

# Uncertainty depth estimation with gated images for 3D reconstruction

Stefanie Walz<sup>1,2</sup>, Tobias Gruber<sup>1,2</sup>, Werner Ritter<sup>1</sup>, Klaus Dietmayer<sup>2</sup>

**Abstract**—Gated imaging is an emerging sensor technology for self-driving cars that provides high-contrast images even under adverse weather influence. It has been shown that this technology can even generate high-fidelity dense depth maps with accuracy comparable to scanning LiDAR systems. In this work, we extend the recent *Gated2Depth* framework with aleatoric uncertainty providing an additional confidence measure for the depth estimates. This confidence can help to filter out uncertain estimations in regions without any illumination. Moreover, we show that training on dense depth maps generated by LiDAR depth completion algorithms can further improve the performance.

## I. INTRODUCTION

Self-driving vehicles require a very detailed perception of their environment in order to move around safely. While 3D information is crucial to understand the scenery and to detect free space, information about texture enables to classify objects and predict their interaction. The environment can be perceived by a variety of sensors, each of them with benefits and drawbacks. Cameras are low-priced and provide high-resolution images that are essential for capturing textures. However, 3D scene reconstruction from a single image is an ill-posed problem [1] and stereo setups are limited in range resolution [2]. In contrast, laser scanners offer very accurate depth measurements, but at very low spatial resolution and for a high cost [3]. In addition to 3D information, radar systems can measure the velocity of objects. Despite of the huge progress in radar development in recent years, radars are still low resolution and can capture neither texture nor fine 3D structures. While radars still perform reasonable in adverse weather conditions, laser scanners completely fail in scattering environments [4] and standard cameras suffer from strong contrast degeneration [5].

For safe autonomous driving in any conditions, all of these sensors are probably required because there is no single sensor that can solve everything reliably and sensors fail asymmetrically under adverse weather [4], [5], [7]. Therefore, the fusion of multiple sensors is the key for reliable and accurate environment perception. For sensor fusion, it is extremely helpful if every sensor does not only deliver measurements but also a measure about the quality or uncertainty of these measurements. Many stereo algorithms [8], [9] for example provide a confidence metric that rates the provided disparity measurement. Sensor stream uncertainties help the

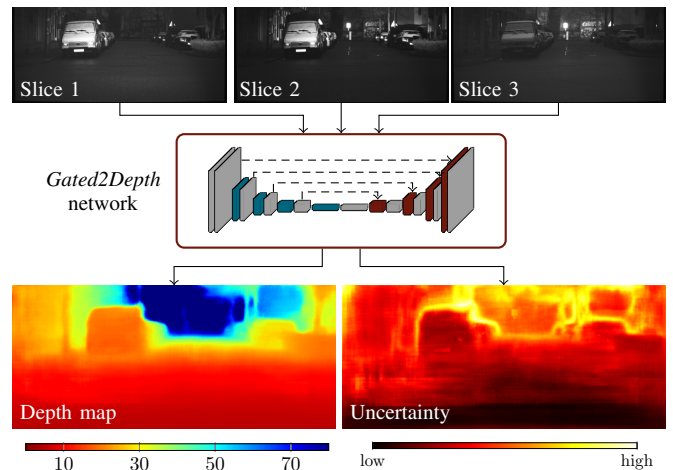


Fig. 1: We extend the recent *Gated2Depth* architecture [6] with an aleatoric uncertainty framework that converts three gated images into a high-resolution depth map and an uncertainty measure.

fusion algorithm to interpret and weight the incoming sensor measurements. As an alternative, information content can be classified by calculating the entropy [7].

Gated vision has been presented as an active imaging solution that significantly reduces backscatter and provides high-contrast images even in fog and rain [5]. Additionally, the scene geometry can be reconstructed at image resolution by processing at least two images with different delays [10]. As gated imaging is an active system, less illuminated areas with low signal-to-noise ratio (SNR) offer less accuracy of 3D sensing [11]. Generative models have been introduced in [6] in order to improve depth estimation in these areas by incorporating context information. Nevertheless, the confidence of the additionally generated information is not yet evaluated and is obviously not constant. For a sensor fusion algorithm, it would be beneficial to know for each pixel if the depth is measured, or generated by context and experience, and how certain the generative model is about its prediction.

In recent years, uncertainty estimation of neural networks based on Bayesian networks [12] has found its way into a wide field of applications such as scene segmentation [13] and object detection [14], [15]. In this work, we extend the generative model for gated depth estimation [6] with additional confidence estimation for better interpretable depth estimates. We show that the estimated depth confidence provides much higher quality than using the SNR of the input images. By filtering out a very small number of extreme outliers with very low confidence, we significantly improve the overall performance of 3D reconstruction.

<sup>1</sup> The authors are with Mercedes-Benz AG, Wilhelm-Runge-Str. 11, 89081 Ulm, Germany. E-mail: stefanie.walz@daimler.com, tobias.gruber@daimler.com, werner.r.ritter@daimler.com

<sup>2</sup> The authors are with the Institute of Measurement, Control and Microtechnology, Ulm University, Albert-Einstein-Allee 41, 89081 Ulm, Germany. E-mail: klaus.dietmayer@uni-ulm.de

## II. RELATED WORK

*a) 3D sensing:* 3D reconstruction from images is one of the fundamental challenges in computer vision. While structure from motion (SfM) [16], [17] exploits the motion of the camera to obtain multiple views for depth estimation by triangulation, stereo vision [8], [9] or multi-view methods [18] rely on a fixed and calibrated camera setup. Multi-view relies on finding correspondences and therefore suffer in texture-less regions and in case of occlusions. There is already work on the confidence of the stereo estimation, either directly in the stereo algorithm [8], [9] or as a post-processing step [19]. Static monocular setups instead leverage visual cues such as shading [20], perspective, or relative size [21], in recent years usually driven by deep neural networks [22], [23], [24], [25], [26], [27]. However, most of the current approaches do not offer a probabilistic interpretation of the depth estimation. Some first approaches for monocular depth prediction with confidence interpretation have been reported in [28], [29]. Gated imaging is related to other active systems such as light detection and ranging (LiDAR) [30] and time-of-flight (ToF) cameras [31]. It was shown in [32] that for ToF cameras it is insufficient to remove inaccurate estimates by simply thresholding low-amplitude values. Uncertainty can be either captured by a combination of distance, amplitude, their temporal and spatial variations, or learned by a regressor such as a random forest [32]. LiDAR depth completion combines a sparse LiDAR point cloud with RGB images in order to generate a high resolution depth map [33], [34], [35], [36]. A first approach where confidence maps are additionally learned during depth completion has been presented in [37].

*b) Uncertainty estimation:* Bayesian modeling offers the possibility to capture epistemic and aleatoric uncertainty in deep learning frameworks. While epistemic uncertainty represents uncertainty in the model parameters, aleatoric uncertainty depicts noise inherent in the observation. Utilizing Bayesian neural networks (BNNs), epistemic uncertainty can be modeled by inferring a posterior distribution over the model weights [38]. This is realized by replacing deterministic weights with stochastic weights following prior distributions. Generally, dropout variational inference is utilized to approximate BNNs [39]. That means that dropout is performed not only during training but also at test time in order to sample the posterior distribution. This approach was applied to camera localization estimation [40], semantic segmentation [41], and open-set object detection tasks [42]. Aleatoric uncertainty is modeled by placing a distribution over the output of the model. It is utilized to improve object detection [43] and to adjust weights of multi-task loss functions automatically [44]. Besides individual modeling of aleatoric and epistemic uncertainty, both uncertainties can be captured simultaneously. Kendall *et al.* [13] introduced a framework for classification and regression tasks, which can model either aleatoric or epistemic uncertainty alone or both together. This approach is further developed for two-stage [45] and one-stage [14], [15] object detection.

In this work, we model only aleatoric uncertainty since it can be extracted without time-consuming dropout sampling enabling our proposed framework to be run as a real-time application.

*c) Gated depth estimation:* Active gated imaging requires a sensitive image sensor and an illumination source. The synchronization of sensor gate and light source delivers the reflectivity of a scene in a certain range which enables the view through scattered environments. Heckman and Hodgson [46] were the first to take advantage of this visualization technique by using it to extend the range of visibility underwater. The active gated imaging technique offers two-dimensional images of the scene. To get a three-dimensional output, multiple gated images must be captured with different delays. Various methods have been developed for gated depth estimation, such as the time-slicing method. In this method, several two-dimensional gated images must be recorded in sequence where the gate delay is increased after each image. The resulting gate delay profiles are used to estimate the depth pixel-wise by threshold the rising or falling edge, determining the maximum, or by computing the weighted average of the profile [47]. Additionally, Andersson [48] resolved the depth by least-squares parameter fitting and data feature positioning. However, high-range accuracy requires small scanning step sizes resulting in a significant rise in capture time and processing effort. To overcome this problem, gain modulation and super-resolution depth mapping have been introduced. The gain modulation method exploits a pulsed laser and an intensified camera to recover the depth information with a gain-modulated and gain-constant image. Gain-modulated images are generated by linearly [49] or exponentially [50] increasing the gain of the intensifier during the gated time. This ensures independence of the laser pulse shape. The super-resolution depth mapping method exploits the knowledge of the range intensity profiles (RIPs) that are given by the convolution of the illuminating laser pulse and the sensor gate. To determine the depth, at least two gated images with distinct delay and overlapping RIPs are required. The use of trapezoid- [10] or triangular-shaped RIPs [51] enables range estimation by exploiting the strong linear dependencies of overlapping RIPs. Additionally, there is the opportunity to determine the depth by finding a transformation rule from pixel intensity to depth by fitting a 5th order ratio-polynomial model [52]. Optionally, the mapping between pixel intensity and depth can be learned by neural networks [53] or regression trees [54]. This enables the flexible adaption of gated settings to any scene to get an optimal image. The aforementioned methods estimate the depth pixel-wise, meaning the range of non-illuminated and saturated pixels cannot be determined correctly. Therefore, an image-based method has been developed that exploits a convolutional neural network to utilize semantic context across gated images [6]. However, the existing framework does not provide any measurement to confirm the exact estimation of non-illuminated and saturated pixels. For this reason, aleatoric uncertainty is introduced in this work to determine the correctness of the network's output.

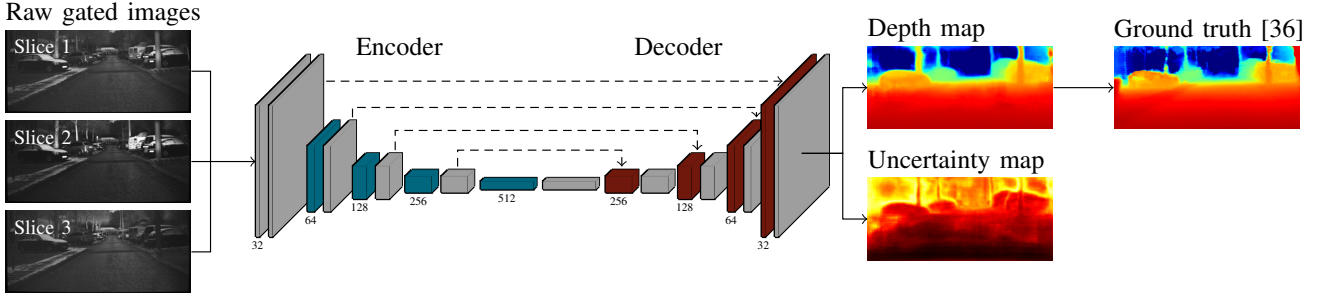


Fig. 2: We extend the *Gated2Depth* architecture [6] with an uncertainty output and train the network based on image-guided interpolated LiDAR ground truth [36].

### III. METHOD

#### A. Gated depth estimation

A gated viewing system consists of a diffused flash illuminator and a synchronized gated image sensor enabling high-contrast imaging in low-light, at night, and in adverse weather conditions [5]. Suppose that the laser pulse  $p_i(t)$  is reflected by a dominating lambertian reflector with albedo  $\alpha$  at distance  $r$  and the camera gate  $g_i(t)$  is delayed by  $\xi_i$ . Then, the RIP  $C_i(r)$  for each gated setting  $i \in \{1, 2, 3\}$  is defined by the correlation of  $p_i$  and  $g_i$ , namely

$$C_i(r) = \int_{-\infty}^{\infty} g_i(t - \xi_i) p_i\left(t - \frac{2r}{c}\right) \beta(r) dt, \quad (1)$$

where  $c$  and  $\beta(r)$  denote the speed of light and the distance-dependent atmospheric influence, respectively [47]. Note that the RIP is defined to be independent of the scene albedo  $\alpha$  and thus the final measured intensity  $z_{i,uv}$  on the image sensor at pixel position  $(u, v)$  is obtained by

$$z_{i,uv} = \alpha C_i(r_{uv}) + \eta_p(\alpha C_i(r_{uv})) + \eta_g = f_i(r_{uv}). \quad (2)$$

We follow the Poissonian-Gaussian noise model of Foi *et al.* [55] with Poissonian photon shot noise  $\eta_p$  and Gaussian read-out noise  $\eta_g$ .

The task of gated depth estimation is to recover the range  $r_{uv}$  from multiple gated measurements  $\mathbf{z}_{uv} = [z_{1,uv}, z_{2,uv}, z_{3,uv}]$  which basically means finding the inverse function  $f^{-1} : \mathbb{R}^3 \rightarrow \mathbb{R}$  that minimizes  $|\hat{r}_{uv} - r_{uv}|$  with

$$\hat{r} = f^{-1}(\mathbf{z}) = f^{-1}(\mathbf{f}(r_{uv})), \quad (3)$$

where the function  $\mathbf{f}(r_{uv}) = [f_1(r_{uv}), f_2(r_{uv}), f_3(r_{uv})]$  describes a vector of modulated noisy gated images that depends on the distance, see Eq. (2). In previous works, this inverse function  $f^{-1}$  has been learned with a fully-connected neural network [53] or a regression tree [54] and is applied pixel by pixel. However, pixel-based gated depth estimation fails in regions with low SNR, saturation, shadows, multipath, and blooming effects because no spatial correlation is exploited. *Gated2Depth* [6] is a fully convolutional encoder-decoder network that is able to generate full-resolution depth maps from gated images exploiting semantic context to fill these failure regions. In this work, we rely on the same network architecture from *Gated2Depth* as shown in Fig. 2 and extend this work with a novel uncertainty measure. The

network is a variant of the popular U-net [56] with skip-connections that consists of an encoder with four pairs of convolutions followed by a max pooling operation, and a decoder with four additional convolutions and transposed convolutions after each pair. More details can be found in [6].

#### B. Loss functions

The main goal of training is to penalize the absolute differences between ground truth depth  $r$  and its depth estimate  $\hat{r}$ . This can be achieved by the L1 loss

$$\mathcal{L}_{L1}(r, \hat{r}) = \frac{1}{N} \sum_{u,v} \|r_{uv} - \hat{r}_{uv}\|, \quad (4)$$

where  $(u, v)$  denotes the pixel position in the image and  $N = uv$  the number of pixels. Nevertheless, the main challenge of training full-image depth regression is the lack of dense ground truth data. In state-of-the-art automotive datasets such as KITTI, only sparse LiDAR measurements are available. Even accumulating consecutive LiDAR point clouds generates depth maps with only 16% coverage [57]. There exists a variety of approaches that tackle this problem by introducing specific loss functions during training that enforce dense depth output. Both multi-scale loss  $\mathcal{L}_{L1,m}$  and smoothness loss  $\mathcal{L}_s$  enforce the network to generate full-image depth either by upsampling the sparse ground truth for smaller variants of the output ( $\mathcal{L}_{L1,m}$ ) or by penalizing large depth differences between neighboring pixels ( $\mathcal{L}_s$ ). Note that the multi-scale loss  $\mathcal{L}_{L1,m}$  extends and therefore replaces  $\mathcal{L}_{L1}$ . Semantic understanding from synthetic datasets with dense ground truth depth can be transferred to the real-world training by adding an adversarial loss  $\mathcal{L}_{adv}$  based on a frozen synthetic discriminator that has been trained on synthetic dense ground truth [6]. This discriminator penalizes unrealistic looking depth maps. We exactly follow the formal definitions of  $\mathcal{L}_{L1,m}$ ,  $\mathcal{L}_s$ , and  $\mathcal{L}_{adv}$  as described in [6]. The final loss is given by

$$\mathcal{L} = \mathcal{L}_{L1,m} + \lambda_{adv} \mathcal{L}_{adv} + \lambda_s \mathcal{L}_s, \quad (5)$$

where  $\lambda_{adv}$  and  $\lambda_s$  denote tunable hyperparameters for weighting the loss components.

#### C. Learning from densified ground truth

Since LiDAR depth completion methods have shown impressive results in interpolating sparse depth measurements

guided by intensity images, we propose a novel training method for *Gated2Depth* that relies on this densified depth as ground truth. We apply the popular Sparse-to-Dense framework [36] on the LiDAR point clouds and RGB images of the training set. This certainly limits the performance of *Gated2Depth* to the performance of the depth completion approach. However, gated depth estimation aims to replace these expensive LiDAR systems with cost-sensitive hardware and intelligent post-processing, and achieving depth completion performance is already sufficient.

#### D. Introducing uncertainty

Bayesian modeling enables the extraction of epistemic and aleatoric uncertainty of neural networks. Epistemic uncertainty represents the uncertainty in the model parameters, which results from an insufficient amount of training data and vanishes for an infinite number of data. However, aleatoric uncertainty depicts observation noise of the input and remains stable independent of the input quantity. In [13], they showed that epistemic uncertainty can identify inputs that deviate from the training dataset, whereas aleatoric uncertainty is appropriate for real-time applications, as no expensive dropout sampling is required. Our research will be applied in cars, where fast sensing of the environment is essential. Hence, only aleatoric uncertainty is implemented into our framework.

To capture aleatoric uncertainty, the output  $r_{uv}$  of the neural network is modeled as a likelihood function. For the given depth regression task, we propose a likelihood function  $p(r_{uv})$  that follows a Laplacian distribution with mean  $\hat{r}_{uv}$  and variance  $\hat{\sigma}_{uv}$ :

$$p(r_{uv}) = \frac{1}{2\hat{\sigma}_{uv}} \exp\left[-\frac{\|r_{uv} - \hat{r}_{uv}\|}{\hat{\sigma}_{uv}}\right]. \quad (6)$$

Thereby,  $\hat{r}_{uv}$  is the estimated depth for an input  $\mathbf{z}_{uv}$  and  $\hat{\sigma}_{uv}$  represents the corresponding aleatoric uncertainty. Both depth  $\hat{r}_{uv}$  and uncertainty  $\hat{\sigma}_{uv}$  are modeled as explicit outputs of a single neural network, which is parameterized by its model weights. Instead of a Laplacian distribution, a Gaussian distribution can be utilized as a likelihood function. However, this has led to worse results in our regression problem and follows the experiences that a L1 loss is more effective for depth regression [58]. To find the model parameters that explain the model best, the likelihood function has to be maximized, which corresponds to the minimization of the negative log-likelihood, given by

$$-\log p(r_{uv}) = \frac{\|r_{uv} - \hat{r}_{uv}\|}{\hat{\sigma}_{uv}} + \log \hat{\sigma}_{uv} + \log 2. \quad (7)$$

Therefore, the negative log-likelihood averaged over each pixel  $(u, v)$  is considered as aleatoric loss function for training the neural network, namely

$$\mathcal{L}_{L1,aleatoric} = \frac{1}{N} \sum_{u,v} \|r_{uv} - \hat{r}_{uv}\| e^{-s_{uv}} + s_{uv}, \quad (8)$$

where  $N$  is the number of pixels of an output image and  $s_{uv} = \log \hat{\sigma}_{uv}$  is the log variance that is numerically more

stable as it avoids division by zero. Additionally, the constant term  $\log 2$  in Eq.(7) is neglected. When aleatoric uncertainty is applied on multiple scaled versions of the output, the multi-scale aleatoric loss  $\mathcal{L}_{L1,aleatoric,m}$  is given by

$$\mathcal{L}_{L1,aleatoric,m} = \sum_{i=0}^{M-1} \lambda_{m_i} \mathcal{L}_{L1,aleatoric}(r^{(i)}, \hat{r}^{(i)}, \hat{s}^{(i)}), \quad (9)$$

where  $r^{(i)}$ ,  $\hat{r}^{(i)}$ , and  $\hat{s}^{(i)}$  are scaled versions of  $r$ ,  $\hat{r}$ , and  $\hat{s}$ . To train and test models without aleatoric uncertainty,  $s_{uv}$  has to be set to zero and the aleatoric loss  $\mathcal{L}_{L1,aleatoric,m}$  converges to the multi-scale L1 loss  $\mathcal{L}_{L1,m}$ .

In conclusion, to introduce aleatoric uncertainty into the *Gated2Depth* architecture, we add an uncertainty map as second output and train the whole network with an aleatoric loss  $\mathcal{L}_{L1,aleatoric,m}$  that replaces the multi-scale loss  $\mathcal{L}_{L1,m}$ . Finally, the full loss is obtained by

$$\mathcal{L} = \mathcal{L}_{L1,aleatoric,m} + \lambda_{adv} \mathcal{L}_{adv} + \lambda_s \mathcal{L}_s. \quad (10)$$

#### E. Uncertainty filtering

Gated depth estimation provides depth information at pixel level. However, low SNR or saturated pixels can infer depth only from context which is not always possible. The uncertainty estimation provides a measure of how confident the model is about the depth estimation. To show the benefit of our proposed confidence measure, we introduce *uncertainty filtering* where depth estimates are filtered out when the uncertainty value is above a threshold  $t$ . This results in an overall better performance because obviously wrong estimates with high uncertainty are neglected. As a baseline, we use a simple SNR filtering based on the laser illumination, assuming that depth estimation from pixels with low illumination variance is insufficient. Given a set of gated input pixels  $\mathbf{z}_{uv} = [z_{1,uv}, z_{2,uv}, z_{3,uv}]$  and a predefined threshold  $\vartheta$ , pixels with low illumination variance are defined as the ones that satisfy  $\max(\mathbf{z}_{uv}) - \min(\mathbf{z}_{uv}) < \vartheta$ .

## IV. DATASETS

Since gated imaging is an emerging sensor technology, there are not many state-of-the-art datasets that provide gated images. Gruber *et al.* [6] have presented the first long-range gated dataset in real-world automotive scenarios that consists of 14,277 samples recorded in Northern Europe. They have equipped a vehicle with a Brightway Vision BrightEye gated camera with two vertical-cavity surface-emitting laser (VSCSEL) illuminators at 808 nm in the bumper. In addition to gated images ( $1920 \times 1024$ , 10 bit), stereo images ( $1920 \times 1024$ , 12 bit) and LiDAR point clouds (64 lines) from a Velodyne HDL64-S3 laser scanner are provided. For experiments with full-resolution ground truth depth, we rely on 9,804 simulated gated images based on the Grand Theft Auto V (GTA V) computer game [6]. We follow the same training, validation and test splits as in *Gated2Depth* [6] in order to make our contribution comparable.

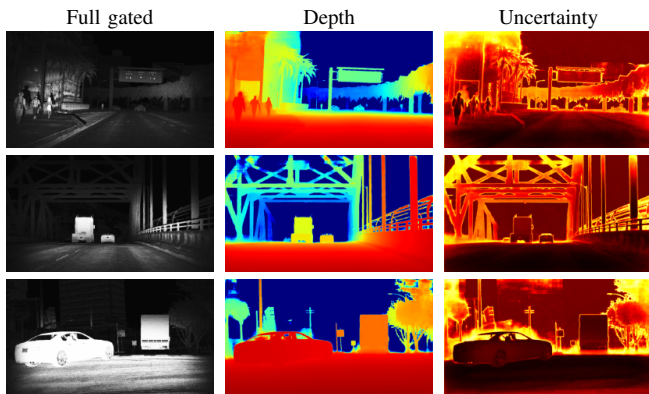


Fig. 3: Qualitative examples for the synthetic dataset. For color coding, we refer to Fig. 1.

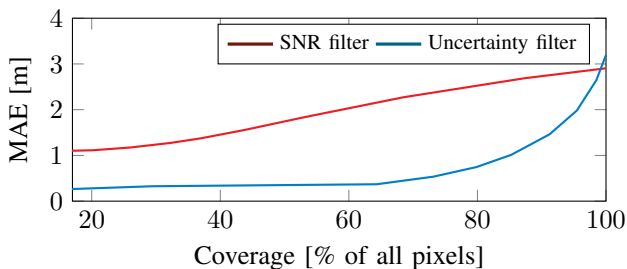


Fig. 4: Accuracy measured in MAE with respect to the depth coverage for *synthetic data*. By increasing the filter thresholds, unreliable depth estimations are filtered out resulting in a better overall performance.

## V. ASSESSMENT

### A. Experimental setup

Besides the implementation of uncertainty, our framework provides further changes compared to the original *Gated2Depth* network [6]. To enlarge the field-of-view of the input images, the original crop of 150 pixels at each side is replaced by a simple upper crop of 152 pixels, which was necessary due to missing LiDAR measurements in this area. Furthermore, instead of training two separate models for day and night, a single daytime-independent model is built, which simplifies the subsequent application in cars.

The presented method is implemented in *tensorflow* and trained with an *adam* optimizer and a learning rate of 0.0001. We have empirically chosen  $\lambda_s$  and  $\lambda_{adv}$  to 10. All models utilized here are trained for 15 epochs on a GeForce GTX TITAN Xp graphics processing unit (GPU), which took roughly 13 hours for the synthetic and 20 hours for the real dataset. According to the validation mean absolute error (MAE), we have selected the best performing epoch and hyperparameters. For the synthetic dataset the error is evaluated in a range from 3-150m and for the real dataset from 3-80m due to the limited range of the applied LiDAR system. We rely on the popular metrics root-mean-squared error (RMSE), MAE, scale invariant logarithmic error (SIlog) and the thresholds  $\delta_i < 1.25^i$  for  $i \in \{1, 2, 3\}$  as defined in [6].

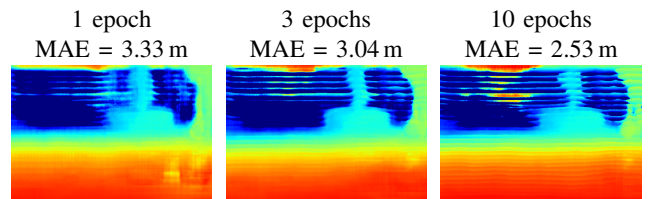


Fig. 5: This figure shows how training on sparse LiDAR ground truth depth generates horizontal patterns for an increasing number of epochs, although the MAE on the validation set decreases.

TABLE I: Ablation study for different loss functions and ground truth data types trained on real data.

Ground truth	Loss $\mathcal{L}$	MAE [m]		
		LiDAR	DC	LiDAR+DC
LiDAR	$\mathcal{L}_{L1}$	2.57	3.66	3.12
LiDAR	$\mathcal{L}_{L1,m}$	2.66	3.74	3.20
LiDAR	$\mathcal{L}_{L1} + \lambda_s \mathcal{L}_s$	2.98	3.31	3.15
LiDAR	$\mathcal{L}_{L1} + \lambda_{adv} \mathcal{L}_{adv}$	2.56	3.88	3.22
LiDAR	$\mathcal{L}_{L1,m} + \lambda_s \mathcal{L}_s + \lambda_{adv} \mathcal{L}_{adv}$	2.85	3.17	<b>3.01</b>
DC	$\mathcal{L}_{L1}$	3.08	2.64	<b>2.86</b>
DC	$\mathcal{L}_{L1,m} + \lambda_s \mathcal{L}_s + \lambda_{adv} \mathcal{L}_{adv}$	3.25	2.81	3.03

### B. Results on simulated data

For the synthetic dataset, the conventional and the aleatoric model are trained from scratch without adversarial and smoothness loss, because dense ground truth is available. Qualitative examples for the model with uncertainty are illustrated in Fig. 3. The bright areas in the uncertainty maps indicate high uncertainty and can be perceived especially at contours of objects and in non-illuminated areas of the image. Fig. 4 compares SNR and uncertainty filtering based on the estimated uncertainty. The performance of the aleatoric model can be significantly improved compared to the conventional one by filtering only a small number of pixels. To reduce the MAE by half, only about 10% of the pixels of the aleatoric model must be filtered whereas the conventional model requires filtering of about 60% to achieve such an error reduction. Note that the conventional model is slightly better than the one with uncertainty which explains the different starting points at 100% coverage.

### C. Ablation study for loss and ground truth

When training on real data with sparse LiDAR ground truth without any countermeasures, horizontal patterns in the estimated depth maps occur as Fig. 5 illustrates. For an increasing number of epochs, these horizontal patterns get even worse, although the MAE decreases due to evaluation on the sparse LiDAR points only. To quantify the horizontal patterns, we additionally evaluate on the depth completion (DC) ground truth. This dense ground truth is obtained by Sparse-to-Dense [36], a LiDAR depth completion method based on RGB images and sparse LiDAR points. While the error on the LiDAR ground truth reflects the accuracy of the depth estimation, the DC error measures the strength of the horizontal patterns.

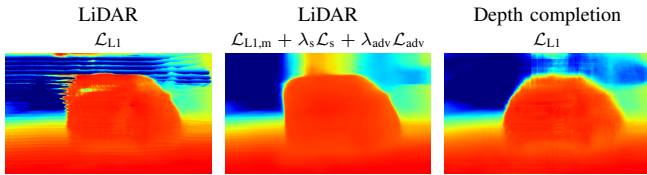


Fig. 6: Horizontal stripe patterns that occur when training on sparse ground truth LiDAR can be removed by additional loss components or extended ground truth annotations.

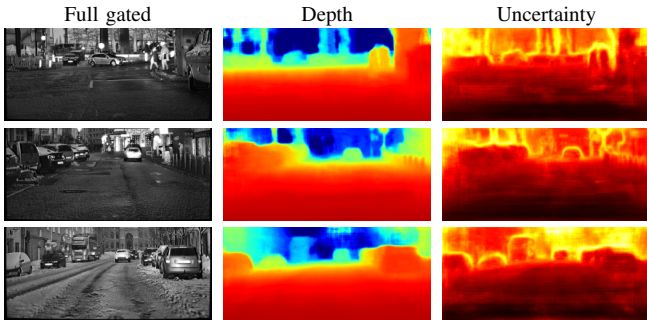


Fig. 7: Qualitative examples for the real dataset. In particular, high uncertainty arises at object edges and shadows above the objects.

Multi-scale loss, smooth loss, and adversarial loss are well-known approaches to handle sparse ground truth depth. Additionally, we propose a training on depth completed ground truth as another countermeasure. We evaluated different loss combinations for LiDAR and DC ground truth after training for 10 epochs with synthetic model initialization. To ensure accurate and smooth depth maps, the average of LiDAR and DC error is computed to compare the different approaches. The ablation study in Tab. I shows that applying a multi-scale loss with additional smoothness and adversarial loss delivers the best performance for models trained on sparse LiDAR ground truth. When trained on DC ground truth, a simple L1 loss generates the best result and additional loss components do not help. Fig. 6 illustrates the significant reduction of the horizontal patterns by using a multi-loss function (multi-scale, smooth, adversarial) or depth completed ground truth. According to the average of LiDAR and DC MAE metric, the model trained on DC ground truth shows a slightly better performance. Moreover, models with dense ground truth can be trained for longer, since no horizontal patterns are generated.

#### D. Results on real data

We train a model with uncertainty and one without uncertainty (baseline) to investigate how the incorporation of uncertainty changes the overall performance. Qualitative examples for our proposed uncertainty model in Fig. 7 exhibit high uncertainty for object contours and non-illuminated areas, e.g. due to shadows above each object. Tab. II indicates that the introduction of uncertainty comes at no additional costs as there is no loss in performance. The model with uncertainty has even a slightly better performance than the baseline model without uncertainty. The benefit of our proposed uncertainty measure is shown by comparing SNR

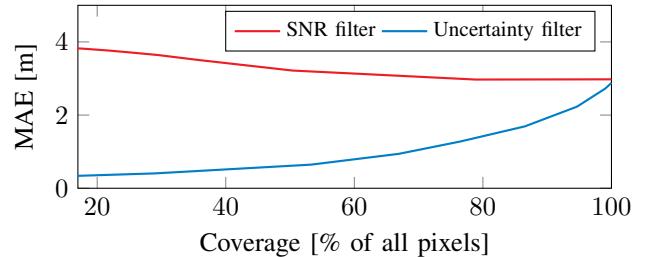


Fig. 8: Accuracy measured in MAE with respect to the depth coverage for real data. We create the points by increasing the filter threshold.

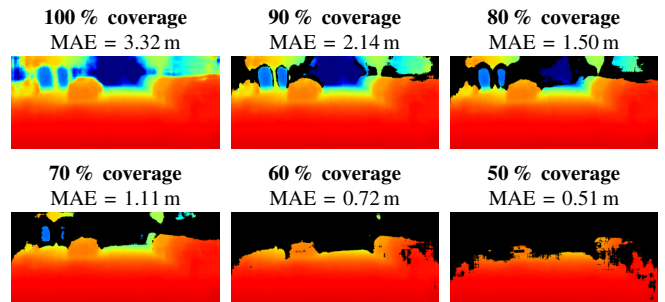


Fig. 9: Uncertainty filtering for varying coverage.

and uncertainty filtering. As Fig. 8 shows, SNR filtering does not help to get rid of erroneous measurements. On the contrary, SNR filtering probably removes pixels with good depth estimates and therefore decreases the overall performance. By applying the novel uncertainty filtering, the MAE can be significantly lowered. To halve the MAE, only about 20% of the pixels have to be filtered out. Fig. 9 demonstrates an exemplary filtering process for different pixel coverages. The MAE can be reduced from 3.32 m to less than 1 m at 50% coverage. However, when filtering out too many pixels, only the foreground is preserved and thus the semantic context of the image disappears. Hence, we decided that obtaining about 80% of the image pixels is a good choice since the MAE is reduced by half and individual objects can still be extracted.

#### E. Comparison with state-of-the-art methods

We follow [6] and compare our extended Gated2Depth method (*G2D+*) with state-of-the-art methods for 3D environment perception, such as monocular depth estimation [25], stereo vision [9], [60] and LiDAR depth completion [36]. *Monodepth* [25] and *PSMnet* [60] are finetuned on the real gated dataset. Finetuning of *Sparse-to-Dense* [36] is not possible because neither dense nor semi-dense ground truth depth is available. While Tab. II provides the evaluation metrics on the whole test dataset (day+night), the radar chart in Fig. 11 visualizes the results in a normalized representation. Each metric is normalized to the range [0, 1] such that the best approach is at 1 and the worst approach at 0.1. The results clearly show that traditional stereo [9] and monocular depth estimation [25] show worst performance, while LiDAR depth completion [36] provides the best results. However, LiDAR depth completion relies on sparse ground truth depth input from LiDAR and it is hard for other approaches to achieve this performance. Note that

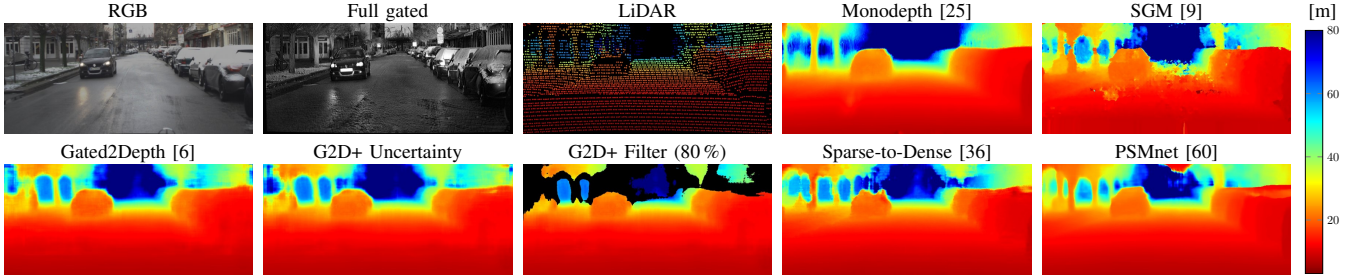


Fig. 10: Experimental daytime results for a variety of state-of-the-art methods.

TABLE II: Comparison of our proposed framework *G2D+* and state-of-the-art methods on the real test dataset according to common metrics as utilized in [59]. Note that Sparse-to-Dense requires ground truth input.

Method	RMSE [m]	MAE [m]	SIlog	$\delta_1$ [%]	$\delta_2$ [%]	$\delta_3$ [%]	Compl.
MONODEPTH [25]	9.59	4.70	25.35	82.32	92.23	95.68	1.00
PSMNET [60]	6.71	2.45	19.35	92.65	95.85	97.27	1.00
SGM [9]	10.83	5.26	37.82	77.38	86.36	90.52	1.00
SPARSE-TO-DENSE [36]	5.77	1.64	18.39	94.91	96.38	97.39	1.00
GATED2DEPTH [6]	7.08	2.98	21.78	89.86	94.91	96.76	1.00
G2D+ UNCERTAINTY	7.00	2.88	21.22	90.28	94.92	96.80	1.00
G2D+ FILTER (80 %)	3.84	1.49	13.78	95.53	98.06	98.75	0.87

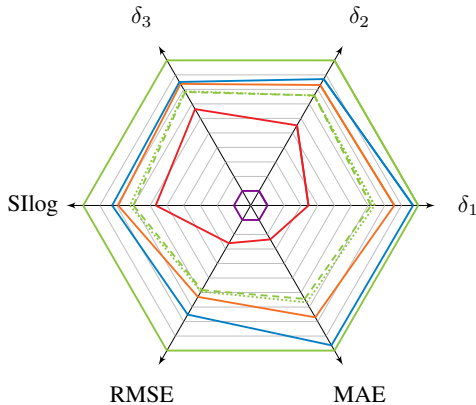
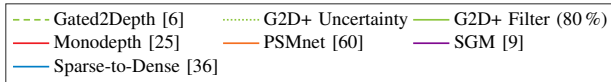


Fig. 11: Kiviatt diagram that visualizes the results from Table II. We normalize each metric such that the best approach is at 1 and the worst approach at 0.1.

compared to *Gated2Depth* in [6], we do not use separate models for day and night and we evaluate on a significantly larger image crop. Gated depth estimation without any filter shows in this setting only similar performance as deep stereo [60]. However, the additional uncertainty maps can be used to filter out unreliable depth estimates and *Gated2Depth* with uncertainty filter even outperforms depth completion performance.

## VI. CONCLUSIONS

This work extends the recent *Gated2Depth* method with aleatoric uncertainty that provides additional confidence information for each depth estimate. We show in the appli-

cation of uncertainty filtering, how this uncertainty measure can help to filter out unreliable depth estimates increasing the overall system performance. In an ablation study, we show that our proposed training on RGB guided interpolated ground truth depth is superior to conventional multi-loss approaches. Exciting future research includes the application of uncertainty maps into sensor fusion approaches, either for object detection or scene understanding enabling the integration of gated viewing systems into recent sensor setups of safe self-driving cars.

This work has received funding from the European Union under the H2020 ECSEL Programme as part of the DENSE project, contract number 692449.

## REFERENCES

- [1] N. Smolyanskiy, A. Kamenev, and S. Birchfield, “On the importance of stereo for accurate depth estimation: An efficient semi-supervised deep neural network approach,” in *Conf. on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 1007–1015.
- [2] E. R. Davies, *Machine vision: Theory, algorithms, practicalities*. Elsevier, 2004.
- [3] D. Chan, H. Buisman, C. Theobalt, and S. Thrun, “A noise-aware filter for real-time depth upsampling,” 2008.
- [4] M. Bijelic, T. Gruber, and W. Ritter, “A benchmark for LiDAR sensors in fog: Is detection breaking down?” in *IEEE Intelligent Vehicle Symposium*. IEEE, 2018, pp. 760–767.
- [5] —, “Benchmarking image sensors under adverse weather conditions for autonomous driving,” in *IEEE Intelligent Vehicle Symposium*. IEEE, 2018, pp. 1773–1779.
- [6] T. Gruber, F. Julca-Aguilar, M. Bijelic, and F. Heide, “Gated2depth: Real-time dense lidar from gated images,” *Proc. of the IEEE Int. Conf. on Computer Vision*, 2019.
- [7] M. Bijelic, T. Gruber, F. Mannan, F. Kraus, W. Ritter, K. Dietmayer, and F. Heide, “Seeing through fog without seeing fog: Deep multimodal sensor fusion in unseen adverse weather,” *arXiv preprint arXiv:1902.08913*, 2020.
- [8] D. Scharstein and R. Szeliski, “A taxonomy and evaluation of dense two-frame stereo correspondence algorithms,” *Int. Journal of Computer Vision*, vol. 47, no. 1-3, pp. 7–42, 2002.
- [9] H. Hirschmüller, “Accurate and efficient stereo processing by semi-global matching and mutual information,” in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, vol. 2. IEEE, 2005, pp. 807–814.
- [10] M. Laurenzis, F. Christnacher, and D. Monnin, “Long-range three-dimensional active imaging with superresolution depth mapping,” *Optics Letters*, vol. 32, no. 21, pp. 3146–3148, 2007.
- [11] B. Göhler and P. Lutzmann, “Range accuracy of a gated-viewing system as a function of the number of averaged images,” in *Electro-Optical Remote Sensing, Photonic Technologies, and Applications*, vol. 8542, 2012, pp. 37–44.
- [12] Y. Gal and Z. Ghahramani, “Bayesian convolutional neural networks with Bernoulli approximate variational inference,” *arXiv preprint arXiv:1506.02158*, 2015.
- [13] A. Kendall and Y. Gal, “What uncertainties do we need in Bayesian deep learning for computer vision?” in *Advances in Neural Information Processing Systems*, 2017, pp. 5574–5584.
- [14] F. Kraus and K. Dietmayer, “Uncertainty estimation in one-stage object detection,” in *IEEE Int. Conf. on Intelligent Transportation Systems*, 2019, pp. 53–60.

- [15] J. Choi, D. Chun, H. Kim, and H. Lee, "Gaussian YOLOv3: An accurate and fast object detector using localization uncertainty for autonomous driving," in *Proc. of the IEEE Int. Conf. on Computer Vision*, 2019.
- [16] R. Ranftl, V. Vineet, Q. Chen, and V. Koltun, "Dense monocular depth estimation in complex dynamic scenes," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2016, pp. 4058–4066.
- [17] B. Ummerhofer, H. Zhou, J. Uhrig, N. Mayer, E. Ilg, A. Dosovitskiy, and T. Brox, "DeMoN: Depth and motion network for learning monocular stereo," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017.
- [18] R. Hartley and A. Zisserman, *Multiple view geometry in computer vision*. Cambridge university press, 2003.
- [19] F. Tosi, M. Poggi, A. Benincasa, and S. Mattoccia, "Beyond local reasoning for stereo confidence estimation with deep learning," in *Proc. of the IEEE European Conf. on Computer Vision*, 2018, pp. 319–334.
- [20] R. J. Woodham, "Photometric method for determining surface orientation from multiple images," *Optical Engineering*, vol. 19, no. 1, p. 191139, 1980.
- [21] A. Saxena, S. H. Chung, and A. Y. Ng, "Learning depth from single monocular images," in *Advances in neural information processing systems*, 2006, pp. 1161–1168.
- [22] D. Eigen, C. Puhrsch, and R. Fergus, "Depth map prediction from a single image using a multi-scale deep network," in *Advances in Neural Information Processing Systems*, 2014, pp. 2366–2374.
- [23] I. Laina, C. Rupprecht, V. Belagiannis, F. Tombari, and N. Navab, "Deeper depth prediction with fully convolutional residual networks," in *Int. Conf. on 3D Vision (3DV)*, 2016, pp. 239–248.
- [24] R. Garg, B. V. Kumar, G. Carneiro, and I. Reid, "Unsupervised CNN for single view depth estimation: Geometry to the rescue," in *Proc. of the IEEE European Conf. on Computer Vision*, 2016, pp. 740–756.
- [25] C. Godard, O. Mac Aodha, and G. J. Brostow, "Unsupervised monocular depth estimation with left-right consistency," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017.
- [26] Y. Kuznetsov, J. Stückler, and B. Leibe, "Semi-supervised deep learning for monocular depth map prediction," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2017, pp. 2215–2223.
- [27] R. Garg, N. Wadhwa, S. Ansari, and J. T. Barron, "Learning single camera depth estimation using dual-pixels," *arXiv preprint arXiv:1904.05822*, 2019.
- [28] X. Yang, Y. Gao, H. Luo, C. Liao, and K.-T. Cheng, "Bayesian denet: Monocular depth prediction and frame-wise fusion with synchronized uncertainty," *IEEE Trans. on Multimedia*, 2019.
- [29] C. Liu, J. Gu, K. Kim, S. G. Narasimhan, and J. Kautz, "Neural RGB-D sensing: Depth and uncertainty from a video camera," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2019, pp. 10986–10995.
- [30] J. A. Christian and S. Cryan, "A survey of LiDAR technology and its use in spacecraft relative navigation," in *AIAA Guidance, Navigation, and Control (GNC) Conf.*, 2013, p. 4641.
- [31] R. Horaud, M. Hansard, G. Evangelidis, and C. M n n r, "An overview of depth cameras and range scanners based on time-of-flight technologies," *Machine vision and applications*, vol. 27, no. 7, pp. 1005–1020, 2016.
- [32] M. Reynolds, J. Doboř, L. Peel, T. Weyrich, and G. J. Brostow, "Capturing time-of-flight data with confidence," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*. IEEE, 2011, pp. 945–952.
- [33] Z. Chen, V. Badrinarayanan, G. Drozdov, and A. Rabinovich, "Estimating depth from RGB and sparse sensing," *arXiv preprint arXiv:1804.02771*, 2018.
- [34] M. Jaritz, R. De Charette, E. Wirbel, X. Perrotton, and F. Nashashibi, "Sparse and dense data with CNNs: Depth completion and semantic segmentation," in *Int. Conf. on 3D Vision (3DV)*, 2018, pp. 52–60.
- [35] F. Ma and S. Karaman, "Sparse-to-dense: Depth prediction from sparse depth samples and a single image," in *IEEE Int. Conf. on Robotics and Automation*, 2018, pp. 1–8.
- [36] F. Ma, G. V. Cavalheiro, and S. Karaman, "Self-supervised sparse-to-dense: Self-supervised depth completion from LiDAR and monocular camera," in *IEEE Int. Conf. on Robotics and Automation*, 2019.
- [37] W. Van Gansbeke, D. Neven, B. De Brabandere, and L. Van Gool, "Sparse and noisy LiDAR completion with RGB guidance and uncertainty," *arXiv preprint arXiv:1902.05356*, 2019.
- [38] D. J. MacKay, "A practical Bayesian framework for backpropagation networks," *Neural computation*, vol. 4, no. 3, pp. 448–472, 1992.
- [39] Y. Gal and Z. Ghahramani, "Dropout as a Bayesian approximation: Representing model uncertainty in deep learning," in *Proc. of the Int. Conf. on Machine Learning*, 2016, pp. 1050–1059.
- [40] A. Kendall and R. Cipolla, "Modelling uncertainty in deep learning for camera relocalization," in *IEEE Int. Conf. on Robotics and Automation*, 2016, pp. 4762–4769.
- [41] A. Kendall, V. Badrinarayanan, and R. Cipolla, "Bayesian segnet: Model uncertainty in deep convolutional encoder-decoder architectures for scene understanding," *arXiv preprint arXiv:1511.02680*, 2015.
- [42] D. Miller, L. Nicholson, F. Dayoub, and N. S nderhauf, "Dropout sampling for robust object detection in open-set conditions," in *IEEE Int. Conf. on Robotics and Automation*, 2018, pp. 3243–3249.
- [43] D. Feng, L. Rosenbaum, F. Timm, and K. Dietmayer, "Leveraging heteroscedastic aleatoric uncertainties for robust real-time LiDAR 3D object detection," in *IEEE Intelligent Vehicle Symposium*, 2019, pp. 1280–1287.
- [44] A. Kendall, Y. Gal, and R. Cipolla, "Multi-task learning using uncertainty to weigh losses for scene geometry and semantics," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 7482–7491.
- [45] D. Feng, L. Rosenbaum, and K. Dietmayer, "Towards safe autonomous driving: Capture uncertainty in the deep neural network for LiDAR 3D vehicle detection," in *IEEE Int. Conf. on Intelligent Transportation Systems*, 2018, pp. 3266–3273.
- [46] P. Heckman and R. T. Hodgson, "2.7—underwater optical range gating," *IEEE Journal of Quantum Electronics*, vol. 3, no. 11, pp. 445–448, 1967.
- [47] J. Busck and H. Heiselberg, "High accuracy 3D laser radar," in *Laser Radar Technology and Applications IX*, vol. 5412, 2004, pp. 257–263.
- [48] P. Andersson, "Long-range three-dimensional imaging using range-gated laser radar images," *Optical Engineering*, vol. 45, no. 3, pp. 1–10, 2006.
- [49] Z. Xiuda, Y. Huimin, and J. Yanbing, "Pulse-shape-free method for long-range three-dimensional active imaging with high linear accuracy," *Optics letters*, vol. 33, no. 11, pp. 1219–1221, 2008.
- [50] C. Jin, X. Sun, Y. Zhao, Y. Zhang, and L. Liu, "Gain-modulated three-dimensional active imaging with depth-independent depth accuracy," *Optics Letters*, vol. 34, no. 22, pp. 3550–3552, 2009.
- [51] W. Xinwei, L. Youfu, and Z. Yan, "Triangular-range-intensity profile spatial-correlation method for 3D super-resolution range-gated imaging," *Applied Optics*, vol. 52, no. 30, pp. 7399–406, 2013.
- [52] M. Laurenzis, F. Christnacher, N. Metzger, E. Bacher, and I. Zielenski, "Three-dimensional range-gated imaging at infrared wavelengths with super-resolution depth mapping," in *SPIE Infrared Technology and Applications XXXV*, vol. 7298, 2009.
- [53] T. Gruber, M. Kokhova, W. Ritter, N. Haala, and K. Dietmayer, "Learning super-resolved depth from active gated imaging," in *IEEE Int. Conf. on Intelligent Transportation Systems*. IEEE, 2018, pp. 3051–3058.
- [54] A. Adam, C. Dann, O. Yair, S. Mazor, and S. Nowozin, "Bayesian time-of-flight for realtime shape, illumination and albedo," *IEEE Trans. on Pattern Analysis and Machine Intelligence*, vol. 39, no. 5, pp. 851–864, 2017.
- [55] A. Foi, M. Trimeche, V. Katkovnik, and K. Egiazarian, "Practical poissonian-gaussian noise modeling and fitting for single-image raw-data," *IEEE Trans. on Image Processing*, vol. 17, no. 10, pp. 1737–1754, 2008.
- [56] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *Int. Conf. on Medical image computing and computer-assisted intervention*, 2015, pp. 234–241.
- [57] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, "Sparsity invariant CNNs," in *Int. Conf. on 3D Vision (3DV)*, 2017, pp. 11–20.
- [58] M. Carvalho, B. Le Saux, P. Trouv -Peloux, A. Almansa, and F. Champagnat, "On regression losses for deep depth estimation," in *IEEE Int. Conf. on Image Processing*, 2018, pp. 2915–2919.
- [59] T. Gruber, M. Bijelic, F. Heide, W. Ritter, and K. Dietmayer, "Pixel-accurate depth evaluation in realistic driving scenarios," in *Int. Conf. on 3D Vision (3DV)*, 2019, pp. 95–105.
- [60] J.-R. Chang and Y.-S. Chen, "Pyramid stereo matching network," in *Proc. of the IEEE Conf. on Computer Vision and Pattern Recognition*, 2018, pp. 5410–5418.