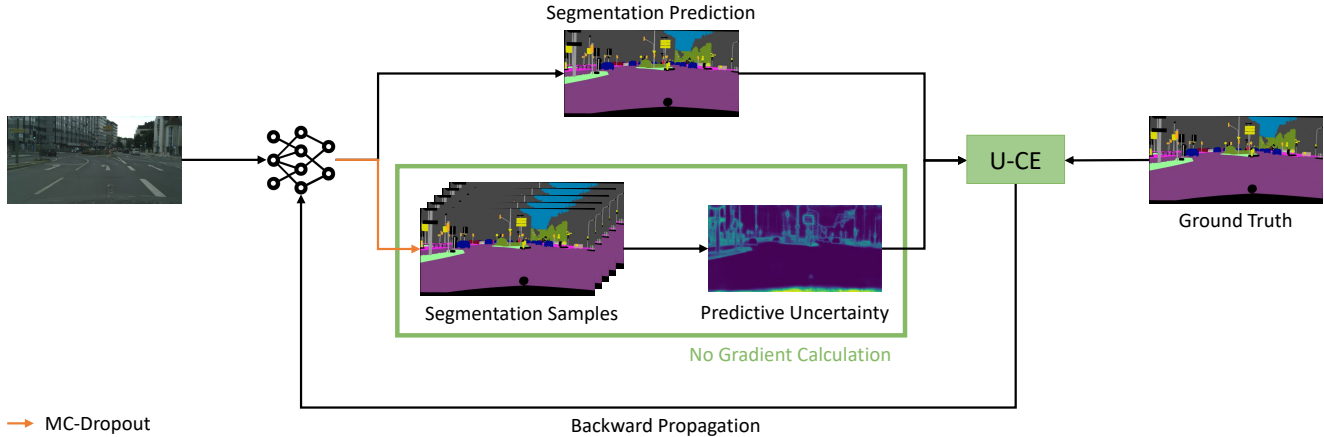


Figure 2. A schematic overview of the training process of U-CE. U-CE integrates the predictive uncertainties of a Monte Carlo Dropout (MC-Dropout) model into the training process to enhance segmentation performance. In comparison to most applications of Monte Carlo Dropout, U-CE utilizes the uncertainties not only at test time but also dynamically during training by applying pixel-wise weighting to the regular cross-entropy loss.

Some of the earlier work on more effective training has originally been designed for object detection. For example, Lin *et al.* [28] introduced the focal loss that down-weights the contribution of easy examples to shift the focus more towards hard examples. Another closely related technique is online hard example mining by Shrivastava *et al.* [42]. They propose to automatically select hard examples to only learn from them and completely ignore the easy examples. By now, both methods have been successfully adapted for semantic segmentation [17, 47].

Another line of work focuses on the identification and compensation of ambiguities and label noise. Kaiser *et al.* [19] propose adding a learned bias to a network’s logits and introducing a novel uncertainty branch to induce the compensation bias only to relevant regions. However, unlike U-CE, their approach does not utilize uncertainties to make training more robust, rather they aim to avoid new noise during data annotation.

More closely related to our work, Bischke *et al.* [3] and Bressan *et al.* [5] propose to leverage uncertainties to improve training on imbalanced aerial image datasets. The former use the per-class uncertainty of the model together with the median frequency to balance training [3]. We argue that dynamically weighting each pixel individually during training, which is what U-CE does, is even more valuable. The latter utilize pixel-wise weights, but only consider the class and labeling uncertainty [5] instead of the predictive uncertainties like U-CE.

In addition to these methods, Chen *et al.* [6] propose to transform the embeddings of the last layer from Euclidean space into Hyperbolic space to dynamically weight pixels based on the hyperbolic distance, which they interpret as uncertainty. Similarly, Bian *et al.* [2] propose an uncertainty estimation and segmentation module to estimate uncertain-

ties that they use to improve the segmentation performance. Unlike U-CE, however, these two works do not incorporate traditional uncertainty quantification methods into training.

In contrast to existing literature, U-CE fully utilizes predictive uncertainties dynamically during training. By pixel-wise uncertainty weighting of the cross-entropy loss, U-CE harnesses valuable insights from the uncertainties to guide the optimization process. This approach enables more effective training, resulting in models that are naturally capable of predicting meaningful uncertainties after training while also improving their segmentation performance.

3. Methodology

In the following, we provide an overview of U-CE, explain our novel uncertainty-aware cross-entropy loss and outline the implementation details.

3.1. Overview

The central idea of U-CE is to incorporate predictive uncertainties into the training process to enhance segmentation performance. As depicted in Figure 2, we propose two simple yet highly effective adaptations to the regular training process:

1. During training, we sample from the posterior distribution with Monte Carlo Dropout to obtain predictive uncertainties alongside the regular segmentation prediction.
2. We apply pixel-wise weighting to the regular cross-entropy loss based on the collected uncertainties.

To compute predictive uncertainties during training, we choose Monte Carlo Dropout. It is straightforward to implement, requires minimal tuning, and is computationally

Algorithm 1 PyTorch-like Pseudocode for a single training step with U-CE

```
x, y = batch
ŷ = model(x)
ŷβ = torch.empty(size = [β, ŷ.size()])
for β and with no_grad do
    ŷβ[β] = model(x)
end for
p = torch.mean(softmax(ŷβ), dim = 0)
q = torch.std(softmax(ŷβ), dim = 0)
σ = q[torch.argmax(p, dim = 1)]
loss = (1 + σ)α · torch.ce(ŷ, y, reduction = "none")
return torch.sum(loss)/torch.numel(loss)
```

more efficient than Deep Ensembles. However, it is worth noting that other uncertainty quantification methods could also be utilized for U-CE. Exploring these alternatives is an interesting avenue for future work, which we will discuss in Section 5.

3.2. Uncertainty-aware Cross-Entropy

Segmentation Sampling. In contrast to typical usage of Monte Carlo Dropout, U-CE incorporates the sampling process from the posterior distribution not only at test time but also during training. To compute the necessary uncertainties for our uncertainty-aware cross-entropy loss, we perform β sampling iterations at each training step. This generates β segmentation samples in addition to the regular segmentation prediction. Notably, gradient computation is disabled during the sampling process as it is unnecessary for backward propagation, which relies solely on the regular segmentation prediction. By disabling gradient computation during sampling, we reduce the additional computational overhead of U-CE in terms of training time and GPU memory usage.

Uncertainty-aware Cross-Entropy Loss. The final objective function of U-CE builds upon the well-known categorical cross-entropy loss and can be defined as:

$$L_{u-ce} = -\frac{1}{N} \sum_{n=1}^N w_n \sum_{c=1}^C y_{n,c} \cdot \log(p_{n,c}), \quad (1)$$

where L_{u-ce} is the uncertainty-aware cross-entropy loss for a single image, N is the number of pixels in the image, C is the number of classes, $y_{n,c}$ is the respective ground truth label, $p_{n,c}$ is the respective predicted softmax probability, and w_n represents the pixel-wise uncertainty weight. It is worth noting that Equation 1 simplifies to the regular cross-entropy loss by setting w_n to one for all pixels.

Pixel-wise Uncertainty Weight. The pixel-wise uncertainty weight w_n can be formulated as:

$$w_n = (1 + \sigma_n)^\alpha, \quad (2)$$

where σ_n denotes the predictive uncertainty, and α controls the influence of the uncertainties in an exponential manner. The predictive uncertainty σ represents the standard deviation of the softmax probabilities of the predicted class of the segmentation samples.

Pseudocode. Finally, Algorithm 1 shows how to use U-CE in the training step in a simplified way. As mentioned earlier, the two key adaptations to the regular training process are sampling with Monte Carlo Dropout during training and applying pixel-wise uncertainty weighting to the regular cross-entropy loss.

4. Experiments

In this section, we conduct an extensive range of experiments to demonstrate the value of incorporating predictive uncertainties into the training process. Firstly, we provide quantitative results comparing regular CE to U-CE under diverse settings. Secondly, we analyze qualitative examples. Lastly, we provide multiple ablation studies.

4.1. Setup

Architecture. For all of our experiments, we employ DeepLabv3+ [7] as the decoder and either a ResNet-18 or ResNet-101 [15] as the encoder. Both backbones are commonly used for semantic segmentation [34, 49], making our work highly comparable and serving as an excellent baseline for future research.

Monte Carlo Dropout. In order to convert our architectures into Monte Carlo Dropout models, we add a dropout layer after each of the four residual block layers of the ResNets, inspired by Kendall *et al.* [20] and Gustafsson *et al.* [14].

Training. For all training processes, we use a Stochastic Gradient Descent (SGD) optimizer [40] with a base learning rate of 0.01, momentum of 0.9, and weight decay of 0.0001. Additionally, we multiply the learning rate of the decoder and segmentation head by ten. Finally, we employ polynomial learning rate scheduling to decay the initial learning rate during the training process, following the formula:

$$lr = lr_{base} \cdot \left(1 - \frac{iteration}{total\ iterations}\right)^{0.9}, \quad (3)$$

where lr is the current learning rate, and lr_{base} is the initial base learning rate. In all training processes, we use a batch size of 16 and train on four NVIDIA A100 GPUs with 40 GB of memory using mixed precision [33].

Datasets. All of our experiments are based on either the Cityscapes dataset [8] or the ACDC dataset [41]. Both datasets are publicly available street scene datasets aimed at

	Encoder	200 Epochs			500 Epochs		
		CE	U-CE $_{\alpha=1}$	U-CE $_{\alpha=10}$	CE	U-CE $_{\alpha=1}$	U-CE $_{\alpha=10}$
Dropout (0%)	RN18	70.0	-	-	72.0	-	-
Dropout (10%)	RN18	69.4	69.6	71.6	72.3	72.3	74.2
Dropout (20%)	RN18	69.0	69.5	71.8	71.9	72.6	73.5
Dropout (30%)	RN18	68.2	69.0	71.0	71.9	72.4	74.1
Dropout (40%)	RN18	66.6	67.7	70.5	71.1	71.1	73.7
Dropout (50%)	RN18	64.3	65.3	69.6	69.0	69.4	72.6
Dropout (0%)	RN101	74.6	-	-	76.1	-	-
Dropout (10%)	RN101	74.8	75.1	76.1	76.3	76.6	77.5
Dropout (20%)	RN101	74.6	74.8	76.6	76.3	77.0	77.7
Dropout (30%)	RN101	74.5	74.7	76.1	76.4	76.6	77.5
Dropout (40%)	RN101	74.7	74.0	75.8	76.1	76.5	78.2
Dropout (50%)	RN101	74.1	73.7	75.9	76.6	76.6	77.3

Table 1. Quantitative comparison between regular CE and U-CE on the Cityscapes dataset [8] for different dropout ratios. The provided numbers represent the mIoU \uparrow . Best respective results are marked in **bold**.

advancing the current state-of-the-art in autonomous driving. The former consists of 2975 training images, 500 validation images, and 1525 test images. The latter contains 1600 training images, 406 validation images, and 2000 test images. Although both datasets share the same 19 evaluation classes and a void class, the ACDC dataset exclusively focuses four adverse conditions: fog, nighttime, rain, and snow.

Data Augmentations. To prevent overfitting, we apply a common data augmentation strategy for all training procedures, regardless of the dataset or architecture used. The strategy includes the following steps:

1. Random scaling with a factor between 0.5 and 2.0.
2. Random cropping with a crop size of 768×768 pixels.
3. Random horizontal flipping with a flip chance of 50%.

Evaluation. Since both test splits are withheld for benchmarking purposes, we utilize the validation images for testing in all our experiments. Unless otherwise specified, we only report single forward pass results based on the original validation images without resizing or sampling for a fair comparison between all of the models. Also, we set the number of segmentation samples β to ten by default.

Metrics. For quantitative evaluations, we primarily report the mean Intersection over Union (mIoU), also known as the Jaccard Index, to measure the segmentation performance. In addition to the mIoU, we also utilize the Expected Calibration Error (ECE) [37] to evaluate the calibration as well as the mean class-wise predictive uncertainty (mUnc) to quantitatively compare the resulting uncertainties.

4.2. Quantitative Evaluation

Tables 1 and 2 outline a quantitative comparison between regular CE and our proposed U-CE loss using two different α values for various dropout ratios and training lengths on

	Encoder	200 Epochs			500 Epochs		
		CE	U-CE $_{\alpha=1}$	U-CE $_{\alpha=10}$	CE	U-CE $_{\alpha=1}$	U-CE $_{\alpha=10}$
Dropout (0%)	RN18	56.3	-	-	62.2	-	-
Dropout (10%)	RN18	55.5	56.4	60.0	62.1	62.8	65.0
Dropout (20%)	RN18	54.6	56.1	60.5	61.5	62.0	65.0
Dropout (30%)	RN18	52.2	54.3	59.2	59.6	61.6	64.3
Dropout (40%)	RN18	48.9	50.8	58.2	56.8	58.8	63.9
Dropout (50%)	RN18	47.7	49.3	56.3	53.3	56.0	62.4
Dropout (0%)	RN101	65.0	-	-	68.8	-	-
Dropout (10%)	RN101	64.5	65.3	67.0	68.4	69.3	69.9
Dropout (20%)	RN101	64.1	65.0	65.8	68.5	68.7	70.2
Dropout (30%)	RN101	62.7	64.3	65.3	68.4	68.5	69.9
Dropout (40%)	RN101	61.1	63.1	65.4	67.8	67.8	70.0
Dropout (50%)	RN101	58.0	60.2	63.7	66.0	67.4	70.2

Table 2. Quantitative comparison between regular CE and U-CE on the ACDC dataset [41] for different dropout ratios. The provided numbers represent the mIoU \uparrow . Best respective results are marked in **bold**.

	Encoder	200 Epochs			500 Epochs		
		mIoU \uparrow	ECE \downarrow	mUnc	mIoU \uparrow	ECE \downarrow	mUnc
CE	RN18	69.0	0.035	0.088	71.9	0.025	0.088
U-CE $_{\alpha=1}$	RN18	69.5	0.036	0.089	72.6	0.027	0.088
U-CE $_{\alpha=10}$	RN18	71.8	0.029	0.085	73.5	0.018	0.084
CE	RN101	74.6	0.026	0.080	76.3	0.041	0.076
U-CE $_{\alpha=1}$	RN101	74.8	0.024	0.079	77.0	0.041	0.076
U-CE $_{\alpha=10}$	RN101	76.6	0.022	0.073	77.7	0.040	0.073

Table 3. A more detailed quantitative comparison between regular CE and U-CE on the Cityscapes dataset [8] using a dropout ratio of 20%. The provided numbers represent the mIoU \uparrow , ECE \downarrow , and mUnc.

the Cityscapes [8] and ACDC [41] datasets. Remarkably, U-CE $_{\alpha=10}$ achieves the highest mIoU across all dropout ratios, even outperforming the baseline models that do not use dropout in most cases. Notably, U-CE $_{\alpha=10}$ achieves a maximum improvement of up to 9.3% over regular CE when training on ACDC [41] for 200 epochs using a ResNet-18 with a dropout ratio of 40%. On average, U-CE $_{\alpha=10}$ outperforms CE by 2.0% on Cityscapes [8] and by 4.6% on ACDC [41]. Interestingly, U-CE $_{\alpha=1}$ also matches or improves upon regular CE training in most cases. On average, U-CE $_{\alpha=1}$ outperforms CE by 0.3% on Cityscapes and by 1.3% on ACDC.

Table 3 provides additional information on the ECE and mUnc for CE and U-CE using a dropout ratio of 20%. In comparison to regular CE and U-CE $_{\alpha=1}$, which exhibit similar results, U-CE $_{\alpha=10}$ not only improves segmentation performance but also yields slightly better calibrated networks, as measured by the ECE. Moreover, the mUnc is also slightly lower for U-CE $_{\alpha=10}$.

Overall, Tables 1, 2 and 3 provide strong evidence for the effectiveness of leveraging predictive uncertainties in the training process.

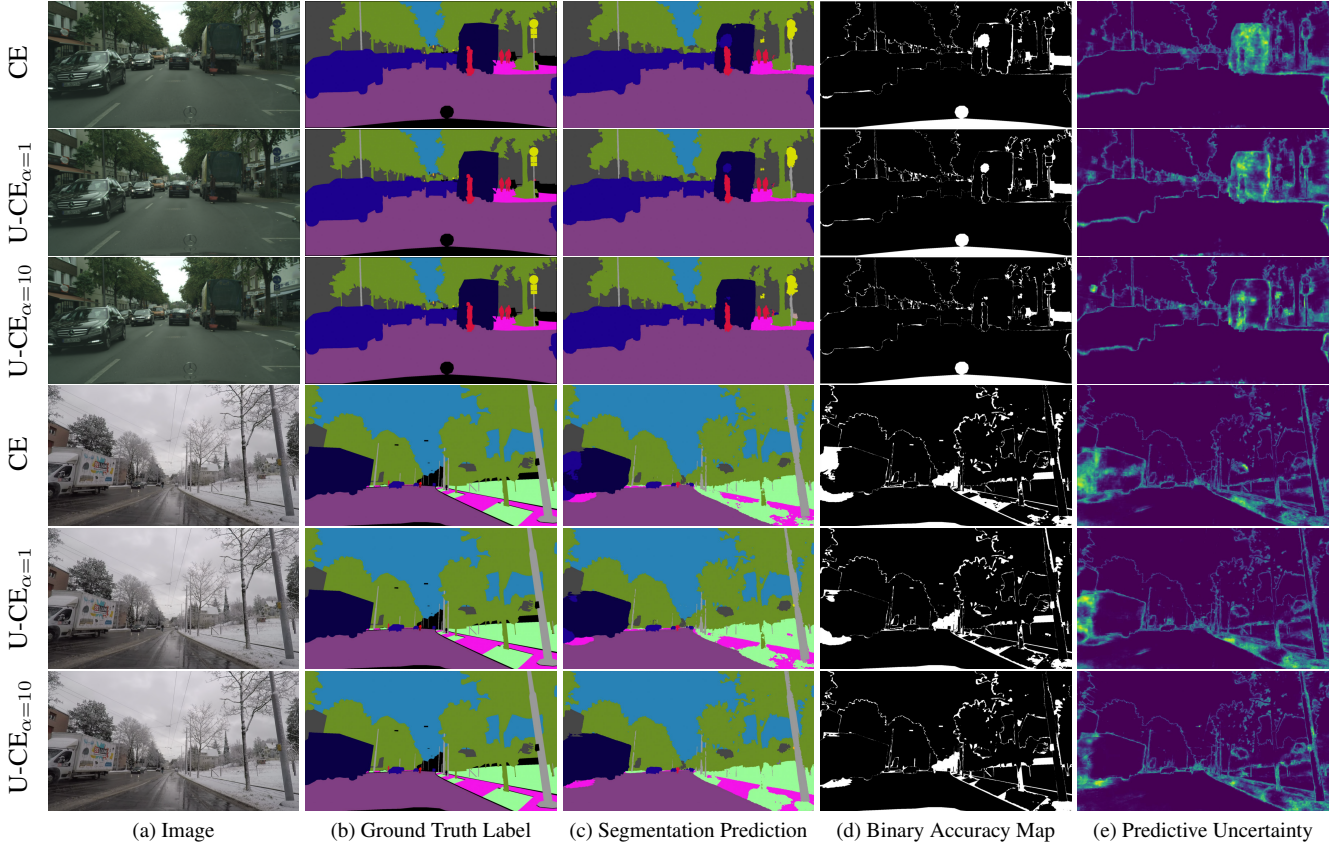


Figure 3. Example images from the Cityscapes and ACDC validation set (a), corresponding ground truth labels (b), the model’s segmentation predictions (c), a binary accuracy map (d), and the predictive uncertainty (e). White pixels in the binary accuracy map are either incorrect predictions or void classes, which appear black in the ground truth label. For the uncertainty prediction, brighter pixels represent higher predictive uncertainties. The first three rows depict results from models with a ResNet-18 backbone and dropout ratio of 20%, trained for 200 epochs on Cityscapes [8]. The last three rows show examples from models using a ResNet-101 backbone and a dropout ratio of 20%, trained for 500 epochs on the ACDC dataset [41].

4.3. Qualitative Evaluation

In addition to the quantitative evaluation, we also provide qualitative examples in Figure 3 showing the original input image, the corresponding ground truth label, the model’s segmentation prediction, a binary accuracy map, and the student’s predictive uncertainty. The first three rows depict results from models with a ResNet-18 backbone and a dropout ratio of 20%, trained for 200 epochs with CE, $U-CE_{\alpha=1}$, $U-CE_{\alpha=10}$ on Cityscapes [8]. The last three rows show examples from models using a ResNet-101 backbone and a dropout ratio of 20%, trained for 500 epochs on the ACDC dataset [41]. The binary accuracy map visualizes incorrectly predicted pixels and void classes in white, and correctly predicted pixels in black.

Generally, for large areas and well-represented classes like road, building, sky, and car, all models perform exceptionally well with minimal errors. Furthermore, there is a strong correlation between the binary accuracy map and the predictive uncertainty, indicating that all models provide

meaningful uncertainties.

Nonetheless, there are nuanced differences between the models. For example, in the first two rows of Figure 3, which represent models trained CE and $U-CE_{\alpha=1}$, there are noticeable misclassifications on top of the human standing in front of the truck. Naturally, this area is also accompanied with high uncertainties. In contrast, the model trained with $U-CE_{\alpha=10}$ exhibits significantly fewer difficulties, resulting in a better segmentation prediction and lower uncertainties.

A similar situation is observable in the last three rows, showing examples from the more challenging ACDC dataset [41]. Here, the model trained with regular CE struggles to correctly segment the truck on the left as well as differentiate between the sidewalk and the terrain on the right side of the image. The model trained with $U-CE_{\alpha=1}$ does slightly better in these areas, but is equally uncertain. Only the model trained with $U-CE_{\alpha=10}$ successfully classifies the truck and differentiates between the sidewalk and the terrain

α	1	2	4	6	8	10	12	14	16
RN18 (Cityscapes)	69.5	70.0	70.7	71.2	71.5	71.8	71.0	47.0	70.9
RN101 (Cityscapes)	74.8	75.2	75.6	76.1	76.4	76.6	76.3	75.8	72.6
RN18 (ACDC)	56.1	56.9	57.6	58.8	58.8	60.5	60.3	60.1	37.5
RN101 (ACDC)	65.0	65.0	65.7	65.5	66.0	65.8	66.7	64.5	19.9

Table 4. Ablation study on the impact of α . The provided numbers represent the mIoU \uparrow . Best respective results are marked in **bold**.

β	0	2	6	10	14	18
CE	69.0 (1:49)	-	-	-	-	-
U-CE $_{\alpha=10}$	-	71.1 (1:52)	71.6 (2:01)	71.6 (2:27)	71.6 (2:53)	71.7 (3:17)

Table 5. Ablation study on the number of segmentation samples β . In addition to the mIoU \uparrow , we provide the training time in hours:minutes \downarrow in paranthesis.

	Random Flipping	Random Scaling	mIoU \uparrow
CE	\times	\times	66.1
	\checkmark	\times	67.0
	\times	\checkmark	68.6
	\checkmark	\checkmark	69.0
U-CE $_{\alpha=1}$	\times	\times	65.8
	\checkmark	\times	67.8
	\times	\checkmark	69.1
	\checkmark	\checkmark	69.5
U-CE $_{\alpha=10}$	\times	\times	69.6
	\checkmark	\times	70.1
	\times	\checkmark	71.8
	\checkmark	\checkmark	71.8

Table 6. Ablation study on the impact of various data augmentations strategies. We use random cropping with a crop size of 768×768 pixels as a baseline for all strategies.

lr_{base}	10^{-1}	10^{-2}	10^{-3}	10^{-4}	10^{-5}
CE	50.5	69.0	55.9	35.6	18.9
U-CE $_{\alpha=1}$	56.0	69.5	57.6	36.9	19.3
U-CE $_{\alpha=10}$	2.0	71.8	65.0	47.6	25.3

Table 7. Ablation study on the base learning rate lr_{base} . The provided numbers represent the mIoU \uparrow . Best results are marked in **bold**.

decently. Consequently, the predictive uncertainty is also lower in these areas.

In summary, the qualitative findings presented in Figure 3 concur with our quantitative evaluation, manifesting the efficacy of U-CE across different datasets and architectures.

4.4. Ablation Studies

In addition to the quantitative and qualitative evaluation, we also present multiple ablation studies. Unless otherwise noted, we confined all of the ablation studies to models that use a ResNet-18 as the backbone, have a dropout ratio of 20%, and were trained for 200 epochs.

Impact of α . The most influential hyperparameter of U-CE is α as it exponentially controls the weighting of the CE

loss. Table 4 demonstrates the impact of different α values on the mIoU for both backbones, ResNet-18 (RN18) and ResNet-101 (RN101), on both Cityscapes and ACDC. Evidently, the segmentation performance consistently improves as α increases until it reaches ten, which stands as the best value in three out of four cases across the two datasets and architectures. Thus, using ten as the default value for α seems to be a fair estimation to achieve the best results, not only for the mentioned cases but potentially for other applications as well. Further increasing α leads to a degradation in mIoU. Additionally, training becomes more unstable as models overly focus on uncertain pixels, resulting in some models failing to converge properly. Nonetheless, U-CE exhibits robustness against changes in α , offering a wide range of valid hyperparameters that lead to improved segmentation results compared to regular CE training.

Impact of β . Table 5 exhibits another ablation study on the number of segmentation samples β . Interestingly, there is no clear benefit of sampling more often than six times, especially with regard to the training time. As indicated by the training times, U-CE $_{\beta=6}$ increases the necessary training time by approximately 10%, whereas U-CE $_{\beta=10}$ extends it by roughly 35%. For comparison, Gal and Ghahramani [10] recommend sampling ten times to get a reasonable estimation of the predictive mean and uncertainty.

Impact of Data Augmentations. The impact of various data augmentation strategies on CE and U-CE is demonstrated in Table 6. The results show that incorporating additional data augmentations on top of the baseline strategy of random cropping with a crop size of 768×768 pixels improves the mIoU across the board. More importantly, this ablation study confirms that U-CE consistently outperforms CE across different data augmentation strategies, indicating its effectiveness in improving segmentation performance.

Impact of lr_{base} . Table 7 shows the ablation study on the base learning rate lr_{base} . The most notable comparison is between regular CE and U-CE $_{\alpha=1}$, which demonstrates that U-CE is not limited to specific learning rates. U-CE $_{\alpha=1}$ consistently outperforms regular CE for all examined base learning rates, despite increasing the training loss by approximately 9% as indicated by the mUnc in Table 3. Moreover, U-CE $_{\alpha=10}$ exceeds the results of CE and U-CE $_{\alpha=1}$ for all base learning rates except 10^{-1} , which caused divergence. Overall, this ablation study confirms the value of leveraging predictive uncertainties during training, irrespective of the learning rate, which is arguably the single most important hyperparameter in deep learning [1].

5. Discussion

In contrast to previous approaches, U-CE fully leverages predictive uncertainties obtained by Monte Carlo Dropout during training. As a result, we manage to train models that not only improve their segmentation performance but are

also naturally capable of predicting meaningful uncertainties after training as well.

While U-CE appears to have no apparent shortcomings, except for a minor increase in training time, we acknowledge the need for a transparent discussion about its potential limitations. Our aim is to effectively guide future work in pushing the boundaries of state-of-the-art techniques, especially in safety-critical applications like autonomous driving.

Limitations. One limitation of U-CE arises in the absence of densely annotated ground truth labels. If most pixels are either labeled as background or designated to be ignored while training, U-CE will likely offer next to no benefit, except for a higher loss around object boundaries. Additionally, U-CE may not contribute to improved segmentation performance if the network is already overfitting the training data. Having said that, the impact of U-CE on generalization needs further examination.

Future Work. With regards to future work, we have multiple suggestions that might be worth investigating. Potentially, the results of U-CE could be further improved if the quality of the uncertainty estimates would be better. Therefore, it would be interesting to integrate Deep Ensembles [23], the state-of-the-art uncertainty quantification method [14, 39, 48], with U-CE, which we could not realize because of computational restraints. On a similar note, it could be worth employing warmup epochs, which we omitted to refrain from introducing another hyperparameter. Additionally, we would like to see α removed from U-CE by incorporating statistical hypothesis testing. This would be beneficial in two ways: Firstly, it would remove the most influential hyperparameter of U-CE. Secondly, and maybe more importantly, it would leverage all of the available uncertainties and not just the predictive uncertainty. Finally, we encourage other researchers to incorporate U-CE into state-of-the-art semantic segmentation approaches and to explore its usefulness in other computer vision tasks that rely on pixel-wise predictions, such as depth estimation.

Overall, we believe that U-CE presents a promising paradigm in semantic segmentation by dynamically leveraging uncertainties to create more robust and reliable models. Despite a minor increase in training time and room for further improvement, we see no reason not to employ U-CE in comparison to regular CE.

6. Conclusion

In this paper, we introduced U-CE, a novel uncertainty-aware cross-entropy loss for semantic segmentation. U-CE incorporates predictive uncertainties, based on Monte Carlo Dropout, into the training process through pixel-wise weighting of the regular cross-entropy loss. As a result, we manage to train models that are naturally capable of predicting meaningful uncertainties after training while si-

multaneously improving their segmentation performance. Through extensive experimentation on the Cityscapes and ACDC datasets using ResNet-18 and ResNet-101 architectures, we demonstrated the superiority of U-CE over regular cross-entropy training.

We hope that U-CE and our thorough discussion of potential limitations and future work contribute to the development of more robust and trustworthy segmentation models, ultimately advancing the state-of-the-art in safety-critical applications and beyond.

Acknowledgment

The authors acknowledge support by the state of Baden-Württemberg through bwHPC.

This work is supported by the Helmholtz Association Initiative and Networking Fund on the HAICORE@KIT partition.

References

- [1] Yoshua Bengio. Practical recommendations for gradient-based training of deep architectures. In *Neural Networks: Tricks of the Trade: Second Edition*, pages 437–478. Springer, 2012.
- [2] Cheng Bian, Chenglang Yuan, Jiexiang Wang, Meng Li, Xin Yang, Shuang Yu, Kai Ma, Jin Yuan, and Yefeng Zheng. Uncertainty-aware domain alignment for anatomical structure segmentation. *Medical Image Analysis*, 64:101732, 2020.
- [3] Benjamin Bischke, Patrick Helber, Damian Borth, and Andreas Dengel. Segmentation of imbalanced classes in satellite imagery using adaptive uncertainty weighted class loss. In *IGARSS 2018-2018 IEEE International Geoscience and Remote Sensing Symposium*, pages 6191–6194. IEEE, 2018.
- [4] Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural network. In Francis Bach and David Blei, editors, *Proceedings of the 32nd International Conference on Machine Learning*, volume 37, pages 1613–1622, Lille, France, 2015. PMLR.
- [5] Patrik Olã Bressan, José Marcato Junior, José Augusto Correia Martins, Maximilian Jaderson de Melo, Diogo Nunes Gonçalves, Daniel Matte Freitas, Ana Paula Marques Ramos, Michelle Taís Garcia Furuya, Lucas Prado Osco, Jonathan de Andrade Silva, et al. Semantic segmentation with labeling uncertainty and class imbalance applied to vegetation mapping. *International Journal of Applied Earth Observation and Geoinformation*, 108:102690, 2022.
- [6] Bike Chen, Wei Peng, Xiaofeng Cao, and Juha Röning. Hyperbolic uncertainty aware semantic segmentation. *arXiv preprint arXiv:2203.08881*, 2022.
- [7] 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)*, September 2018.

- [8] 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)*, June 2016.
- [9] Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep Ensembles: A Loss Landscape Perspective. *arXiv:1912.02757*, 2020.
- [10] Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In Maria Florina Balcan and Kilian Q. Weinberger, editors, *Proceedings of The 33rd International Conference on Machine Learning*, volume 48 of *Proceedings of Machine Learning Research*, pages 1050–1059, New York, New York, USA, 20–22 Jun 2016. PMLR.
- [11] Jochen Gast and Stefan Roth. Lightweight Probabilistic Deep Networks. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 3369–3378, 2018.
- [12] Jakob Gawlikowski, Cedric Rovile Njjeutcheu Tassi, Mohsin Ali, Jongseok Lee, Matthias Humt, Jianxiang Feng, Anna Kruspe, Rudolph Triebel, Peter Jung, Ribana Roscher, Muhammad Shahzad, Wen Yang, Richard Bamler, and Xiao Xiang Zhu. A Survey of Uncertainty in Deep Neural Networks. *arXiv:2107.03342*, 2022.
- [13] Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q. Weinberger. On calibration of modern neural networks. In Doina Precup and Yee Whye Teh, editors, *Proceedings of the 34th International Conference on Machine Learning*, volume 70 of *Proceedings of Machine Learning Research*, pages 1321–1330. PMLR, 06–11 Aug 2017.
- [14] Fredrik K Gustafsson, Martin Danelljan, and Thomas B Schon. Evaluating scalable bayesian deep learning methods for robust computer vision. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*, pages 318–319, 2020.
- [15] 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)*, June 2016.
- [16] Michael Heizmann, Alexander Braun, Markus Glitzner, Matthias Günther, Günther Hasna, Christina Klüver, Jakob Kroß, Erik Marquardt, Michael Overdick, and Markus Ulrich. Implementing machine learning: chances and challenges. *at-Automatisierungstechnik*, 70(1):90–101, 2022.
- [17] Shruti Jadon. A survey of loss functions for semantic segmentation. In *2020 IEEE conference on computational intelligence in bioinformatics and computational biology (CIBCB)*, pages 1–7. IEEE, 2020.
- [18] Laurent Valentin Jospin, Hamid Laga, Farid Boussaid, Wray Buntine, and Mohammed Bennamoun. Hands-on bayesian neural networks—a tutorial for deep learning users. *IEEE Computational Intelligence Magazine*, 17(2):29–48, 2022.
- [19] Timo Kaiser, Christoph Reinders, and Bodo Rosenhahn. Compensation learning in semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3266–3277, 2023.
- [20] Alex Kendall, Vijay Badrinarayanan, and Roberto Cipolla. Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding. *arXiv preprint arXiv:1511.02680*, 2015.
- [21] Alex Kendall and Yarin Gal. What uncertainties do we need in bayesian deep learning for computer vision? In *Proceedings of the 31st International Conference on Neural Information Processing Systems, NIPS’17*, page 5580–5590. Curran Associates Inc., 2017.
- [22] Justin Kruger and David Dunning. Unskilled and unaware of it: how difficulties in recognizing one’s own incompetence lead to inflated self-assessments. *Journal of personality and social psychology*, 77(6):1121, 1999.
- [23] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc., 2017.
- [24] Steven Landgraf, Markus Hillemann, Moritz Aberle, Valentin Jung, and Markus Ulrich. Segmentation of industrial burner flames: A comparative study from traditional image processing to machine and deep learning. *arXiv preprint arXiv:2306.14789*, 2023.
- [25] Steven Landgraf, Kira Wursthorn, Markus Hillemann, and Markus Ulrich. Dudes: Deep uncertainty distillation using ensembles for semantic segmentation. *arXiv preprint arXiv:2303.09843*, 2023.
- [26] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training Confidence-calibrated Classifiers for Detecting Out-of-Distribution Samples. *arXiv:1711.09325*, 2018.
- [27] Christian Lebig, Vaneeda Allken, Murat Seçkin Ayhan, Philipp Berens, and Siegfried Wahl. Leveraging uncertainty information from deep neural networks for disease detection. *Scientific Reports*, 7(1):17816, 2017.
- [28] Tsung-Yi Lin, Priya Goyal, Ross Girshick, Kaiming He, and Piotr Dollár. Focal loss for dense object detection. In *Proceedings of the IEEE international conference on computer vision*, pages 2980–2988, 2017.
- [29] Jeremiah Liu, Zi Lin, Shreyas Padhy, Dustin Tran, Tania Bedrax Weiss, and Balaji Lakshminarayanan. Simple and principled uncertainty estimation with deterministic deep learning via distance awareness. *Advances in Neural Information Processing Systems*, 33:7498–7512, 2020.
- [30] Antonio Loquercio, Mattia Segu, and Davide Scaramuzza. A General Framework for Uncertainty Estimation in Deep Learning. *IEEE Robotics and Automation Letters*, 5(2):3153–3160, 2020.
- [31] David J. C. MacKay. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3):448–472, 1992.
- [32] Rowan McAllister, Yarin Gal, Alex Kendall, Mark van der Wilk, Amar Shah, Roberto Cipolla, and Adrian Weller. Concrete Problems for Autonomous Vehicle Safety: Advantages of Bayesian Deep Learning. In *Proceedings of the Twenty-Sixth International Joint Conference on Artificial Intelligence*, pages 4745–4753, Melbourne, Australia, 2017.

- International Joint Conferences on Artificial Intelligence Organization.
- [33] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, et al. Mixed precision training. *arXiv preprint arXiv:1710.03740*, 2017.
- [34] Shervin Minaee, Yuri Boykov, Fatih Porikli, Antonio Plaza, Nasser Kehtarnavaz, and Demetri Terzopoulos. Image segmentation using deep learning: A survey. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 44(7):3523–3542, 2022.
- [35] Jishnu Mukhoti and Yarin Gal. Evaluating bayesian deep learning methods for semantic segmentation. *arXiv preprint arXiv:1811.12709*, 2018.
- [36] Jishnu Mukhoti, Andreas Kirsch, Joost van Amersfoort, Philip HS Torr, and Yarin Gal. Deep deterministic uncertainty: A new simple baseline. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 24384–24394, 2023.
- [37] Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015.
- [38] Daniel W Otter, Julian R Medina, and Jugal K Kalita. A survey of the usages of deep learning for natural language processing. *IEEE transactions on neural networks and learning systems*, 32(2):604–624, 2020.
- [39] Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D. Sculley, Sebastian Nowozin, Joshua Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model's uncertainty? evaluating predictive uncertainty under dataset shift. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett, editors, *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- [40] Herbert Robbins and Sutton Monro. A Stochastic Approximation Method. *The Annals of Mathematical Statistics*, 22(3):400–407, 1951.
- [41] Christos Sakaridis, Dengxin Dai, and Luc Van Gool. ACDC: The adverse conditions dataset with correspondences for semantic driving scene understanding. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2021.
- [42] Abhinav Shrivastava, Abhinav Gupta, and Ross Girshick. Training region-based object detectors with online hard example mining. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 761–769, 2016.
- [43] Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15(56):1929–1958, 2014.
- [44] Carsten Steger, Markus Ulrich, and Christian Wiedemann. *Machine Vision Algorithms and Applications*. John Wiley & Sons, 2018.
- [45] Markus Ulrich and Markus Hillemann. Generic hand–eye calibration of uncertain robots. In *2021 IEEE International Conference on Robotics and Automation (ICRA)*, pages 11060–11066, 2021.
- [46] Joost Van Amersfoort, Lewis Smith, Yee Whye Teh, and Yarin Gal. Uncertainty estimation using a single deep deterministic neural network. In *International conference on machine learning*, pages 9690–9700. PMLR, 2020.
- [47] Yuchao Wang, Haochen Wang, Yujun Shen, Jingjing Fei, Wei Li, Guoqiang Jin, Liwei Wu, Rui Zhao, and Xinyi Le. Semi-supervised semantic segmentation using unreliable pseudo-labels. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 4248–4257, June 2022.
- [48] Kira Wursthorn, Markus Hillemann, and Markus Ulrich. Comparison of uncertainty quantification methods for CNN-based regression. *The International Archives of the Photogrammetry, Remote Sensing and Spatial Information Sciences*, XLIII-B2-2022:721–728, 2022.
- [49] Rongyu Zhang, Lixuan Du, Qi Xiao, and Jiaming Liu. Comparison of backbones for semantic segmentation network. In *Journal of Physics: Conference Series*, volume 1544, page 012196. IOP Publishing, 2020.