

Efficient Multi-task Uncertainties for Joint Semantic Segmentation and Monocular Depth Estimation

Steven Landgraf Markus Hillemann Theodor Kapler Markus Ulrich
 Institute of Photogrammetry and Remote Sensing (IPF)
 Karlsruhe Institute of Technology

(steven.landgraf, markus.hillemann, markus.ulrich)@kit.edu
 theodor.kapler@student.kit.edu

Abstract

Quantifying the predictive uncertainty emerged as a possible solution to common challenges like overconfidence or lack of explainability and robustness of deep neural networks, albeit one that is often computationally expensive. Many real-world applications are multi-modal in nature and hence benefit from multi-task learning. In autonomous driving, for example, the joint solution of semantic segmentation and monocular depth estimation has proven to be valuable. In this work, we first combine different uncertainty quantification methods with joint semantic segmentation and monocular depth estimation and evaluate how they perform in comparison to each other. Additionally, we reveal the benefits of multi-task learning with regard to the uncertainty quality compared to solving both tasks separately. Based on these insights, we introduce EMUFormer, a novel student-teacher distillation approach for joint semantic segmentation and monocular depth estimation as well as efficient multi-task uncertainty quantification. By implicitly leveraging the predictive uncertainties of the teacher, EMUFormer achieves new state-of-the-art results on Cityscapes and NYUv2 and additionally estimates high-quality predictive uncertainties for both tasks that are comparable or superior to a Deep Ensemble despite being an order of magnitude more efficient.

1. Introduction

Because of their unparalleled performance in fundamental perception tasks like semantic segmentation [46] or monocular depth estimation [9], deep neural networks are increasingly being deployed in real-time and safety-critical applications like autonomous driving [44], industrial inspection [21, 62], and automation [31]. However, they often suffer from overconfidence [17], lack explainability [16], and struggle to distinguish between in-domain and out-of-

domain samples [34], which is of paramount importance for applications where prediction reliability is crucial. Since incorrect predictions can lead to severe consequences, previous work suggests that quantifying the uncertainty inherent to a model’s prediction is a promising endeavour to make such applications safer [32–35, 40, 49, 50]. In autonomous driving, for instance, the car could provide feedback to the driver when it is uncertain or preemptively make risk-averse predictions based on the uncertainty.

In recent years, a number of promising uncertainty quantification methods have been proposed to make deep neural networks more robust [1, 12, 30, 37, 42, 50, 63, 64]. Unfortunately, these methods either introduce technical complexity or require computationally expensive sampling from a stochastic process to estimate the uncertainty of a prediction. Additionally, they do not consider that real-world applications, like robotics [52] or autonomous driving [5], are multi-modal in nature and benefit from multi-task learning, especially within the context of semantic segmentation and monocular depth estimation [5, 52]. Although there have been successful attempts at making uncertainty quantification methods more efficient through the concept of knowledge distillation [2, 23, 33, 57, 60], they have thereby either focused on semantic segmentation [2, 23, 33, 57] or monocular depth estimation [57, 60]. This represents a notable research gap in the current literature.

In this work, we conduct a comprehensive series of experiments to study multi-task uncertainties and propose a novel student-teacher distillation approach for joint semantic segmentation and monocular depth estimation as well as efficient multi-task uncertainty quantification. Our contributions can be summarized as follows:

- We propose a novel student-teacher distillation approach for **Efficient Multi-task Uncertainties** for joint semantic segmentation and monocular depth estimation with a modern Vision-Transformer, which we call **EMUFormer**.

	Seg.	Pred. Unc.	Depth	Pred. Unc.	Parameters	FLOPs	FPS
a) SegFormer-B2 [69]	✓	×	×	×	27.3M	72.6G	55.3
b) DepthFormer-B2	×	✓	×	×	27.3M	72.1G	57.1
c) SegDepthFormer-B2	✓	×	✓	×	30.5M	120.1G	44.8
DE of a)	✓	✓	×	×	273.6M	726.4G	5.6
DE of b)	×	×	✓	✓	273.5M	720.8G	7.2
DE of c)	✓	✓	✓	✓	305.1M	1201.1G	4.9
EMUFormer-B2 (Ours)	✓	✓	✓	✓	30.5M	120.1G	44.8

Table 1. Overview of the segmentation (Seg.), depth estimation (Depth) and uncertainty quantification (Pred. Unc.) capabilities as well as the respective number of parameters, FLOPs and FPS for different single-task and multi-task models and their respective Deep Ensemble (DE) versions with 10 members. SegFormer [69] and DepthFormer represent single-task models, whereas SegDepthFormer and EMUFormer depict multi-task models. B2 represents the medium-sized encoder of SegFormer, which was used for all models. Results are based on single-scale inference conducted on the NYUv2 [59] dataset using an NVIDIA A100 GPU.

- We show that by implicitly leveraging the predictive uncertainties during training, EMUFormer can achieve new state-of-the-art results on Cityscapes and NYUv2.
- We combine different uncertainty quantification methods with joint semantic segmentation and monocular depth estimation and evaluate how they perform in comparison to each other.
- We reveal the benefits of multi-task learning with regard to the uncertainty quality compared to solving semantic segmentation and monocular depth estimation separately.

As Table 1 demonstrates, EMUFormer estimates high-quality predictive uncertainties for both tasks that are comparable to the Deep Ensemble teacher despite being an order of magnitude more efficient.

2. Related Work

In this section, we summarize the related work on joint semantic segmentation and monocular depth estimation, uncertainty quantification, and knowledge distillation.

2.1. Joint Semantic Segmentation and Monocular Depth Estimation

Semantic segmentation and monocular depth estimation are both fundamental problems in image understanding that involve pixel-wise predictions based on a single input image. Motivated by the strong correlation and complementary properties of the two tasks, multiple previous works have focused on solving both tasks in a joint manner [3, 4, 14, 20, 26, 27, 29, 36, 38, 39, 48, 52, 67, 70, 71]. To limit the scope of this literature review, we refrain from covering other multi-task approaches with joint representation

sharing [72] or methods that leverage the depth map to improve the semantic segmentation prediction [24, 66].

In their pioneering work, Wang et al. [67] propose a unified framework for semantic segmentation and monocular depth prediction through joint training and applying a two-layer hierarchical conditional random field to enforce synergy between global and local predictions. Similarly, Liu et al. [38] use a conditional random field that fuses the feature maps from both tasks. In contrast, Mousavian et al. [48] train parts of the model for each task separately and then fine-tune the full model on both tasks with a single loss function. On a similar note, Xu et al. [70] propose a multi-task prediction-and-distillation network, which first predicts a set of intermediate auxiliary tasks. These intermediate outputs are then utilized as multi-modal input for the final task - a concept also followed by Vandenhende et al. [65]. The idea of knowledge distillation is also used by Nekrasov et al. [52], primarily focusing on real-time estimation without specifically delving into uncertainty quantification. Jiao et al. [27] introduce an attention-driven loss that does not treat all pixels in an image equally to mutually improve semantic segmentation and monocular depth estimation. In a similar way, Bruggemann et al. [4] and Liu et al. [39] build on the idea of introducing attention mechanisms into the architecture to improve results. Comparably, Gao et al. [14] propose a shared attention block with contextual supervision next to a feature-sharing module and a consistency loss. In a follow-up work, they extend their approach by incorporating confidences into their losses to improve the performance [15]. Similarly, Kendall et al. [29] utilize the homoscedastic uncertainty, which they define as a task-dependent uncertainty that captures the relative confidence between tasks, to weight the individual losses. Finally, there are multiple works [20, 26, 36] that propose specialized architectures, where they either improve the feature extraction by separating the relevant features for one task from the features which are relevant for both tasks [36] or exploit geometric constraints by integrating the information of the objectness [20] or apply a randomly-weighted training strategy to balance the losses and gradients impartially and dynamically [26].

Remarkably, most of the discussed approaches use out-of-date architectures and require complex adaptations to either the model, the training process, or both. In order to push the state-of-the-art forward, we adapt a modern Vision-Transformer-based architecture similar to Xu et al. [71]. In order to maintain methodological simplicity and transparency of the results, we refrain from introducing cross-task attention mechanisms, contrastive self-supervised learning algorithms, and the loss weighting strategy of [29], and nevertheless achieve superior results. However, these strategies could also be applied to our method, potentially further improving the results.

2.2. Uncertainty Quantification

A large variety of uncertainty quantification methods [1, 12, 30, 37, 42, 50, 63, 64] have been developed to compensate for the above-mentioned shortcomings of deep neural networks. The predictive uncertainty can be decomposed into aleatoric and epistemic uncertainty [11]. Aleatoric uncertainty captures the irreducible data uncertainty, which, for example, can be introduced by image noise or noisy labels as a result of imprecise measurements. Epistemic uncertainty accounts for the model uncertainty, which can be reduced by using more or better training data [11, 28]. Disentangling these two uncertainty components can be essential for applications such as active learning [13] or the detection of out-of-distribution samples [56]. For instance, active learning benefits from avoiding inputs with high aleatoric uncertainty unless they exhibit high epistemic uncertainty, which is vital for model improvement [13, 28].

Most well-known uncertainty quantification methods require multiple forward passes at test time, making them computationally expensive. For instance, Gal and Ghahramani [12] propose Monte Carlo Dropout (MCD) as an approximation of a stochastic Gaussian process. While dropout is usually only used for regularization during training [61], MCD applies this technique during test time to sample from the posterior distribution of the predictions at test time. Although MCD is easy to implement and thus very popular, Deep Ensembles [30] are commonly regarded as the state-of-the-art approach for uncertainty quantification across varying tasks [19, 54, 68]. They consist of an ensemble of trained models that generate diverse predictions due to the introduction of randomness through random weight initialization or different data augmentations during training [10].

Multiple forward passes at test time render the aforementioned methods impracticable or even unusable for real-time applications because of their high computational cost. Consequently, there has been an increased interest in deterministic single forward-pass methods that demand less overhead. For example, Van Amersfoort et al. [64] and Liu et al. [37] consider distance-aware output layers for quantifying the predictive uncertainty. Even though these methods provide a computationally more efficient approach, they are not competitive with the current state-of-the-art and require significant modifications to the training process [50]. By using Gaussian Discriminant Analysis post-training for feature-space density estimation, Mukhoti et al. [50] simplify the aforementioned approaches. Although they manage to perform on par with a Deep Ensemble in some settings, their method requires performing Gaussian Discriminant Analysis after training, which adds complexity. In contrast, Valdenegro-Toro [63] proposes a simple, yet effective approximation to Deep Ensembles, where the ensemble covers only a subset of layers instead of the whole

model. These so-called Deep Sub-Ensembles (DSE) enable a trade-off between uncertainty quality and computational cost [63].

To the best of our knowledge, quantifying predictive uncertainties in joint semantic segmentation and monocular depth estimation has not been explored yet. To this end, we compare multiple uncertainty quantification methods for this task and investigate how multi-task learning influences the quality of uncertainty estimates in comparison to solving both tasks separately.

2.3. Knowledge Distillation

Knowledge distillation, introduced by Hinton et al. [22], involves transferring the knowledge from a complex model (teacher) to a typically smaller model (student), aiming to enhance the student’s performance on a given task by imitating the predictions of the teacher [22] or transferring knowledge from intermediate features [55]. More recent work has adapted the concept of knowledge distillation to enable real-time uncertainty quantification. While some previous work employs MCD to estimate uncertainties for the student to learn [2, 18, 57], the majority proposes to use a Deep Ensemble [7, 23, 33, 43, 60]. Among these, Deng et al. [7] are the only ones to consider a multi-task problem by looking at emotion recognition.

To enable real-time uncertainty quantification in joint semantic segmentation and monocular depth estimation, we propose EMUFormer, a novel student-teacher distillation approach that aims to preserve both prediction and uncertainty quality without introducing a speed-penalty during inference.

3. Methodology

In the following, we provide an overview of the methodology of this paper, describe the baseline models that we use to analyse the uncertainties of joint semantic segmentation and monocular depth estimation. We will also explain our student-teacher distillation approach for efficient multi-task uncertainties.

3.1. Overview

This paper can broadly be categorized into two parts: First, we evaluate how multi-task learning influences the uncertainty quality. Second, we propose EMUFormer, a novel student-teacher distillation approach for efficient multi-task uncertainties.

Multi-task Uncertainty Evaluation. Drawing from the related work on uncertainty quantification (Section 2.2), we evaluate Deep Ensembles (DEs) [30], Monte Carlo Dropout (MCD) [12], and Deep Sub-Ensembles (DSEs) [63]. The choice is motivated by their simplicity, ease of implementation, parallelizability, minimal tuning requirements,

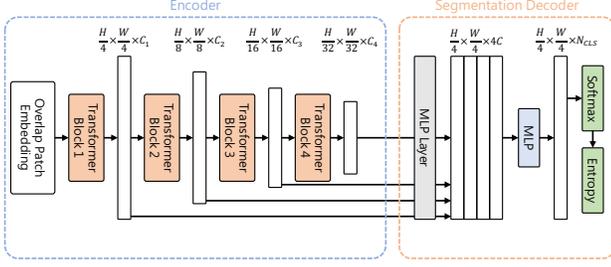


Figure 1. A schematic overview of the SegFormer [69] architecture. The model consists of two main modules: A hierarchical Transformer-based encoder that generates high-resolution coarse features and low-resolution fine features and a lightweight all-MLP segmentation decoder.

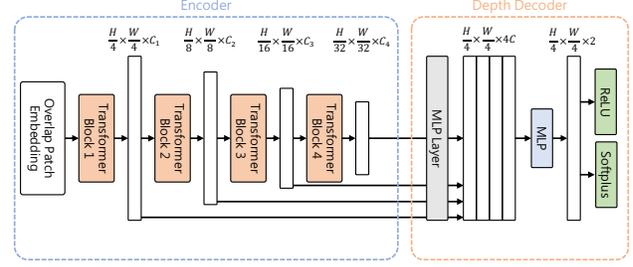


Figure 2. A schematic overview of our DepthFormer architecture. Being derived from SegFormer [69], it consists of two main modules: A hierarchical Transformer-based encoder that generates high-resolution coarse features and low-resolution fine features and a lightweight all-MLP depth decoder.

and representation of the current state-of-the-art in uncertainty quantification. Moreover, applying these approaches to both semantic segmentation and monocular depth estimation is straightforward, which is not the case for the other aforementioned uncertainty quantification approaches [1, 37, 50, 64].

To explore the impact of multi-task learning on uncertainty quality, we conduct all of the evaluations using three models:

1. **SegFormer** [69]: An efficient semantic segmentation Vision Transformer.
2. **DepthFormer**: An efficient monocular depth estimation model Vision Transformer.
3. **SegDepthFormer**: A joint model addressing both semantic segmentation and monocular depth estimation.

We derive the latter two, DepthFormer and SegDepthFormer, from the SegFormer [69] architecture. Key modifications will be explained in Section 3.2.2 and 3.2.3 respectively.

EMUFormer. In order to achieve efficient multi-task uncertainties without sacrificing neither prediction performance nor uncertainty quality, we propose EMUFormer. EMUFormer applies student-teacher distillation as a two-step framework: First, we train an adequate teacher with ground truth labels that is able to quantify high-quality uncertainties. Subsequently, we train a student with the same ground truth labels while distilling the teacher’s uncertainties.

3.2. Baseline Models

Hereinafter, we go over the three baseline models, SegFormer [69], DepthFormer, and SegDepthFormer. For all of the three models, we will shortly describe their architecture, illustrate the training criterion, and how we obtain a measurement for the uncertainty. While these models are capable of estimating the aleatoric uncertainty [28, 30], they

are not able to quantify the more complete predictive uncertainty, which includes the epistemic uncertainty. For this, one of the aforementioned uncertainty quantification methods has to be used.

3.2.1 SegFormer

Architecture. For the semantic segmentation task, we use SegFormer [69], a modern Transformer-based architecture that stands out because of its high efficiency and performance. Thus, it is particularly suitable for real-time uncertainty quantification. As depicted in Figure 1, SegFormer consists of two main modules: A hierarchical Transformer-based encoder that generates high-resolution coarse features and low-resolution fine features and a lightweight all-MLP segmentation decoder. The latter fuses the multi-level features of the encoder to produce a final segmentation prediction with the softmax activation function, which can be formulated as:

$$p(z) = \frac{e^{z_i}}{\sum_{k=1}^K e^{z_k}}, \quad (1)$$

where $p(z)$ are the class probabilities of the softmax function that exponentiates each of the K elements of the input vector x , often referred to as logits, and then normalizes the results to obtain a probability distribution. Since SegFormer [69] only outputs logits at a $\frac{H}{4} \times \frac{W}{4}$ resolution given an input image of size $H \times W$, we use bilinear interpolation [69] before applying the softmax function on z to obtain the original resolution for the final segmentation prediction.

Training Criterion. For the objective function during training, we use the well-known categorical Cross-Entropy loss

$$\mathcal{L}_{\text{CE}} = -\frac{1}{N} \sum_{n=1}^N \sum_{c=1}^C y_{n,c} \cdot \log(p(z)_{n,c}), \quad (2)$$

where \mathcal{L}_{CE} is the Cross-Entropy loss for a single image, N is the number of pixels in the image, C is the number of

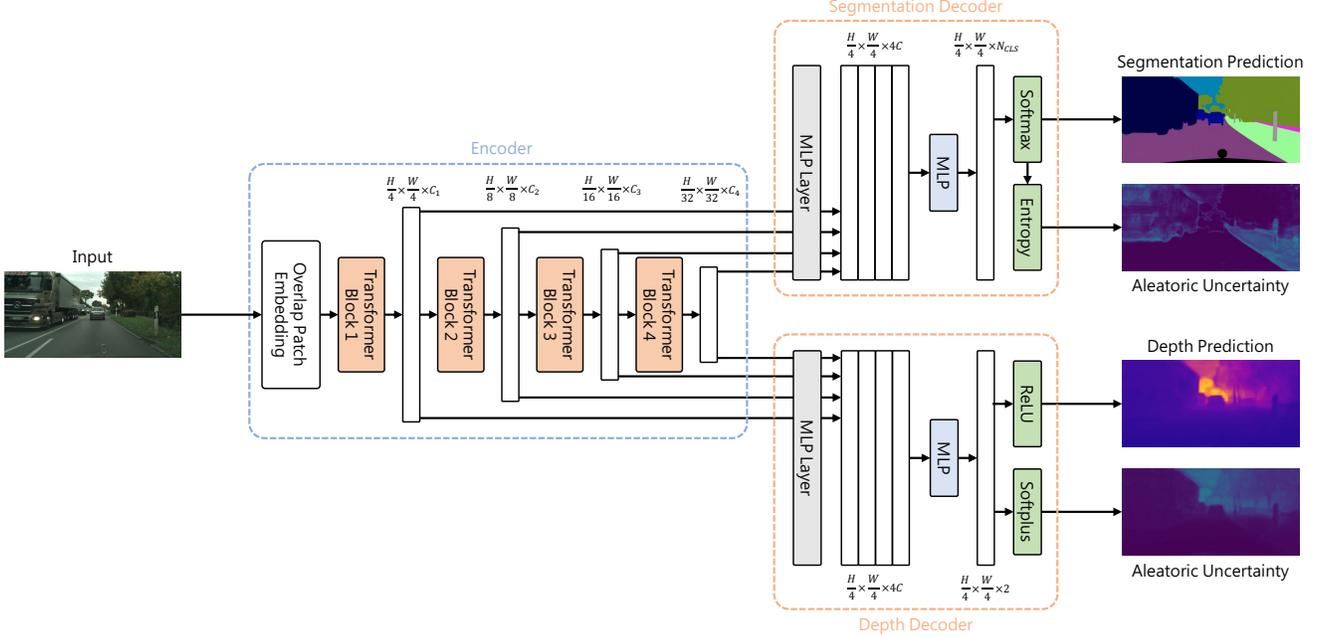


Figure 3. A schematic overview of the SegDepthFormer architecture. The model combines the SegFormer [69] architecture with a lightweight all-MLP depth decoder.

classes, $y_{n,c}$ is the corresponding ground truth label, and $p(z)_{n,c}$ is the predicted softmax probability.

Aleatoric Uncertainty. We compute the predictive entropy

$$H(p(z)) = - \sum_{c=1}^C p(z)_c \cdot \log(p(z)_c) , \quad (3)$$

which serves as the aleatoric uncertainty [28].

3.2.2 DepthFormer

Architecture. Inspired by the efficiency and performance of SegFormer [69], we propose DepthFormer for monocular depth estimation. As Figure 2 shows, we use the same hierarchical Transformer-based encoder as SegFormer to generate high-level and low-level features. Similarly, those multi-level features are fused in an all-MLP decoder. In contrast to the segmentation decoder, the depth decoder differs by having two output channels: one for the predictive mean $\mu(z)$ and one for the predictive variance $s^2(z)$ [40].

Predictive Mean. The first output channel uses a Rectified Linear Unit (ReLU) output activation function

$$\mu(z) = \max(0, z) , \quad (4)$$

which serves as the predictive mean for monocular depth estimation.

Predictive Variance. The second output channel applies a Softplus activation

$$s^2(z) = \log(1 + e^z) , \quad (5)$$

which is a smooth approximation of the ReLU function with the advantage of being differentiable, also at $z = 0$. Empirically, we found Softplus to work better than ReLU for the predictive variance, following the work by Lakshminarayanan et al. [30].

Training Criterion. For regression tasks, neural networks typically output only a predictive mean $\mu(z)$ and the parameters are, in the most straightforward approach, optimized by minimizing the mean squared error (MSE). However, the MSE does not cover uncertainty. Therefore, we follow the approach of Nix and Weigend [53] instead: By treating the neural networks prediction as a sample from a Gaussian distribution with the predictive mean $\mu(z)$ and corresponding predictive variance $s^2(z)$, we can minimize the Gaussian Negative Log-Likelihood (GNLL) loss, which can be formulated as:

$$\mathcal{L}_{\text{GNLL}} = \frac{1}{2} \left(\frac{(y - \mu(z))^2}{s^2(z)} + \log(s^2(z)) \right) , \quad (6)$$

where y is the the ground truth depth.

Aleatoric Uncertainty. Through GNLL minimization, DepthFormer does not only optimize the predictive means, but also inherently learns the corresponding variances, which can be interpreted as the aleatoric uncertainty [28, 40].

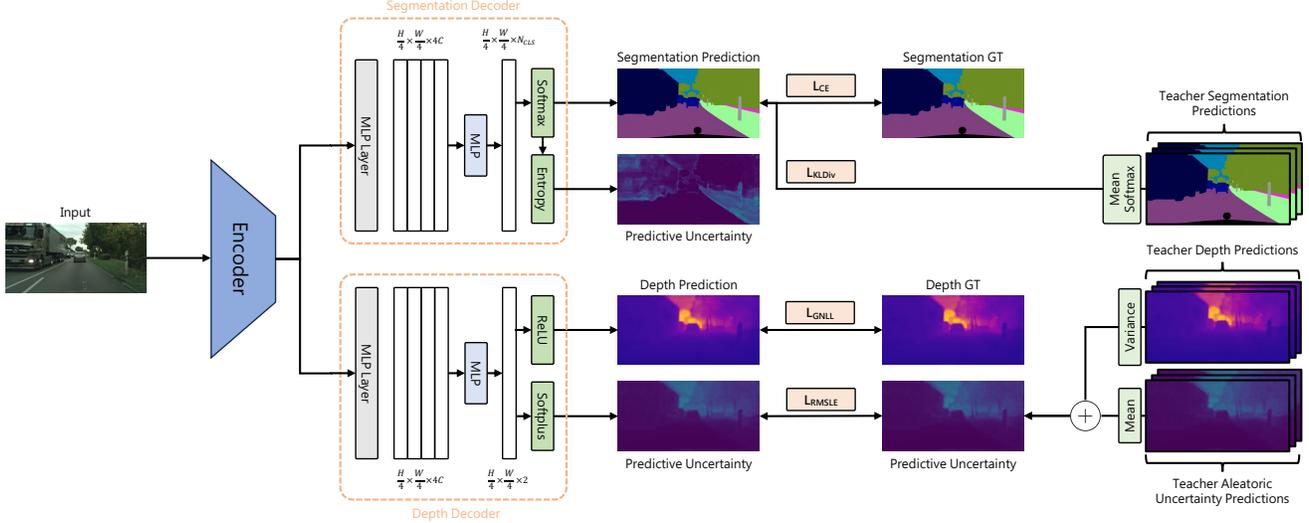


Figure 4. A schematic overview of EMUFormer. In comparison to our proposed SegDepthFormer, EMUFormer utilizes two additional losses that distill the predictive uncertainties of the teacher into the student model.

3.2.3 SegDepthFormer.

Architecture. In order to solve semantic segmentation and monocular depth estimation in a joint manner, we propose SegDepthFormer. The architecture, which is shown in Figure 3, comprises three modules: a hierarchical Transformer-based encoder, an all-MLP segmentation decoder, and an all-MLP depth decoder. The encoder and segmentation decoder are adapted from SegFormer [69] (Section 3.2.1), while the depth decoder is from DepthFormer (Section 3.2.2). Both decoders fuse the multi-level features obtained through the shared encoder to predict a final segmentation mask and a pixel-wise depth estimation, respectively.

Training Criterion. SegDepthFormer is trained to minimize the weighted sum of the two previously described objective functions:

$$\mathcal{L} = \mathcal{L}_{CE} + w_1 \mathcal{L}_{GNLL}, \quad (7)$$

where w_1 is a simple weighting factor. Because both loss values are of similar magnitude, we set $w_1 = 1$. However, tuning w_1 might slightly improve SegDepthFormer’s performance.

Aleatoric Uncertainty. The respective aleatoric uncertainty is obtained by computing the predictive entropy $H(p(z))$ (see Equation 3) for the segmentation task or by the predictive variance $s^2(z)$ (see Equation 5), which is learned implicitly through the optimization of \mathcal{L}_{GNLL} .

3.3. EMUFormer

In the following, we explain our student-teacher distillation framework for efficient multi-task uncertainties, which we call EMUFormer. Our objective with EMUFormer is threefold:

1. Achieve state-of-the-art joint semantic segmentation and monocular depth estimation results.
2. Estimate well-calibrated predictive uncertainties for both tasks.
3. Avoid introducing additional computational overhead during inference.

In order to achieve these goals, EMUFormer employs a two-step student-teacher distillation framework:

1. Training a teacher with ground truth labels.
2. Training the student with ground truth labels while distilling the teacher’s predictive uncertainties.

Teacher. Although our framework is flexible with regard to the type of teacher, we use a DE that is known for producing high-quality estimates [19, 54, 68].

Student. We propose employing the SegDepthFormer architecture for the student model due to its simplicity, performance, and efficiency. In principle, though, any architecture capable of outputting a semantic segmentation mask along with a predictive mean and variance for monocular depth estimation is suitable.

Distillation Approach. To efficiently estimate predictive uncertainties for semantic segmentation and monocular depth estimation, EMUFormer utilizes student-teacher distillation. Figure 4 shows a schematic overview of EMUFormer. The training is performed with two additional uncertainty-related losses compared to the regular SegDepthFormer. To compute both predictive uncertainties we compute multiple prediction samples from the teacher. Additionally, we add color jittering as an additional data augmentation to the teacher’s input \tilde{x} . Pre-

vious work showed that this is helpful when the training dataset is used for training and distillation to prevent the student from underestimating the epistemic uncertainty of the teacher [33, 57]. The color jitter causes the teacher’s uncertainty distribution on the training dataset to be more closely related to the test-time distribution.

Segmentation Uncertainty Loss. The segmentation uncertainty knowledge of the teacher model is transferred into the student model by using the Kullback-Leibler divergence loss:

$$\mathcal{L}_{\text{KL}} = \sum_{c=1}^C q_c(\tilde{z}) \cdot \log \left(\frac{q_c(\tilde{z})}{p_c(z)} \right), \quad (8)$$

where \tilde{z} are the logits based on the perturbed input image, $q_c(\tilde{z})$ is the teacher’s mean softmax probability map, and $p_c(z)$ is the student’s softmax probability map. Minimizing this loss ensures that the student learns to match the well-calibrated softmax probabilities provided by the teacher, allowing the predictive entropy $H(p(z))$ (see Equation 3) to capture the underlying predictive uncertainty.

Depth Uncertainty Loss. Because it is not possible to match two distributions for the unbound uncertainties in the regression task, we introduce the root mean squared logarithmic error (RMSLE) for the depth uncertainty distillation:

$$\mathcal{L}_{\text{RMSLE}} = \sqrt{\frac{1}{N} \sum_{n=1}^N (\log(\sigma_n^2(\tilde{z}) + 1) - \log(s_n^2(z) + 1))^2}, \quad (9)$$

where $\sigma_n^2(\tilde{z})$ is the teacher’s predictive uncertainty and $s_n^2(z)$ is the student’s predictive uncertainty estimate. The natural logarithm penalizes underestimations more than overestimations, thereby providing special attention to the pixels with higher uncertainties. Minimizing the depth uncertainty loss trains the student to mimic the predictive uncertainty of the teacher. Consequently, the second output channel of the decoder does not only output the aleatoric uncertainty anymore, but rather the more meaningful predictive uncertainty, which additionally covers the epistemic uncertainty.

We follow Loquercio et al. [40] to calculate the predictive uncertainty of the teacher with:

$$\sigma^2(\tilde{z}) = \frac{1}{T} \sum_{t=1}^T s_t^2(\tilde{z}) + \frac{1}{T} \sum_{t=1}^T (\mu_t(\tilde{z}) - \bar{\mu}(\tilde{z}))^2, \quad (10)$$

where T is the number of prediction samples from the teacher, $v(\tilde{z})$ is the predictive variance (see Equation 5, $\mu(\tilde{z})$ is the predictive mean of a sample, and $\bar{\mu}(\tilde{z})$ is the mean predictive mean across all samples.

Training Criterion. In summary, EMUFormer is trained to minimize the weighted sum of four objective functions:

$$\mathcal{L} = \mathcal{L}_{\text{CE}} + w_1 \mathcal{L}_{\text{GNLL}} + w_2 \mathcal{L}_{\text{KL}} + w_3 \mathcal{L}_{\text{RMSLE}}. \quad (11)$$

By setting $w_1 = w_3 = 1$ and $w_2 = 10$, we obtain good results across all of our experiments. However, depending on the application, tuning these hyperparameters may further enhance performance.

4. Experimental Setup

Datasets. We conduct all experiments on Cityscapes [6] and NYUv2 [59]. Cityscapes, with 2975 training and 500 validation images, is a popular urban street scene benchmark dataset. Notably, the depth values are based on the disparity of stereo camera images. NYUv2 contains 795 training and 654 testing images of indoor scenes.

Data Augmentations. Regardless of the trained model, we apply a very common data augmentation strategy:

1. Random scaling with a factor between 0.5 and 2.0.
2. Random cropping with a crop size of 768×768 pixels on Cityscapes and 480×640 pixels on NYUv2.
3. Random horizontal flipping with a flip chance of 50%.

Implementation Details. For all training processes, we use AdamW [41] optimizer with a base learning rate of 0.00006 and employ a polynomial rate scheduler:

$$lr = lr_{\text{base}} \cdot \left(1 - \frac{\text{iteration}}{\text{total iterations}}\right)^{0.9}, \quad (12)$$

where lr is the current learning rate and lr_{base} is the initial base learning rate. Besides, we use a batch size of 8 and train on four NVIDIA A100 GPUs with 40 GB of memory using mixed precision [45]. The encoders of the baseline models are initialized with weights pre-trained on ImageNet [8] and then trained for 250 epochs on Cityscapes and for 100 epochs on NYUv2, respectively. In contrast, EMUFormer is initialized with the weights of a pre-trained SegDepthFormer and fine-tuned for 100 epochs on both datasets. Unless otherwise noted, we use the SegFormer-B2 [69] backbone for all experiments. We do not adopt any of the widely-used methods such as OHEM [58], auxiliary losses, class imbalance compensation, or sliding window testing to keep our approach as simple and transparent as possible.

Metrics. For quantitative evaluations of the semantic segmentation task, we report the mean Intersection over Union (mIoU), also known as the Jaccard Index. Additionally, we use the Expected Calibration Error (ECE) [51] to evaluate the calibration of the softmax probabilities. For the monocular depth estimation task, we use the common root mean squared error (RMSE). Finally, we employ the following uncertainty evaluation metrics proposed by Mukhoti and Gal [49]:

1. $p(\text{accurate}|\text{certain})$: The probability that the model is accurate on its output given that the uncertainty is below a certain threshold.

2. $p(\text{uncertain}|\text{inaccurate})$: The probability that the uncertainty of the model exceeds a certain threshold given that the prediction is inaccurate.
3. $PAvPU$: The combination of both cases, i.e. $\text{accurate}|\text{certain}$ and $\text{inaccurate}|\text{uncertain}$.

Although these metrics have originally been proposed for semantic segmentation [49], we also use them to evaluate the depth regression uncertainties. Since one cannot simply determine whether a depth prediction is accurate, we apply the following formula:

$$\max\left(\frac{\mu(z)}{y}, \frac{y}{\mu(z)}\right) = \delta_1 < 1.25, \quad (13)$$

where $\mu(z)$ is the predicted depth value of a pixel and y is the corresponding ground truth depth [47]. δ_1 serves as a standard metric for quantifying the accuracy of monocular depth estimation models, using 1.25 as the threshold to determine whether a depth prediction is accurate or not. In contrast, δ_2 and δ_3 are less strict, typically utilizing thresholds of 1.25^2 and 1.25^3 , respectively.

For the sake of simplicity and to simulate real-world employment, we set the uncertainty threshold to the mean uncertainty of a given image for all evaluations.

Monte Carlo Dropout. MCD depends primarily on the number of dropout layers, where they are inserted inside the network, and most-importantly the dropout rate. Since the original SegFormer [69] already applies dropout layers throughout the entire network, we follow their work and only consider two dropout rates, 20% and 50%. We sample ten times to obtain the prediction and predictive uncertainty [12, 19, 57].

Deep Sub-Ensemble. Consistent with the DEs and MCD, we train the DSE with ten decoder heads for each task on top of a shared encoder [63]. During training, we only optimize a single decoder head per training batch and alternate between them. Thereby, we aim to introduce as much randomness as possible, analogous to the training of DEs. For inference, we utilize all decoder heads, of course.

Deep Ensemble. DEs achieve the best results if they are trained to explore diverse modes in function space, which we accomplish by randomly initializing all decoder heads, by using random augmentations, and by applying random shuffling of the training data points [10, 30]. Unless otherwise noted, we report results of a DE with ten members, following the suggestions of previous work [10, 30, 33].

Predictions. Regardless of the uncertainty quantification method, we report the results of the mean prediction. For the semantic segmentation task, we compute the mean softmax probability of all samples. For the monocular depth estimation task, we first apply ReLU (see Equation 4) and then compute the mean depth of the corresponding samples.

Uncertainty. For the semantic segmentation task, we compute the predictive entropy (see Equation 3) based on the mean softmax probabilities as a measure for the predictive uncertainty [49]. For the depth estimation task, however, we calculate the predictive uncertainty based on the mean predictive variance and the variance of the depth predictions of the samples (see Equation 10) [40].

5. Joint Uncertainty Evaluation

In this section, we describe the results of our joint uncertainty evaluation quantitatively. We compare combinations of the baseline models SegFormer, DepthFormer, and SegDepthFormer with the uncertainty quantification methods MCD, DSE, and DEs. Tables 2 and 3 contain a detailed quantitative comparison for the different combinations. The focus particularly lies on the uncertainty quality.

Single-task vs. Multi-task. Looking at the differences between the single-task models, SegFormer and DepthFormer, and the multi-task model, SegDepthFormer, the single-task models generally deliver slightly better prediction performance. However, SegDepthFormer exhibits greater uncertainty quality for the semantic segmentation task in comparison to SegFormer. This is particularly evident for $p(\text{uncertain}|\text{inaccurate})$ on Cityscapes. For the depth estimation task, there is no significant difference in terms of uncertainty quality.

Baseline Models. As expected, the baseline models have the lowest inference times, being 5 to 30 times faster without using any uncertainty quantification method. While their prediction performance turns out to be quite competitive, only beaten by DEs, they show poor calibration and uncertainty quality for semantic segmentation. Surprisingly, the uncertainty quality for the depth estimation task is very decent, often only surpassed by the DE.

Monte Carlo Dropout. The use of MCD causes a significantly higher inference time compared to the respective baseline model. Additionally, leaving dropout activated during inference to sample from the posterior has a detrimental effect on the prediction performance, particularly with a 50% dropout ratio. Nevertheless, MCD outputs well-calibrated softmax probabilities and uncertainties, although the results should be interpreted with caution because of the deteriorated prediction quality.

Deep Sub-Ensemble. Across both datasets, DSEs show comparable prediction performance compared with the baseline models. Notably, DSEs consistently demonstrate a high uncertainty quality across all metrics, particularly in the segmentation task on Cityscapes.

Deep Ensemble. In accordance to previous work [19, 54, 68], DEs emerge as state-of-the-art, delivering the best prediction performance and mostly superior uncertainty quality. At the same time, DEs suffer from the highest computational cost.

		Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
		mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
Baseline	SegFormer	0.772	0.033	0.882	0.395	0.797	-	-	-	-	17.90 \pm 0.47
	DepthFormer	-	-	-	-	-	7.452	0.749	0.476	0.766	17.59 \pm 0.82
	SegDepthFormer	0.738	0.028	0.913	0.592	0.826	7.536	0.745	0.472	0.762	22.04 \pm 0.27
MCD (20%)	SegFormer	0.759	0.007	0.883	0.424	0.780	-	-	-	-	177.13 \pm 0.64
	DepthFormer	-	-	-	-	-	7.956	0.749	0.555	0.739	139.32 \pm 0.78
	SegDepthFormer	0.738	0.020	0.911	0.592	0.803	7.370	0.761	0.523	0.757	202.23 \pm 0.39
MCD (50%)	SegFormer	0.662	0.028	0.883	0.485	0.760	-	-	-	-	176.98 \pm 0.53
	DepthFormer	-	-	-	-	-	21.602	0.181	0.366	0.431	139.81 \pm 1.20
	SegDepthFormer	0.640	0.021	0.906	0.616	0.782	8.316	0.733	0.558	0.723	203.82 \pm 0.81
DSE	SegFormer	0.772	0.037	0.890	0.456	0.797	-	-	-	-	132.30 \pm 3.16
	DepthFormer	-	-	-	-	-	7.036	0.762	0.467	0.772	91.82 \pm 2.01
	SegDepthFormer	0.749	0.009	0.931	0.696	0.844	7.441	0.751	0.463	0.766	212.11 \pm 8.44
DE	SegFormer	0.784	0.033	0.887	0.416	0.798	-	-	-	-	667.51 \pm 2.89
	DepthFormer	-	-	-	-	-	7.222	0.759	0.486	0.771	626.79 \pm 2.05
	SegDepthFormer	0.755	0.015	0.917	0.609	0.828	7.156	0.763	0.493	0.773	743.23 \pm 32.95

Table 2. Quantitative comparison on the Cityscapes dataset [6] between the three baseline models paired with MCD, DSE, and DEs, respectively. Best results are marked in **bold**.

		Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
		mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
Baseline	SegFormer	0.470	0.159	0.768	0.651	0.734	-	-	-	-	18.09 \pm 0.41
	DepthFormer	-	-	-	-	-	0.554	0.786	0.449	0.610	17.51 \pm 0.87
	SegDepthFormer	0.466	0.151	0.769	0.659	0.733	0.558	0.776	0.446	0.594	22.31 \pm 0.23
MCD (20%)	SegFormer	0.422	0.102	0.767	0.706	0.724	-	-	-	-	222.67 \pm 0.61
	DepthFormer	-	-	-	-	-	0.605	0.741	0.478	0.568	139.58 \pm 0.52
	SegDepthFormer	0.433	0.093	0.771	0.710	0.725	0.610	0.731	0.450	0.560	251.25 \pm 0.81
MCD (50%)	SegFormer	0.273	0.083	0.705	0.722	0.713	-	-	-	-	223.25 \pm 0.82
	DepthFormer	-	-	-	-	-	0.978	0.516	0.492	0.526	139.27 \pm 0.69
	SegDepthFormer	0.272	0.084	0.702	0.721	0.711	0.837	0.576	0.473	0.525	251.98 \pm 0.60
DSE	SegFormer	0.469	0.092	0.776	0.681	0.726	-	-	-	-	180.42 \pm 3.93
	DepthFormer	-	-	-	-	-	0.547	0.782	0.423	0.596	91.66 \pm 0.26
	SegDepthFormer	0.461	0.077	0.776	0.692	0.723	0.584	0.738	0.403	0.573	261.69 \pm 5.10
DE	SegFormer	0.486	0.125	0.782	0.675	0.734	-	-	-	-	715.97 \pm 7.55
	DepthFormer	-	-	-	-	-	0.524	0.808	0.475	0.613	624.30 \pm 2.07
	SegDepthFormer	0.481	0.122	0.783	0.682	0.733	0.552	0.785	0.453	0.590	788.76 \pm 2.00

Table 3. Quantitative comparison on the NYUv2 dataset [59] between the three baseline models paired with MCD, DSE, and DEs, respectively. Best results are marked in **bold**.

6. Efficient Multi-task Uncertainties

In this section, we conduct several experiments to demonstrate the efficiency and efficacy of EMUFormer. We begin by comparing EMUFormer’s performance with its DE teacher for multiple backbones. Subsequently, we compare our results with previous state-of-the-art approaches, followed by qualitative examples. Lastly, we provide an ablation study on the impact of the GNLL loss.

6.1. Quantitative Evaluation

Baseline vs. Teacher vs. Student. We present a comprehensive analysis in Tables 4 and 5 by com-

paring SegDepthFormer (baseline), SegDepthFormer DE (teacher), and EMUFormer (student). EMUFormer emerges as the standout performer, surpassing the baseline SegDepthFormer model across all metrics on both datasets, with only a single exception. Remarkably, this performance is achieved while maintaining an equivalent inference time. Remarkably, EMUFormer even outperforms the SegDepthFormer DE, which served as its teacher and has approximately 33 times higher inference time, in most cases. In terms of prediction performance, EMUFormer gives slightly worse segmentation results compared to the DE. However, it notably excels in the depth estimation task, especially on Cityscapes [6], which is a phenomenon we

	Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
	mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
SegDepthFormer (Baseline)	0.738	0.028	0.913	0.592	0.826	7.536	0.745	0.472	0.762	22.04 \pm 0.27
SegDepthFormer (DE)	0.755	0.015	0.917	0.609	0.828	7.156	0.763	0.493	0.773	743.23 \pm 32.95
EMUFormer	0.752	0.012	0.923	0.658	0.811	6.983	0.772	0.491	0.783	22.04 \pm 0.27

Table 4. Quantitative comparison on the Cityscapes dataset [6] between the baseline SegDepthFormer, a SegDepthFormer Deep Ensemble, which acts as the teacher with ten members, and our EMUFormer. Best results are marked in **bold**.

	Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
	mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
SegDepthFormer (Baseline)	0.466	0.151	0.769	0.659	0.733	0.558	0.776	0.446	0.594	22.31 \pm 0.23
SegDepthFormer (DE)	0.481	0.122	0.783	0.682	0.733	0.552	0.785	0.453	0.590	788.76 \pm 2.00
EMUFormer	0.475	0.129	0.787	0.692	0.737	0.514	0.810	0.440	0.633	22.31 \pm 0.23

Table 5. Quantitative comparison on the NYUv2 dataset [59] between the baseline SegDepthFormer, a SegDepthFormer Deep Ensemble, which acts as the teacher with ten members, and our EMUFormer. Best results are marked in **bold**.

		Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
		mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
B0 Backbone	SegFormer (DE)	0.689	0.037	0.888	0.486	0.779	-	-	-	-	273.20 \pm 1.38
	DepthFormer (DE)	-	-	-	-	-	8.452	0.692	0.414	0.719	236.13 \pm 0.70
	SegDepthFormer (DE)	0.651	0.045	0.912	0.634	0.803	8.495	0.692	0.425	0.718	317.47 \pm 15.64
	EMUFormer	0.630	0.023	0.924	0.714	0.791	8.086	0.717	0.473	0.732	9.58 \pm 0.07
B5 Backbone	SegFormer (DE)	0.809	0.032	0.896	0.435	0.819	-	-	-	-	1931.01 \pm 12.77
	DepthFormer (DE)	-	-	-	-	-	6588	0.782	0.487	0.791	1892.47 \pm 9.24
	SegDepthFormer (DE)	0.789	0.037	0.928	0.657	0.852	6.664	0.785	0.502	0.792	2018.04 \pm 32.31
	EMUFormer	0.771	0.014	0.934	0.703	0.845	6.157	0.804	0.536	0.799	50.72 \pm 0.45

Table 6. Quantitative comparison on the Cityscapes dataset [6] between the three baseline models as Deep Ensembles and EMUFormer with SegFormer’s B0 and B5 backbone [69]. The respective SegDepthFormer Deep Ensemble served as the teacher for the corresponding EMUFormer. Best results are marked in **bold**.

		Semantic Segmentation					Monocular Depth Estimation				Inference Time [ms]
		mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	
B0 Backbone	SegFormer (DE)	0.376	0.105	0.743	0.701	0.718	-	-	-	-	315.42 \pm 2.41
	DepthFormer (DE)	-	-	-	-	-	0.642	0.720	0.476	0.566	227.92 \pm 2.39
	SegDepthFormer (DE)	0.375	0.097	0.744	0.703	0.718	0.678	0.693	0.466	0.553	346.21 \pm 2.72
	EMUFormer	0.363	0.090	0.743	0.713	0.720	0.674	0.705	0.498	0.558	10.04 \pm 0.06
B5 Backbone	SegFormer (DE)	0.534	0.138	0.792	0.653	0.744	-	-	-	-	1958.46 \pm 36.71
	DepthFormer (DE)	-	-	-	-	-	0.468	0.852	0.505	0.647	1875.53 \pm 12.83
	SegDepthFormer (DE)	0.526	0.133	0.794	0.665	0.743	0.451	0.838	0.478	0.619	2038.26 \pm 13.06
	EMUFormer	0.520	0.134	0.798	0.688	0.744	0.476	0.846	0.467	0.647	52.27 \pm 1.40

Table 7. Quantitative comparison on the NYUv2 dataset [59] between the three baseline models as Deep Ensembles and EMUFormer with SegFormer’s B0 and B5 backbone [69]. The respective SegDepthFormer Deep Ensemble served as the teacher for the corresponding EMUFormer. Best results are marked in **bold**.

observed across multiple experiments (cf. Tables 6, 7, and 8) and which we will discuss in Section 7.

Backbone Size. Tables 6 and 7 display a comprehensive assessment of the influence of the backbone size on Cityscapes [6] and NYUv2 [59]. In this context, we de-

cidated to evaluate the three baseline models as a DE with ten members each in comparison to EMUFormer for the smallest, B0, and the biggest, B5, backbone of SegFormer [69], respectively. The findings broadly align with the earlier observations of Section 5 in terms of single-tasking versus

	NYUv2		Cityscapes	
	mIoU \uparrow	RMSE \downarrow	mIoU \uparrow	RMSE \downarrow
HybridNet A2 [36]	0.343	0.682	0.666	12.09
Mousavian et al. [48]	0.392	0.816	-	-
C-DCNN [38]	0.398	0.628	-	-
BMTAS [3]	0.411	0.543	-	-
Gao et al. [15]	0.419	0.528	-	-
Nekrasov et al. [52]	0.420	0.565	-	-
CI-Net [14]	0.426	0.504	0.701	6.880
Wang et al. [67]	0.442	0.745	-	-
SOSD-Net [20]	0.450	0.514	0.682	-
ATRC [4]	0.463	0.536	-	-
MTI-Net [65]	0.490	0.529	-	-
PAD-Net [70]	0.502	0.582	0.761	-
MTFormer [71]	0.506	<u>0.483</u>	-	-
SegDepthFormer-B2 (Ours)	0.476	0.549	0.763	7.286
SegDepthFormer-B5 (Ours)	<u>0.518</u>	0.499	0.784	<u>6.819</u>
EMUFormer-B2 (Ours)	0.475	0.514	0.752	6.983
EMUFormer-B5 (Ours)	0.520	0.476	<u>0.771</u>	6.157

Table 8. Comparison against previous state-of-the-art approaches in joint semantic segmentation and monocular depth estimation. Best results are marked in **bold**, second best results are underlined.

multi-tasking. More specifically, EMUFormer emerges as the top performer on all segmentation metrics, except for the mIoU where the SegFormer DE gives slightly better results. On the Cityscapes dataset, EMUFormer stands out by delivering the best results for all depth metrics across both backbones. Notably, it achieves this superior performance while maintaining a 20 to 30 times faster inference time compared to the DEs. On NYUv2, the DepthFormer DE performs marginally better on the depth metrics, although EMUFormer remains highly competitive, especially if inference time is considered.

Comparison with SOTA. On both datasets, Cityscapes [6] and NYUv2 [59], EMUFormer-B5 outperforms the previous state-of-the-art in joint semantic segmentation and monocular depth estimation. For instance, on NYUv2 [59], EMUFormer delivers 1.4% higher mIoU and 0.007 lower RMSE than MTFormer [71], which also adopts a modern Vision-Transformer-based architecture. In contrast to our work, however, they rely on cross-task attention mechanisms and on a sophisticated self-supervised pre-training routine, which introduce additional complexity. Our SegDepthFormer-B5 baseline model already achieves very competitive results without such adaptations to the architecture or the training routine. It improves upon previous work in all cases except for RMSE on NYUv2 [59]. Besides, even with the lightweight B2 models, we achieve very decent results in comparison to prior work, offering an alternative for real-time applications.

6.2. Qualitative Evaluation

In addition to the quantitative evaluation, we also provide qualitative examples of EMUFormer-B2 in Figure 5 for Cityscapes [6] and NYUv2 [59].

Cityscapes. On Cityscapes, EMUFormer demonstrates good prediction performance for both tasks. In the segmentation task, its uncertainty prediction proves particularly insightful as the red rectangles highlight. For example, in areas such as the car hood, which is not part of the training labels (indicated by black pixels), the model exhibits high uncertainty, indicating its ability to capture out-of-distribution information or epistemic uncertainty. Similarly, in noisy background areas, the model effectively captures the aleatoric noise. Additionally, the model correctly predicts high uncertainties for challenging areas like the wall on the right of the image, highlighting the utility of uncertainties in identifying potential model errors. In the depth estimation task, analogous to the segmentation task, EMUFormer predicts high uncertainty on the car hood or the sky, which are both areas that are not part of the training ground truth, i.e. areas of high epistemic uncertainty. Furthermore, the uncertainty is appropriately high at object boundaries, indicating sensitivity to significant depth variations.

NYUv2. For the segmentation task, EMUFormer again outputs high uncertainties for pixels that are not part of the ground truth or those that are misclassified, consistently providing useful predictive uncertainties. In the depth estimation task, the uncertainties seem to correlate with the estimated depth, providing an intuitive and helpful indication. This alignment suggests that the model effectively captures the depth prediction quality, particularly as it relates to increasing distances.

In summary, the qualitative evaluation aligns with the quantitative findings of Section 6.1 and highlights EMUFormer’s proficiency in handling both the segmentation and the depth estimation tasks, showcasing its ability to generate meaningful predictive uncertainties that enable more thorough interpretations of the predictions.

6.3. Ablation Studies

Impact of GNLL Loss. EMUFormer is trained to minimize the weighted sum of the following four objective functions:

1. Cross-Entropy loss for the semantic segmentation task.
2. Kullback-Leibler divergence loss for the segmentation uncertainty distillation.
3. Gaussian Negative Log-Likelihood loss for the monocular depth estimation task.
4. Root mean squared error loss for the depth uncertainty distillation.

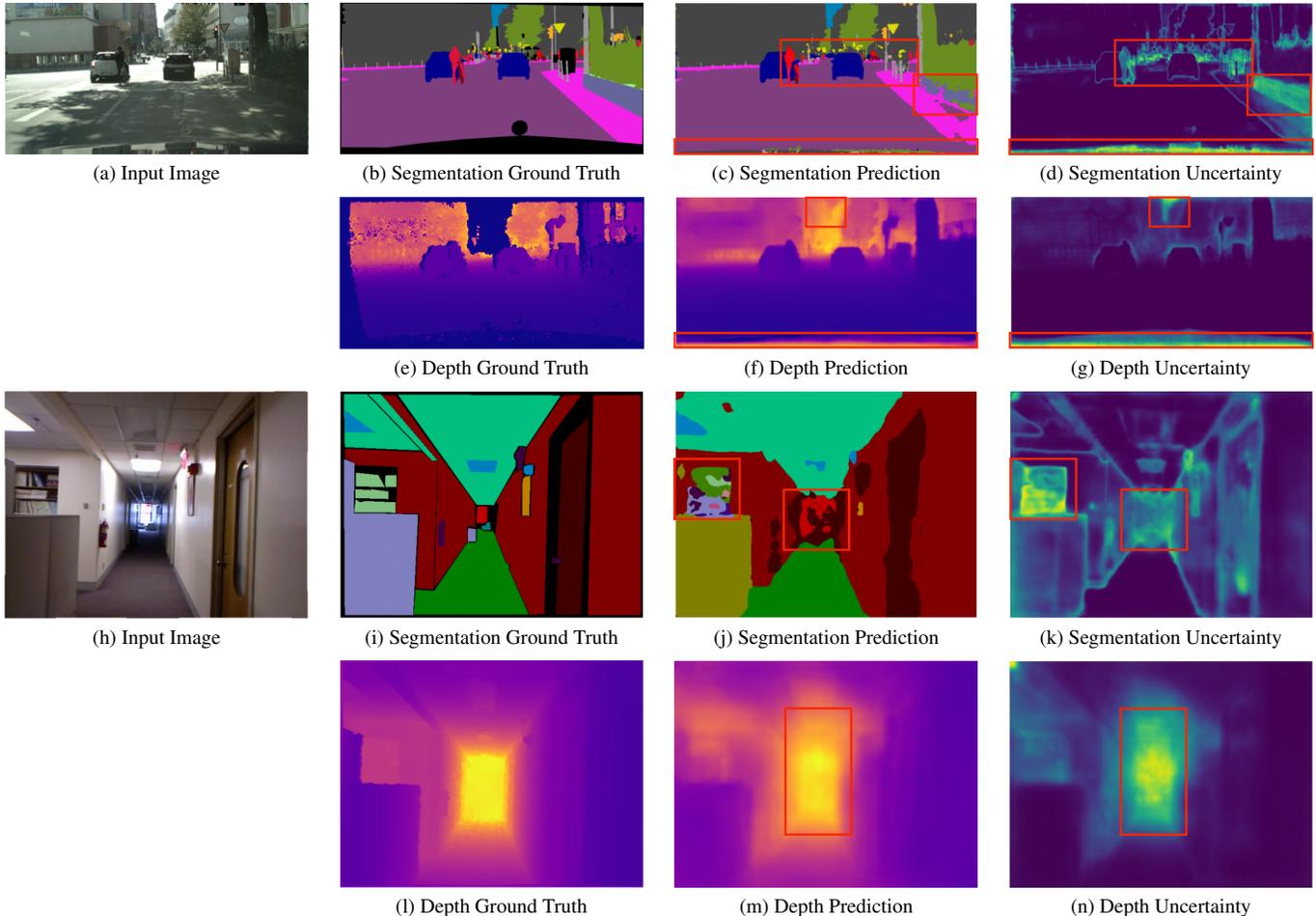


Figure 5. Qualitative examples of our EMUFormer-B2 on the Cityscapes [6] (top) and NYUv2 [59] (bottom) datasets. Red rectangles are added to highlight interesting areas. Best viewed in color.

As described in Section 3.2.2 and shown by Equation 6, GNLL treats every prediction as a sample from a Gaussian distribution with a predictive mean and a corresponding predictive variance. Usually, these variances are solely learned implicitly through the optimization of the predictive means based on the ground truth labels. In the case of EMUFormer, however, the network is also being trained to mimic the predictive uncertainty of the teacher in parallel. Consequently, the depth uncertainty does not need to be learned implicitly, rather it can be used to improve the depth estimation itself. In order to explore this more thoroughly, we performed an ablation study on the impact of the GNLL loss by replacing the GNLL loss with the Mean Squared Error (MSE) loss and the Huber loss [25], respectively.

Tables 9 and 10 show a quantitative comparison of the impact of the respective depth loss for EMUFormer-B2 on the Cityscapes and NYUv2 datasets. On Cityscapes, training with GNLL loss leads to the best performance across

the board, especially with regard to the RMSE for monocular depth estimation. GNLL loss results in a RMSE of 6.983 in comparison to 7.217 and 7.340 for MSE and Huber loss [25], respectively. Similarly, on NYUv2, training with GNLL loss yields the best RMSE with 0.514 versus 0.527 and 0.533 for MSE and Huber loss [25], although at the cost of a very slight deterioration of 0.006 in mIoU. Remarkably, using GNLL loss leads to the highest depth uncertainty quality for both datasets.

Overall, these results show that incorporating the predictive uncertainties using the GNLL loss enhances the performance of EMUFormer for depth estimation and depth uncertainty quantification compared to other loss functions like MSE or Huber loss [25] that do not account for the uncertainty. We consider this a valuable insight and believe that leveraging high-quality predictive uncertainties during the optimization process offers great potential for future work.

	Semantic Segmentation					Monocular Depth Estimation			
	mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow
MSE	0.749	0.014	0.922	0.659	0.810	7.217	0.742	0.446	0.761
Huber [25]	0.748	0.013	0.923	0.657	0.809	7.340	0.743	0.446	0.760
GNLL	0.752	0.012	0.923	0.658	0.811	6.983	0.772	0.491	0.783

Table 9. Ablation study on the impact of the depth loss on the results of EMUFormer-B2 on Cityscapes [6]. Best results are marked in bold.

	Semantic Segmentation					Monocular Depth Estimation			
	mIoU \uparrow	ECE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow	RMSE \downarrow	p(acc/cer) \uparrow	p(inacc/unc) \uparrow	PAvPU \uparrow
MSE	0.481	0.127	0.788	0.690	0.737	0.527	0.788	0.431	0.587
Huber [25]	0.481	0.127	0.788	0.689	0.737	0.533	0.786	0.431	0.587
GNLL	0.475	0.129	0.787	0.692	0.737	0.514	0.810	0.440	0.633

Table 10. Ablation study on the impact of the depth loss on the results of EMUFormer-B2 on NYUv2 [59]. Best results are marked in bold.

7. Discussion

Joint Uncertainty Evaluation. Quantifying the uncertainty in joint segmentation and depth estimation has not been thoroughly examined in prior research. Therefore, we evaluated multiple uncertainty quantification methods with modern Vision-Transformer-based architectures for joint semantic segmentation and monocular depth estimation. In general, single-task models demonstrate slightly better prediction performance, which may arise from multiple factors. For one, the single-task models can optimize all available parameters for their specific task. Additionally, the multi-task models do not exploit any sophisticated adaptations to the architecture or the training process. Unlike previous work [3, 4, 14, 20, 26, 27, 29, 36, 38, 39, 48, 52, 67, 70, 71], we intentionally left out all of the complexities for the joint uncertainty evaluation in order to maintain methodological simplicity and transparency of the results. Interestingly, multi-task models showcase greater uncertainty quality, particularly in the context of the semantic segmentation task. This suggests that jointly training a model to solve multiple tasks can enhance the model’s ability to better quantify its uncertainty. In terms of uncertainty quantification methods, DEs stand out as the preferred choice, demonstrating superior prediction performance and, for the most part, higher uncertainty quality. However, it is crucial to note that this advantage comes at the highest computational cost. Both findings align closely with previous work focusing on the evaluation of uncertainties [19, 54, 68]. Among the more efficient methods, MCD and DSE, the latter exhibits a prediction performance that is comparable with the baseline models while achieving a high uncertainty qual-

ity. This positions DSEs as an attractive alternative to DEs, offering efficiency without significant sacrifices in performance or uncertainty quality.

EMUFormer. In addition to the joint uncertainty evaluation, we also proposed EMUFormer, which employs student-teacher distillation to achieve state-of-the-art results in joint semantic segmentation and monocular depth estimation on Cityscapes [6] and NYUv2 [59]. Notably, it accomplishes this while estimating well-calibrated predictive uncertainties for both tasks, all without introducing any additional computational overhead during inference. Remarkably, EMUFormer even surpasses the performance of its DE teacher in certain cases, despite the latter having ten times the parameters and approximately 30 times higher inference time. The backbone ablation analysis further reinforces the versatility of our proposed method, showcasing its efficacy across different backbone configurations. Most interestingly, however, EMUFormer achieves particularly outstanding performance in the depth estimation task in comparison to the teacher. We primarily attribute this success to the use of the Gaussian Negative Log-Likelihood loss (cf. Section 6.3), which is commonly employed to implicitly learn corresponding variances in addition to the predictive means. In the case of EMUFormer, however, the teacher model already provides high-quality variances through distillation, allowing for a more accurate approximation of the predictive means and their associated uncertainties. Consequently, leveraging uncertainties during the training, either implicitly like EMUFormer, or explicitly like previous work [29, 32], is an interesting venue for future work.

8. Conclusion

In this work, we first combine multiple uncertainty quantification methods with joint semantic segmentation and monocular depth estimation and evaluate how they perform in comparison to each other. Quantitative evaluations revealed that Deep Ensembles stand out as the preferred choice concerning prediction performance and uncertainty quality, although having the highest computational cost. Among the less costly methods, Deep Sub-Ensembles emerge as an attractive alternative to Deep Ensembles, offering efficiency without major sacrifices in prediction performance or uncertainty quality. Additionally, we reveal the benefits of multi-task learning with regard to the uncertainty quality compared to solving both tasks separately. Building on these insights, we propose EMUFormer, a novel student-teacher distillation approach for joint semantic segmentation and monocular depth estimation as well as efficient multi-task uncertainty quantification. By implicitly leveraging the predictive uncertainties of the teacher, EMUFormer achieves new state-of-the-art results on Cityscapes and NYUv2 for both tasks. Notably, EMUFormer also manages to estimate high-quality predictive uncertainties for both tasks that are comparable or superior to a DE despite being an order of magnitude more efficient.

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] Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. Deep evidential regression. *Advances in Neural Information Processing Systems*, 33:14927–14937, 2020. [1](#), [3](#), [4](#)
- [2] Victor Besnier, David Picard, and Alexandre Briot. Learning uncertainty for safety-oriented semantic segmentation in autonomous driving. In *2021 IEEE International Conference on Image Processing (ICIP)*, pages 3353–3357. IEEE, 2021. [1](#), [3](#)
- [3] David Brüggenmann, Menelaos Kanakis, Stamatios Georgoulis, and Luc Van Gool. Automated search for resource-efficient branched multi-task networks. *arXiv preprint arXiv:2008.10292*, 2020. [2](#), [11](#), [13](#)
- [4] David Brüggenmann, Menelaos Kanakis, Anton Obukhov, Stamatios Georgoulis, and Luc Van Gool. Exploring relational context for multi-task dense prediction. In *Proceedings of the IEEE/CVF international conference on computer vision*, pages 15869–15878, 2021. [2](#), [11](#), [13](#)
- [5] Liangfu Chen, Zeng Yang, Jianjun Ma, and Zheng Luo. Driving scene perception network: Real-time joint detection, depth estimation and semantic segmentation. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1283–1291, 2018. [1](#)
- [6] 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. [7](#), [9](#), [10](#), [11](#), [12](#), [13](#)
- [7] Didan Deng, Liang Wu, and Bertram E. Shi. Iterative distillation for better uncertainty estimates in multitask emotion recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 3557–3566, October 2021. [3](#)
- [8] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. ImageNet: A large-scale hierarchical image database. In *2009 IEEE Conference on Computer Vision and Pattern Recognition*, pages 248–255, Miami, FL, 2009. IEEE. [7](#)
- [9] Xingshuai Dong, Matthew A Garratt, Sreenatha G Anavatti, and Hussein A Abbass. Towards real-time monocular depth estimation for robotics: A survey. *IEEE Transactions on Intelligent Transportation Systems*, 23(10):16940–16961, 2022. [1](#)
- [10] Stanislav Fort, Huiyi Hu, and Balaji Lakshminarayanan. Deep Ensembles: A Loss Landscape Perspective. *arXiv:1912.02757*, 2020. [3](#), [8](#)
- [11] Yarin Gal. Uncertainty in deep learning. *Ph.D. thesis, University of Cambridge*, 2016. [3](#)
- [12] 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. [1](#), [3](#), [8](#)
- [13] Yarin Gal, Riashat Islam, and Zoubin Ghahramani. Deep bayesian active learning with image data. In *International conference on machine learning*, pages 1183–1192. PMLR, 2017. [3](#)
- [14] Tianxiao Gao, Wu Wei, Zhongbin Cai, Zhun Fan, Sheng Quan Xie, Xinmei Wang, and Qiuda Yu. Ci-net: A joint depth estimation and semantic segmentation network using contextual information. *Applied Intelligence*, 52(15):18167–18186, 2022. [2](#), [11](#), [13](#)
- [15] Tianxiao Gao, Wu Wei, Xinmei Wang, Qiuda Yu, and Zhun Fan. Predictive uncertainties for multi-task learning network. In *International Conference on Advanced Algorithms and Neural Networks (AANN 2022)*, volume 12285, pages 294–300. SPIE, 2022. [2](#), [11](#)
- [16] Jakob Gawlikowski, Cedrique Rovile Njiteutcheu 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. [1](#)
- [17] 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. **1**
- [18] Corina Gurau, Alex Bewley, and Ingmar Posner. Dropout distillation for efficiently estimating model confidence. *arXiv preprint arXiv:1809.10562*, 2018. **3**
- [19] 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. **3, 6, 8, 13**
- [20] Lei He, Jiwen Lu, Guanghui Wang, Shiyu Song, and Jie Zhou. Ssd-net: Joint semantic object segmentation and depth estimation from monocular images. *Neurocomputing*, 440:251–263, 2021. **2, 11, 13**
- [21] Michael Heizmann, Alexander Braun, Markus Glitzner, Matthias Günther, Günther Hasna, Christina Klüver, Jakob Krooß, Erik Marquardt, Michael Overdick, and Markus Ulrich. Implementing machine learning: chances and challenges. *at-Automatisierungstechnik*, 70(1):90–101, 2022. **1**
- [22] Geoffrey Hinton, Oriol Vinyals, and Jeffrey Dean. Distilling the knowledge in a neural network. In *NIPS Deep Learning and Representation Learning Workshop*, 2015. **3**
- [23] Christopher J. Holder and Muhammad Shafique. Efficient uncertainty estimation in semantic segmentation via distillation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV) Workshops*, pages 3087–3094, October 2021. **1, 3**
- [24] Yaosi Hu, Zhenzhong Chen, and Weiyao Lin. Rgb-d semantic segmentation: a review. In *2018 IEEE International Conference on Multimedia & Expo Workshops (ICMEW)*, pages 1–6. IEEE, 2018. **2**
- [25] Peter J Huber. Robust estimation of a location parameter. In *Breakthroughs in statistics: Methodology and distribution*, pages 492–518. Springer, 1992. **12, 13**
- [26] Naihua Ji, Huiqian Dong, Fanyun Meng, and Liping Pang. Semantic segmentation and depth estimation based on residual attention mechanism. *Sensors*, 23(17):7466, 2023. **2, 13**
- [27] Jianbo Jiao, Ying Cao, Yibing Song, and Rynson Lau. Look deeper into depth: Monocular depth estimation with semantic booster and attention-driven loss. In *Proceedings of the European conference on computer vision (ECCV)*, pages 53–69, 2018. **2, 13**
- [28] 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. **3, 4, 5**
- [29] Alex Kendall, Yarin Gal, and Roberto Cipolla. Multi-task learning using uncertainty to weigh losses for scene geometry and semantics. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7482–7491, 2018. **2, 13**
- [30] 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. **1, 3, 4, 5, 8**
- [31] 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. **1**
- [32] Steven Landgraf, Markus Hillemann, Kira Wursthorn, and Markus Ulrich. U-ce: Uncertainty-aware cross-entropy for semantic segmentation. *arXiv preprint arXiv:2307.09947*, 2023. **1, 13**
- [33] Steven Landgraf, Kira Wursthorn, Markus Hillemann, and Markus Ulrich. Dudes: Deep uncertainty distillation using ensembles for semantic segmentation. *arXiv preprint arXiv:2303.09843*, 2023. **1, 3, 7, 8**
- [34] Kimin Lee, Honglak Lee, Kibok Lee, and Jinwoo Shin. Training Confidence-calibrated Classifiers for Detecting Out-of-Distribution Samples. *arXiv:1711.09325*, 2018. **1**
- [35] Christian Leibig, 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. **1**
- [36] Xiao Lin, Dalila Sánchez-Escobedo, Josep R Casas, and Montse Pardàs. Depth estimation and semantic segmentation from a single rgb image using a hybrid convolutional neural network. *Sensors*, 19(8):1795, 2019. **2, 11, 13**
- [37] 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. **1, 3, 4**
- [38] Jing Liu, Yuhang Wang, Yong Li, Jun Fu, Jianguyun Li, and Hanqing Lu. Collaborative deconvolutional neural networks for joint depth estimation and semantic segmentation. *IEEE transactions on neural networks and learning systems*, 29(11):5655–5666, 2018. **2, 11, 13**
- [39] Shikun Liu, Edward Johns, and Andrew J Davison. End-to-end multi-task learning with attention. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1871–1880, 2019. **2, 13**
- [40] Antonio Loquercio, Mattia Segù, and Davide Scaramuzza. A general framework for uncertainty estimation in deep learning. *IEEE Robotics and Automation Letters*, 5(2):3153–3160, 2020. **1, 5, 7, 8**
- [41] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017. **7**
- [42] David J. C. MacKay. A Practical Bayesian Framework for Backpropagation Networks. *Neural Computation*, 4(3):448–472, 1992. **1, 3**
- [43] Andrey Malinin, Bruno Mlodozieniec, and Mark Gales. Ensemble Distribution Distillation. *arXiv:1905.00076*, 2019. **3**
- [44] 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. 1
- [45] 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. 7
- [46] 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. 1
- [47] Yue Ming, Xuyang Meng, Chunxiao Fan, and Hui Yu. Deep learning for monocular depth estimation: A review. *Neuro-computing*, 438:14–33, 2021. 8
- [48] Arsalan Mousavian, Hamed Pirsiavash, and Jana Koščeká. Joint semantic segmentation and depth estimation with deep convolutional networks. In *2016 Fourth International Conference on 3D Vision (3DV)*, pages 611–619. IEEE, 2016. 2, 11, 13
- [49] Jishnu Mukhoti and Yarin Gal. Evaluating bayesian deep learning methods for semantic segmentation. *arXiv preprint arXiv:1811.12709*, 2018. 1, 7, 8
- [50] 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. 1, 3, 4
- [51] Mahdi Pakdaman Naeni, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Proceedings of the AAAI conference on artificial intelligence*, volume 29, 2015. 7
- [52] Vladimir Nekrasov, Thanuja Dharmasiri, Andrew Spek, Tom Drummond, Chunhua Shen, and Ian Reid. Real-time joint semantic segmentation and depth estimation using asymmetric annotations. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 7101–7107. IEEE, 2019. 1, 2, 11, 13
- [53] David A Nix and Andreas S Weigend. Estimating the mean and variance of the target probability distribution. In *Proceedings of 1994 IEEE international conference on neural networks (ICNN'94)*, volume 1, pages 55–60. IEEE, 1994. 5
- [54] 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. 3, 6, 8, 13
- [55] Adriana Romero, Nicolas Ballas, Samira Ebrahimi Kahou, Antoine Chassang, Carlo Gatta, and Yoshua Bengio. FitNets: Hints for Thin Deep Nets. *arXiv:1412.6550*, 2015. 3
- [56] Adrian Schwaiger, Poulami Sinhamahapatra, Jens Gansloer, and Karsten Roscher. Is uncertainty quantification in deep learning sufficient for out-of-distribution detection? *Aisafety@ ijcai*, 54, 2020. 3
- [57] Yichen Shen, Zhilu Zhang, Mert R. Sabuncu, and Lin Sun. Real-time uncertainty estimation in computer vision via uncertainty-aware distribution distillation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*, pages 707–716, January 2021. 1, 3, 7, 8
- [58] 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. 7
- [59] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *Computer Vision—ECCV 2012: 12th European Conference on Computer Vision, Florence, Italy, October 7-13, 2012, Proceedings, Part V 12*, pages 746–760. Springer, 2012. 2, 7, 9, 10, 11, 12, 13
- [60] Ivor JA Simpson, Sara Vicente, and Neill DF Campbell. Learning structured gaussians to approximate deep ensembles. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 366–374, 2022. 1, 3
- [61] 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. 3
- [62] Carsten Steger, Markus Ulrich, and Christian Wiedemann. *Machine Vision Algorithms and Applications*. John Wiley & Sons, 2018. 1
- [63] Matias Valdenegro-Toro. Sub-ensembles for fast uncertainty estimation in neural networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4119–4127, 2023. 1, 3, 8
- [64] 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. 1, 3, 4
- [65] Simon Vandenhende, Stamatios Georgoulis, and Luc Van Gool. Mti-net: Multi-scale task interaction networks for multi-task learning. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pages 527–543. Springer, 2020. 2, 11
- [66] Changshuo Wang, Chen Wang, Weijun Li, and Haining Wang. A brief survey on rgb-d semantic segmentation using deep learning. *Displays*, 70:102080, 2021. 2
- [67] Peng Wang, Xiaohui Shen, Zhe Lin, Scott Cohen, Brian Price, and Alan L Yuille. Towards unified depth and semantic prediction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2800–2809, 2015. 2, 11, 13
- [68] Kira Wurstthorn, 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. 3, 6, 8, 13

- [69] Enze Xie, Wenhai Wang, Zhiding Yu, Anima Anandkumar, Jose M Alvarez, and Ping Luo. Segformer: Simple and efficient design for semantic segmentation with transformers. *Advances in Neural Information Processing Systems*, 34:12077–12090, 2021. [2](#), [4](#), [5](#), [6](#), [7](#), [8](#), [10](#)
- [70] Dan Xu, Wanli Ouyang, Xiaogang Wang, and Nicu Sebe. Pad-net: Multi-tasks guided prediction-and-distillation network for simultaneous depth estimation and scene parsing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 675–684, 2018. [2](#), [11](#), [13](#)
- [71] Xiaogang Xu, Hengshuang Zhao, Vibhav Vineet, Ser-Nam Lim, and Antonio Torralba. Mtfomer: Multi-task learning via transformer and cross-task reasoning. In *European Conference on Computer Vision*, pages 304–321. Springer, 2022. [2](#), [11](#), [13](#)
- [72] Yu Zhang and Qiang Yang. A survey on multi-task learning. *IEEE Transactions on Knowledge and Data Engineering*, 34(12):5586–5609, 2021. [2](#)