

SVDistNet: Self-Supervised Near-Field Distance Estimation on Surround View Fisheye Cameras

Varun Ravi Kumar, Marvin Klingner, Senthil Yogamani, Markus Bach,
Stefan Milz, Tim Fingscheidt, and Patrick Mäder

Abstract—A 360° perception of scene geometry is essential for automated driving, notably for parking and urban driving scenarios. Typically, it is achieved using surround-view fisheye cameras, focusing on the near-field area around the vehicle. The majority of current depth estimation approaches focus on employing just a single camera, which cannot be straightforwardly generalized to multiple cameras. The depth estimation model must be tested on a variety of cameras equipped to millions of cars with varying camera geometries. Even within a single car, intrinsics vary due to manufacturing tolerances. Deep learning models are sensitive to these changes, and it is practically infeasible to train and test on each camera variant. As a result, we present novel camera-geometry adaptive multi-scale convolutions which utilize the camera parameters as a conditional input, enabling the model to generalize to previously unseen fisheye cameras. Additionally, we improve the distance estimation by pairwise and patchwise vector-based self-attention encoder networks. We evaluate our approach on the Fisheye WoodScape surround-view dataset, significantly improving over previous approaches. We also show a generalization of our approach across different camera viewing angles and perform extensive experiments to support our contributions. To enable comparison with other approaches, we evaluate the front camera data on the KITTI dataset (pinhole camera images) and achieve state-of-the-art performance among self-supervised monocular methods. An overview video with qualitative results is provided at <https://youtu.be/bmX0UcU9wtA>. Baseline code and dataset will be made public¹.

Index Terms—Depth estimation, semantic segmentation, fisheye cameras, surround-view, multi-task learning, neural networks, self-supervised learning.

I. INTRODUCTION

Surround View Cameras Typically, automotive perception systems use multiple cameras, with current systems having at least four cameras. The number is likely to increase to more than ten cameras for future generation systems. Such surround-view cameras are focused on near field sensing, which is typically used for low-speed applications such as parking or traffic jam assistance functions [1]. Fig. 1 shows sample distance estimation images of four cameras mounted on a car covering the entire 360° field of view surrounding the car. Near-field distance estimation is a challenging problem because of distortion and partial visibility of close-by objects [2]. Also, centimeter-level accuracy is required to enable precise



Fig. 1: Our surround-view distance estimation framework is facilitated by employing a single network on images from multiple cameras. A surround-view coverage of geometric information is obtained for an autonomous vehicle by utilizing and post-processing the distance maps from all cameras.

low-speed maneuvers such as parking. Up to now, for the generation of high-quality distance estimates, one network per camera has to be trained, inducing unfeasible computational complexity with an increasing number of cameras.

Motivation of SVDistNet Our work's primary motivation is to propose a solution for surround-view fisheye cameras targeting large-scale industrial deployment. In other words, we aim to develop a model that is applicable to a vast number of automated driving vehicles possessing their own set of cameras. However, the underlying camera intrinsics are not even entirely identical for a given vehicle family. There are differences due to manufacturing processes that lead to a separate calibration of each camera instance. Calibration can vary even after deployment due to high environmental temperatures or aging. As a result, we require a calibration adaptation mechanism in the model. This contrasts with public datasets, which employ a single camera instance for both training and testing. There are 12 different cameras with slight intrinsic variations in the WoodScape dataset [3] to assess this effect. There are 4 camera instances around the vehicle with

V. Ravi Kumar, M. Bach are with Valeo DAR Kronach, Kronach, Germany.
V. Ravi Kumar, S. Milz and P. Mäder are with TU Ilmenau, Germany.
M. Klingner and T. Fingscheidt are with TU Braunschweig, Germany.
S. Yogamani is with Valeo Vision Systems, Tuam, Ireland.
Manuscript received September 13, 2020

¹<https://github.com/valeoai/WoodScape>

varying intrinsics, even for a single instance of a surround-view system. Instead of having four separate models for each camera, we, therefore, propose a single model for all cameras, which has advantages such as (i) a lower memory footprint and a better data transmission rate leading to improved efficiency, (ii) a better regularization during training due to the wider variety of different cameras in a more extensive training dataset, and (iii) an easier to handle maintenance of just a single neural network model instead of multiple ones.

Recently proposed methods for self-supervised depth estimation [4], [5], [6] present a powerful tool for the perception of scene geometry from camera images as they are trainable on arbitrary image sequences. While initial methods mainly focused on pinhole cameras [7], [5], [8], for practical deployment, the more general distance estimation [9], [10] on other camera geometries such as fisheye cameras are of significant interest due to their large field of view. There are many benefits in using raw images instead of undistortion, refer to [10] for details. Moreover, current methods essentially focus on single front-view camera systems, while the question of *how to extend such models to multiple cameras is essentially left open*. Therefore, in this work, we extend our previous synergized multi-task learning method for self-supervised distance estimation and semantic segmentation [11] to a general application on surround-view cameras and propose additional improvements.

To sum up, our contributions are as follows:

- We present a novel camera geometry adaptive multi-scale convolution to incorporate the camera parameters into our self-supervised distance estimation framework. We feed this camera geometry tensor representation to the model as a generic way to adapt to new camera intrinsics.
- We create a training framework for self-supervised distance estimation, which jointly trains and infers images from multiple fisheye cameras and viewpoints.
- We demonstrate a single model for 12 fisheye cameras, which achieves the equivalent result as an individual specialized model that overfits a particular camera model.
- We present an improved version of our network architecture for multi-task learning of self-supervised distance estimation and semantic segmentation. We significantly improve upon our previous works FisheyeDistanceNet [9], UnrectDepthNet [10], and SynDistNet [11].
- We achieve state-of-the-art results on the WoodScape and KITTI datasets among monocular self-supervised depth estimation methods.

II. RELATED WORK

This section provides a survey regarding self-supervised depth estimation and, in particular, its applications to different camera geometries. Also, we discuss the application of self-supervised depth estimation in multi-task learning approaches and approaches using semantic guidance for this task.

Self-Supervised Depth Estimation Training a neural network to predict a depth map by optimizing this depth map as a parameter of geometric projections between stereo images or sequential images has been defined as the task of self-supervised depth estimation [4]. For performance gains, the

following approaches extended the fundamental concept by improved loss functions [12], [5], or network architectures, optimized for the task of depth estimation [13]. Concurrently, it was proposed to use proxy labels from traditional stereo algorithms [14]. To better model the structural information from the temporal context, recurrent neural networks were employed [15], [16], and Casser *et al.* [7], [17] introduced a test-time refinement method to improve the network performance in an online manner. While all of these approaches achieve significant results on pinhole camera images, they are seldom transferred to more complex camera models, *e.g.*, fisheye camera images. In contrast, we exhibit our self-supervised distance estimation method’s applicability on both pinhole and fisheye camera images.

Multi-Task Learning Exemplary tasks, which benefit from multi-task learning, *i.e.*, share some of the network parts, in terms of performance and efficiency are, *e.g.*, depth estimation [18], [19], semantic segmentation [20], [21], [19] and domain adaptation [22]. As an empirical weighting of the losses involves extensive hyperparameter tuning, we follow the uncertainty-based weighting of Kendall *et al.* [18] to weigh the different tasks. To handle dynamic objects during training that violate the static world assumption in the projections of the self-supervised depth estimation, it is common to additionally estimate optical flow, which can also be trained in a self-supervised fashion [23]. Through this simultaneous estimation of both outputs, it is thereby plausible to enforce cross-task consistency [24], to modify the projected image and correct the influence of dynamic objects [25], or to enforce prior known geometric constraints. Although optical flow is the commonly used technique to handle dynamic objects, we consider an alternative method, introduced by [8], [11] because, firstly, semantic segmentation is a task that is readily available in most autonomous driving systems and, secondly, optical flow is computationally more complex and hard to validate due to the extensive effort needed to obtain ground truth.

Semantically-Guided Depth Estimation Complementing self-supervised depth estimation by the additional prediction of semantic or instance segmentation in cross-task guidance approaches has been shown to increase the performance of both prediction modalities [7], [17], [8], [26], [27], [28]. Either the segmentation masks are given as an additional input to the network [28], or they are used to predict relative poses for the single dynamic objects to counterpoise their movement between consecutive frames [7], [17]. While primary approaches used a pre-trained segmentation network, it was also shown that both tasks could be trained simultaneously [29], [30]. Thereby, it is possible to project the segmentation outputs between consecutive frames and enforce temporal semantic consistency [29], [31], or enforce edge consistency across the different output modalities [29], [30]. In this work, we use the segmentation output to identify dynamic objects and handle them accordingly inside the photometric loss as in [8], [11]. Additionally, we propose a multi-task network, where the encoder is based on pairwise and patchwise self-attention networks together with pixel-adaptive convolutions inside the decoder (as in [28]), which can be trained in a single-stage manner, eliminating the need for a pre-trained segmentation

network compared to [28].

Depth Estimation for Different Camera Geometries Employing self-supervised depth estimation techniques to different camera geometries than a single pinhole camera is challenging for two reasons. Firstly, the intrinsic camera parameters are assumed to be constant for all images. Secondly, the viewing angle is assumed to be constant for all camera images, as usually only a front-facing camera is considered (e.g., [4], [5]). Gordon *et al.* [32] predicts not only the relative pose but also the intrinsic camera parameters, enabling the generalization across images from different pinhole cameras. Meanwhile, Facil *et al.* [33] introduces a method to give the intrinsic camera parameters an additional input to the convolution operation, showing the same effect. In this work, we transfer the approach [33] to fisheye camera images, thereby enabling generalization of our distance estimation method across different fisheye cameras and viewing angles.

The transfer of self-supervised depth estimation to the more general self-supervised distance estimation on camera geometries, such as, e.g., fisheye cameras, was introduced by Ravi Kumar *et al.* [9], [10], which focus on front-view cameras. Surround-view depth coverage has also been shown by [34], [35], which applies self-supervised depth estimation techniques to 360° images. In this work, we do not have to impose a 360° image in advance to the model; we enable a surround-view coverage of geometric information for an autonomous vehicle by applying the same neural network to many arbitrary cameras with different viewing angles.

III. MULTI-TASK LEARNING FRAMEWORK

This section describes our multi-task learning framework for collaborative distance and semantic segmentation prediction, which we apply to surround-view camera systems. The learning framework is based on our previous work FisheyeDistanceNet [9], UnrectDepthNet [10], and SynDistNet [11] with extensions to support multiple cameras. To enable comparison, we use the same protocols and losses.

A. Self-Supervised Distance and Pose Estimation

We develop our single-image distance and pose estimation framework for multiple cameras building on our previous work FisheyeDistanceNet [9], where we outlined the basic mechanism of our polynomial projection model for fisheye cameras. Accordingly, the total loss for our distance and pose estimation networks (cf. Fig. 2) consists of a image reconstruction contribution \mathcal{L}_r , minimizing the difference between the reconstructed images $\hat{I}_{t' \rightarrow t}$ and the target frames I_t and a regularizing contribution \mathcal{L}_s , enforcing edge-aware smoothness as outlined in [12]. Additionally, we make use of the cross-sequence distance consistency loss \mathcal{L}_{dc} as well as the scale recovery technique from [9]. The total loss is averaged over image batches, scales, and pixel locations and is given as

$$\mathcal{L}_{tot} = \mathcal{L}_r(I_t, \hat{I}_{t' \rightarrow t}) + \beta \mathcal{L}_s(\hat{D}_t) + \gamma \mathcal{L}_{dc}(\hat{D}_t, \hat{D}_{t'}). \quad (1)$$

with β and γ being constants, weighing the influence of the regularization term \mathcal{L}_s and the cross-sequence distance consistency term \mathcal{L}_{dc} , respectively.

Image Reconstruction Loss The majority of cutting-edge self-supervised depth estimation approaches use heuristic loss functions. The optimal selection of a loss function, on the other hand, is not well defined theoretically. We emphasized the importance of investigating a better reconstruction loss in our previous paper, SynDistNet [11]. We incorporated Barron’s [36] more generic robust loss function, which was used to replace the L_1 term employed in [9], [10], [12], [5]. In the context of distance/depth estimation, we presented the common notion of a per-pixel regression ρ , indicated by

$$\rho(\xi) = \rho\left(I_t - \hat{I}_{t' \rightarrow t}\right) \quad (2)$$

This robust loss function generalizes several common losses, including the L_1 , L_2 , Geman-McClure, Welsch/Leclerc, Cauchy/Lorentzian, and Charbonnier loss functions. Robustness is introduced as a continuous parameter that can be optimized within the loss function to increase regression task performance. The robust loss function ρ_{rob} is as follows:

$$\rho_{\text{rob}}(\xi) = \frac{|\alpha - 2|}{\alpha} \left(\left(\frac{(\xi/c)^2}{|\alpha - 2|} + 1 \right)^{\alpha/2} - 1 \right) \quad (3)$$

As described in [36], one can use a data-driven optimization to automatically adapt the free parameters α and c to the investigated problem. The image reconstruction loss between the target frame I_t and the reconstructed target frame $\hat{I}_{t' \rightarrow t}$ is computed by a mixture of the Structural Similarity loss term and the robust pixel-wise loss term as

$$\rho_{\text{rob}}(\xi) = \rho \left\| (I_t - \hat{I}_{t' \rightarrow t}) \odot \mathcal{M}_{t \rightarrow t'} \right\| \quad (4)$$

$$\tilde{\mathcal{L}}_r(I_t, \hat{I}_{t' \rightarrow t}) = \tau \frac{1 - \text{SSIM}(I_t, \hat{I}_{t' \rightarrow t}, \mathcal{M}_{t \rightarrow t'})}{2} + (1 - \tau) \rho_{\text{rob}}(\xi) \quad (5)$$

where $\tau = 0.85$ is the weighting factor between both loss terms and $\mathcal{M}_{t \rightarrow t'}$ is the binary ego-mask employed from [9]. Finally, we apply the per-pixel minimum reconstruction loss \mathcal{L}_r [5] by a computation over all source images.

$$\mathcal{L}_r = \min_{t' \in \{t+1, t-1\}} \tilde{\mathcal{L}}_r(I_t, \hat{I}_{t' \rightarrow t}) \quad (6)$$

B. Semantic Guidance for Image Reconstruction Loss

To improve distance estimation performance through handling dynamic objects accordingly, in this part, we describe how we train our semantic segmentation baseline and afterward describe the semantic masking technique incorporated from our previous work [11] to exclude their influence on the reconstruction loss.

Semantic Segmentation Baseline: Semantic segmentation aims at assigning a class label s from a subset of classes $s \in \mathcal{S} = \{1, 2, \dots, S\}$. The desired output is then a pixel-wise semantic segmentation mask M_t with the same spatial dimensions as the segmentation network’s input I_t . During training, the network predicts posterior probabilities Y_t , representing the likelihood that a pixel belongs to a class $s \in \mathcal{S}$.

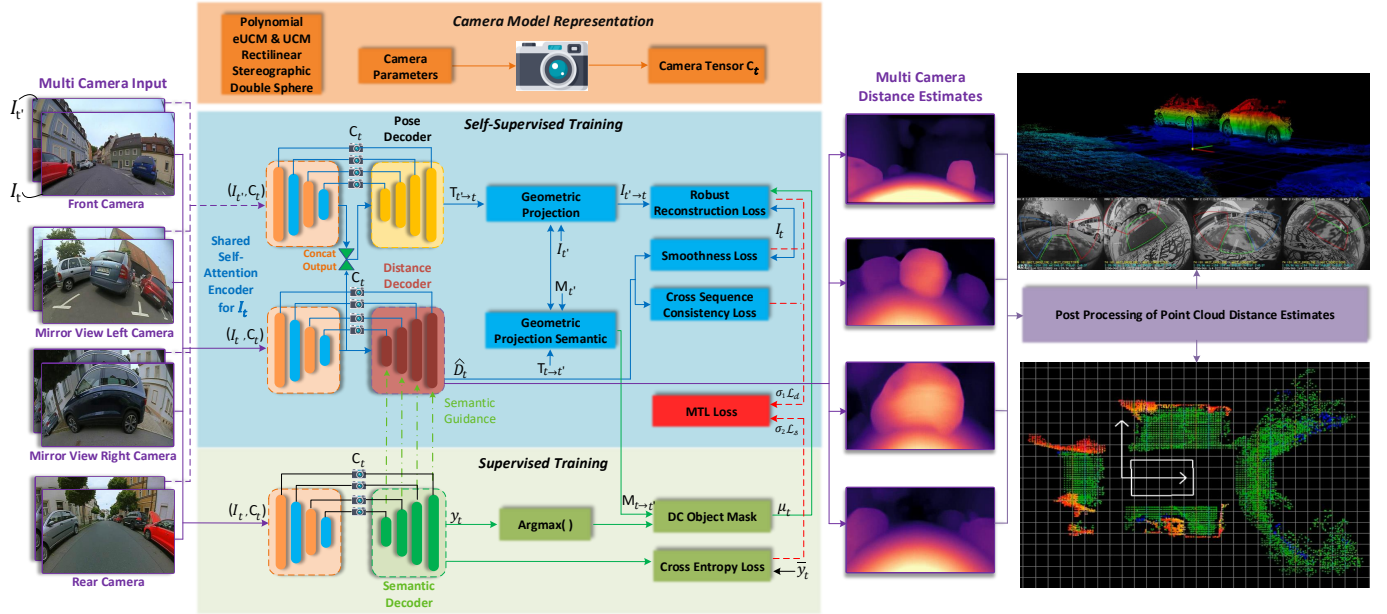


Fig. 2: **High-level overview of our surround-view self-supervised distance estimation framework**, which employs semantic guidance as well as camera-geometry adaptive convolutions (orange blocks). Our framework comprises training units for self-supervised distance estimation (blue blocks) and semantic segmentation (green blocks). The camera tensor C_t (orange block) assists our SVDistNet in producing distance maps across multiple camera-viewpoints and making the network camera independent. As explained in Section IV-A C_t can also be applied to standard camera models. The multi-task loss from 9 weights and optimizes both modalities at the same time. By post-processing the predicted distance maps in 3D space, we can obtain surround-view geometric information using our proposed framework.

A cross-entropy loss is then used to optimize this posterior probability in a supervised manner.

$$\mathcal{L}_{ce} = - \sum_{s \in S} \bar{Y}_{t,s} \cdot \log(Y_{t,s}) \quad (7)$$

which is averaged per pixel and where \bar{Y}_t represents one-hot encoded ground truth labels. The final segmentation mask M_t is obtained by performing a pixel-wise argmax operation on the posterior probabilities $Y_{t,s}$.

Dealing With Dynamic Objects: Because dynamic objects defy the static world assumption, knowledge of their depth/distance is critical in autonomous driving. Therefore, we use the information provided by the semantic segmentation to learn the distance from *non-moving* dynamic objects, while we exclude potentially *moving* dynamic objects. To this end, we employ a pixel-wise mask μ_t containing a 1, if a pixel does neither show a dynamic object from the current frame I_t nor a wrongfully projected dynamic object from the reconstructed frames $\hat{I}_{t' \rightarrow t}$. Otherwise, we set the value of the mask to 0. Accordingly, we have to predict a segmentation mask M_t from the target frame I_t , as well as segmentation masks $M_{t'}$ for the source frames $I_{t'}$. While we can easily detect the dynamic objects in the source frame's segmentation mask $M_{t'}$, the wrongfully projected dynamic objects can be obtained by warping the segmentation masks and employing a nearest-neighbor sampling to the target frame. This yields projected segmentation masks $M_{t' \rightarrow t}$. Finally, we can define the set of dynamic object classes $S_{DC} \subset S$ and reduce the semantic segmentation mask to a binary mask, with the properties as outlined above by computing its elements at location ij as:

$$\mu_{t,ij} = \begin{cases} 1, & M_{t,ij} \notin S_{DC} \wedge M_{t' \rightarrow t,ij} \notin S_{DC} \\ 0, & \text{else} \end{cases} \quad (8)$$

We mask dynamic objects by multiplying the mask with the reconstruction loss from Eq. 4 in a pixel-wise fashion. However, there are images in which these dynamic objects are not moving, *e.g.*, parking cars are perfect to learn the distance of cars, while cars on a highway severely violate the static world assumption. Accordingly, we want to learn the distance of *moving* dynamic objects by learning from images, in which they are detected as *non-moving*. As a result, we detect images with non-moving dynamic objects using a semantic consistency measure introduced in our previous works [8], [11] and only mask the dynamic objects in a fraction *epsilon* of images where dynamic objects are detected as mostly moving.

C. Joint Optimization

Based on Kendall's *et al.* [18] task weighting strategy, the distance estimation and semantic segmentation loss terms from Eq. 1 and Eq. 7, respectively, are weighted. The total optimization objective is derived from [11].

$$\frac{1}{2\sigma_1^2} \mathcal{L}_{tot} + \frac{1}{2\sigma_2^2} \mathcal{L}_{ce} + \log(1 + \sigma_1) + \log(1 + \sigma_2) \quad (9)$$

where σ_1 and σ_2 are learnable parameters introduced to optimize the homoscedastic uncertainty. The goal is to optimize a more substantial uncertainty, resulting in a smaller contribution of the loss from the distance estimation or semantic segmentation task to the total loss. It is important to note that the respective task's weight is reduced by increasing the noise parameter σ . We observe that the noise parameter σ_1 of the distance estimation task is lower than σ_2 of semantic segmentation, and over time, the convergence occurs correspondingly.

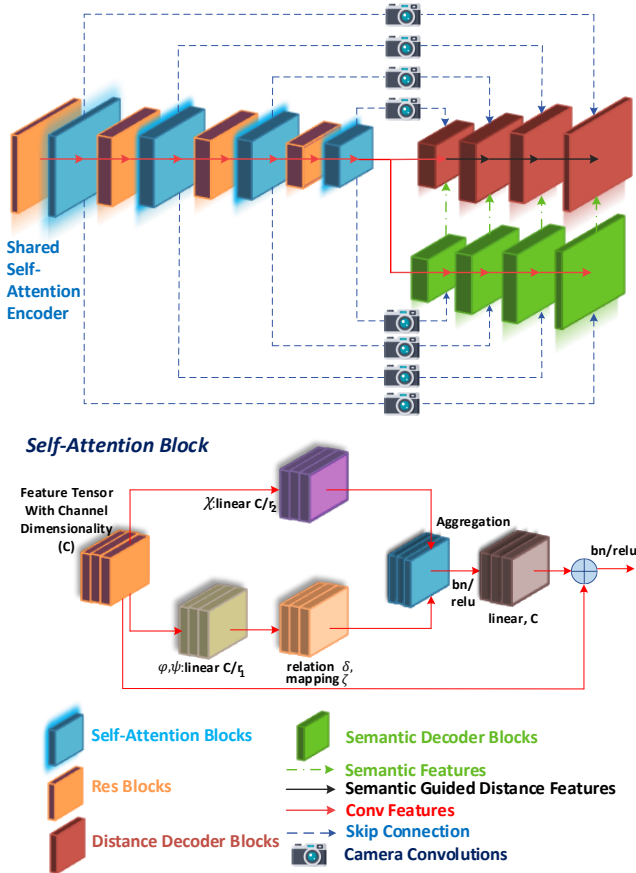


Fig. 3: **Overview of our proposed network architecture** for semantically guided self-supervised distance estimation. It consists of a shared vector-based self-attention encoder and task-specific decoders. Our encoder is a self-attention network with pairwise and patchwise variants, while the decoder uses pixel-adaptive convolutions, which are both complemented by our novel Camera Geometry convolutions.

IV. NETWORK ARCHITECTURE

This section outlines our novel architecture for semantically-guided self-supervised distance estimation utilizing the camera geometry tensor to handle varying camera intrinsics and viewpoints. We also incorporate pairwise and patchwise self-attention encoders from [37] which improves on our baseline SynDistNet [11], which employs a scalar self-attention encoder. Thereby, an adaptation of the weights across spatial and channel dimensions is implemented. We retain our semantic guidance decoder from [11] to facilitate the distance decoder to yield accurate distance predictions. The complete training of both tasks is performed in a one-stage manner.

A. Camera Geometry Tensor C_t

Automated driving systems have a wide variety of cameras, typically around 10, placed in different car locations and with different fields of view. Instead of developing an individual model for each camera, it is highly desirable to develop a single model for all cameras, as discussed in the introduction. This is an unsolved problem, and we aim to solve this by incorporating camera geometry into the distance estimation. We intend to convert all the camera geometry characteristics into a tensor called camera geometry tensor C_t and provide

this as input to the CNN model at both training and inference. From the view of distance estimation, camera intrinsics is the primary model adaptation needed. However, the camera tensor notion is generic, and we plan to extend it to include camera extrinsics and camera motion (visual odometry) for improving other tasks. CAM-Convs [33], which utilizes camera-aware convolutions for pinhole cameras, is the closest work. We extend and generalize this work to arbitrary camera geometries, including fisheye cameras.

The camera geometry adaptive mechanism is fundamental in the training process of the SVDistNet as the four different cameras mounted on the car have different intrinsic parameters and viewpoints. The trained distance and pose estimation networks need to generalize when deployed on a different car with a change in multiple viewpoints and intrinsics. To accomplish this, as shown in Fig. 3, we include the camera geometry tensor C_t in the mapping from RGB features to 3D information for distance estimation and semantic segmentation. In addition, as shown in Fig. 2, we add C_t to the pose encoder.

The tensor C_t can be obtained using a three step procedure: First, we pre-compute pixel-wise coordinates and the angle of the incidence maps to facilitate a more efficient training. To incorporate the knowledge encoded in these camera calibration parameters, the normalized coordinates per pixel are used as additional channels and concatenated, yielding the tensor C_t . Finally, these features are passed along with the encoder features to our SAN *pairwise* and *patchwise* models. Accordingly, we have six additional channels, supplementing the existing decoder channel inputs. In theory, the proposed method can be applied to any given fisheye projection model. We briefly present the standard projection models that our camera geometry tensor supports. The mapping of 3D points to pixels in fisheye lenses requires a radial component $\varrho(\vartheta)$. The polynomial model is the most commonly used, and more recent projection models include UCM (Unified Camera Model) [38] and eUCM (Enhanced UCM) [39]. Rectilinear (representation of a pinhole model) and Stereographic (mapping of a sphere to a plane) are not appropriate for fisheye lenses but are provided for comparison. The recently proposed Double Sphere model [40] has a closed-form inverse with low computational complexity. The following are summaries of the radial distortion models:

- 1) Polynomial: $\varrho(\vartheta) = a_1\vartheta + a_2\vartheta^2 + a_3\vartheta^3 + a_4\vartheta^4$
- 2) UCM: $\varrho(\vartheta) = f \cdot \sin \vartheta / (\cos \vartheta + \xi)$
- 3) eUCM: $\varrho(\vartheta) = f \cdot \frac{\sin \vartheta}{\cos \vartheta + \alpha (\sqrt{\beta \cdot \sin^2 \vartheta + \cos^2 \vartheta} - \cos \vartheta)}$
- 4) Rectilinear: $\varrho(\vartheta) = f \cdot \tan \vartheta$
- 5) Stereographic: $\varrho(\vartheta) = 2f \cdot \tan(\vartheta/2)$
- 6) Double Sphere: $\varrho(\vartheta) = f \cdot \frac{\sin \vartheta}{\alpha \sqrt{\sin^2 \vartheta + (\xi + \cos \vartheta)^2} + (1 - \alpha)(\xi + \cos \vartheta)}$

The distortion coefficients a_1, a_2, a_3, a_4 are used to create the angle of incidence maps (a_x, a_y) , c_x, c_y are used to compute the principal point coordinate maps (cc_x, cc_y) , and the camera's sensor dimensions (width w and height h) are used to formulate the normalized coordinate maps. We compute these maps using the camera intrinsic parameters and include them in our shared self-

attention encoder.

Centered Coordinates (cc) The SAN *pairwise* and *patchwise* operation modules receive the principal point position via the cc_x and cc_y coordinate channels, which are centred at $(0, 0)$. We create cc_x and cc_y channels, as shown below:

$$cc_x = \begin{pmatrix} 0 - c_x \\ 1 - c_x \\ \vdots \\ w - c_x \end{pmatrix} \cdot \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}^{\top}_{(h+1) \times 1} \quad cc_y = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}_{(w+1) \times 1} \cdot \begin{pmatrix} 0 - c_y \\ 1 - c_y \\ \vdots \\ h - c_y \end{pmatrix}^{\top} \quad (10)$$

We concatenate cc_x and cc_y by resizing them with bilinear interpolation to match the input feature's size.

Angle of Incidence Maps (a_x, a_y) Using the focal length f of the camera, the horizontal and vertical angle of incidence maps for the pinhole (Rectilinear) camera model are determined from the cc maps.

$$a_{ch}[i, j] = \arctan\left(\frac{cc_{ch}[i, j]}{f}\right) \quad (11)$$

where ch represents either x or y (see Eq. 10). By computing the inverse of the radial distortion functions $\varrho(\vartheta)$, angle of incidence maps for different fisheye camera models can be calculated. The angle of incidence ϑ is computed numerically for the polynomial model employed in this paper by computing the 4th order polynomial roots of $\varrho = \sqrt{x_I^2 + y_I^2} = a_1\vartheta + a_2\vartheta^2 + a_3\vartheta^3 + a_4\vartheta^4$. To improve training efficiency, we store the pre-calculated roots in a lookup table for all pixel coordinates and generate the a_x and a_y maps by setting $x_I = cc_x[i, j], y_I = 0$ and $x_I = 0, y_I = cc_y[i, j]$, respectively. A lookup table is not necessary for UCM, eUCM, and Double Sphere projection models as they can be inverted analytically.

Normalized Coordinates (nc) Following [41], [33], we additionally introduce two channels of normalized coordinates, whose values span linearly with respect to the image coordinates between -1 and 1 . The channels are unaffected by camera sensor properties and characterize the spatial extent of the content in feature space in each direction (e.g., a \hat{x} channel value close to 1 indicates that the feature vector at this location is near the right border).

B. Semantically-Guided Distance Decoder

To better guide the distance features, we adapt the pixel-adaptive convolutions from [42], [28], thereby distilling the knowledge of the extracted features from the segmentation branch into the distance decoder. This follows the intuition that the distance estimation can profit from the location-specific knowledge encoded in the semantic features, which in return also breaks up the convolutions' spatial invariance. While other semantic guidance approaches relied on a pretrained segmentation network, [7], [28], our semantic segmentation branch is learned simultaneously with the distance estimation, thereby introducing a more favorable one-stage training.

To implement the pixel-adaptive convolutions, we rely on the multi-level extraction of features from the segmentation decoder, as shown in Fig. 3. Accordingly, the input x is processed as follows:

$$x'_{ij} = \sum_{ab \in \mathcal{N}_k(i, j)} K(F_{ij}, F_{ab})W[r_{a-i, b-j}]x_{ab} + B \quad (12)$$

where $r_{a-i, b-j}$ represents the distance between pixel locations ij and ab and $\mathcal{N}_k(i, j)$ defines a neighborhood window of size $k \times k$ centered around the location ij . In Eq. 12, the elements of x which are covered by the neighborhood window $\mathcal{N}_k(i, j)$ make up the input to the convolution, characterized by a kernel function K , weights W , and a bias $B \in \mathbb{R}^1$. The kernel function is necessary to compute the correlation between the semantic features $F \in \mathbb{R}^D$ extracted by the segmentation decoder. Following [28], we define the kernel function as:

$$K(F_{ij}, F_{ab}) = \exp\left(-\frac{1}{2}(F_{ij} - F_{ab})^T \Sigma_{ijab}^{-1}(F_{ij} - F_{ab})\right) \quad (13)$$

where Σ_{ijab} is the covariance matrix between the features F_{ij} and F_{ab} . This matrix is chosen to be diagonal, which means that it can be expressed in terms of a learnable parameter σ for each convolutional filter as $\sigma^2 \cdot \mathbf{1}^D$.

V. IMPLEMENTATION DETAILS

To minimise the training objective function, we use Pytorch to implement the models described in Section IV and the Ranger (RAdam [43] + LookAhead [44]) optimizer. RAdam uses an active rectifier to adjust Adam's [45] adaptive momentum, resulting in an automated warm-up that is tailored to the dataset of interest. Meanwhile, LookAhead introduces a set of "slow weights", which are updated based on the previous weight updates of the normal optimizer (RAdam). The combination of both methods is highly synergistic, resulting in improved results. We train the model for 10 epochs because we now have 4x the data compared to training only the front camera images on a 24GB Titan RTX with a batch size of 20. We use a learning rate of 4×10^{-4} , which we reduce to 10^{-5} after 7 epochs for another 3 epochs. The distance decoder's sigmoid output σ is converted to distance using $D = m \cdot \sigma + n$. Depth is calculated for the pinhole model as $D = 1/(m \cdot \sigma + n)$, where n and m are constants chosen such that D is constrained between 0.1 and 100. As described in Section IV-A (also cf. Fig. 3), we shuffle all images from the surround-view cameras with diverse viewpoints and present them to the distance and pose models along with their corresponding intrinsics to build the camera geometry tensor C_t . We limited the length of the sequences to three. The lengths of the sequences considered for pose estimation results are 2, 3, and 5. The semantic segmentation is trained supervised by applying the cross-entropy loss and optimized jointly with the distance estimation (see Fig. 2).

Comparison to Convolutional Networks Pairwise models match or outperform convolutional baselines while requiring comparable or lower parameters and FLOP budgets. In terms of computational complexity, the patchwise models outperform CNN's. For example, the patchwise SAN10 with 11.8M params and 1.9G FLOPS outperforms ResNet50 with 25.6M params and 4.1G FLOPS, while having a 54% lower parameter and 44% lower FLOP count. SAN10-patch models' parameter count is almost equivalent to ResNet18 with 11.7M params and 1.8G FLOPS, whereas SAN15-patch with 20.5M params and 3.3G FLOPS is equivalent to ResNet50's parameter count. The SAN10-pairwise with 10.5M params and 2.2G FLOPS

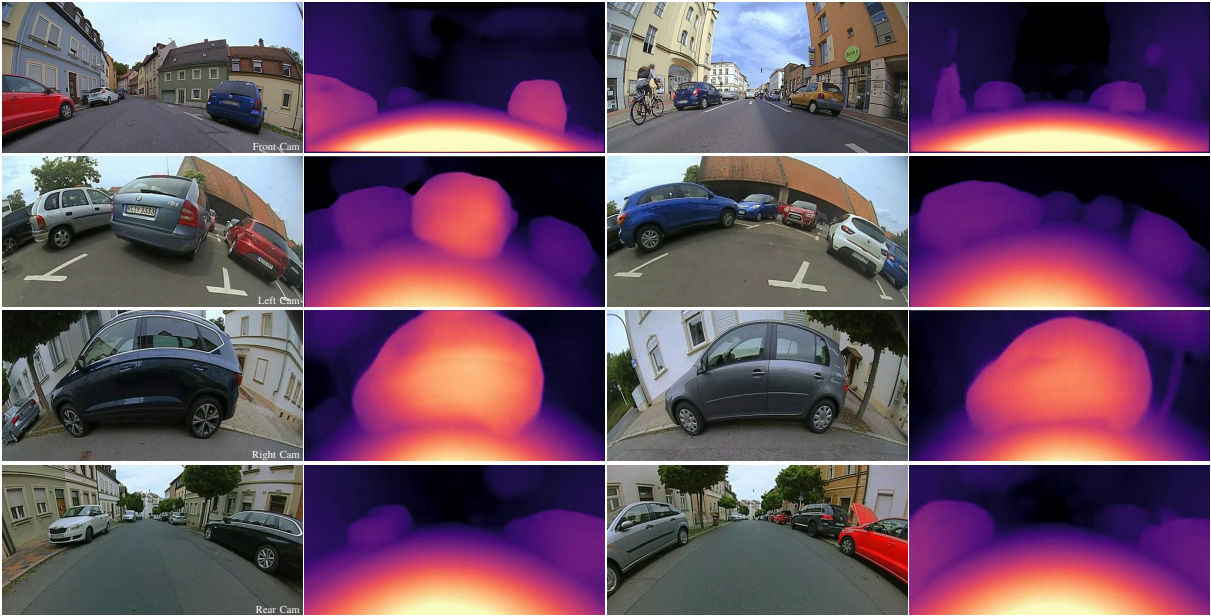


Fig. 4: **Distance estimation results** of the same network evaluated on four different fisheye cameras of a surround-view camera system. One can see that our SVDistNet model generalizes well across different viewing angles and consistently produces high-quality distance outputs.

Method	Seg	Dist.	MTL	MTL (Synergy)	Seg. (mIOU)	Distance (RMSE)
SynDistNet [11]	✓	✗	✗	✗	76.8	✗
	✗	✓	✗	✗	✗	2.316
	✓	✓	✓	✗	78.3	2.128
	✓	✓	✗	✓	81.5	1.714
SVDistNet (SAN10-patch)	✓	✗	✗	✗	77.1	✗
	✗	✓	✗	✗	✗	2.153
	✓	✓	✓	✗	78.6	1.861
	✓	✓	✗	✓	82.3	1.532

TABLE I: **Effect of our multi-task training approach** on semantic segmentation and distance estimation when using a ResNet 18 backbone for SynDistNet [11] and self-attention network SAN-10 for SVDistNet. We observe that both outputs improve through the multi-task training and mainly the distance estimation performance profits from the synergized semantic guidance.

outperforms ResNet18 with a 9% lower parameter count and 22% higher FLOP count. By performing a single-scale image depth prediction rather than the multi-scale in our previous works [9], [10], we could leverage the use of a more robust loss function over L_1 to reduce training times on SAN10.

VI. EXPERIMENTAL EVALUATION

A. Effect of Multi-Task Learning

In Table I, we compare our SVDistnet with previous work [11] on the Woodscape dataset [3]. Specifically, we evaluate our model using the patchwise self-attention encoder. First, we create single-task baselines using the various methods mentioned above, and then we create an important shared encoder-based multi-task learning (MTL) baseline for these tasks. Due to the use of an improved self-attention encoder and the camera geometry tensor, the MTL results and single-task benchmark for distance estimation outperform SynDistNet [11]. However, we observe only minimal gain for

the semantic segmentation task. In the final experiment, we include the synergy between the distance and segmentation decoders. For these diverse tasks, we observe that the content and channel adaptive self-attention encoder features can be learned jointly, wherein the captured semantic features, used along with pixel-adaptive convolutions, guide the distance decoder to capture better geometric features. We use ablation experiments to further dissect these findings.

B. Ablation Study on the Single Contributions

For our ablation analysis, we consider two self-attention encoder variants: pairwise and patchwise. Table II displays the distance estimation results for these variants. We show the impact of certain contributions and their importance in our SVDistNet framework compared to our previous work SynDistNet [11]. As a start, we replace the L_1 loss with the general robust loss function and test it on the self-attention encoder’s patchwise variant, which results in significant accuracy improvements. To filter all the dynamic objects that contaminate the reconstruction loss, we use the dynamic object mask obtained by projecting semantic segmentation predictions as described in Section III-B. Furthermore, this contribution has the potential to solve the infinite distance problem.

We see a significant increase in accuracy after adding the camera geometry tensor to this setting. The training on multiple cameras with varying intrinsics and viewing angles adds an additional regularization. The network’s aforementioned training strategy is thereby more camera-independent and leads to a better generalization to images captured with varying cameras. To improve our multi-task setup even further, we use a distance and segmentation decoder synergy. We apply semantic features from the supervised task to the distance decoder’s geometric features while still training the distance

Network	RL	Self-Attn	SEM	Mask	CGT	Abs Rel	Sq Rel	RMSE	RMSE _{log}	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
						lower is better				higher is better		
FisheyeDistanceNet [9]	✗	✗	✗	✗	✗	0.152	0.768	2.723	0.210	0.812	0.954	0.974
SynDistNet(ResNet18) [11]	✓	✗	✗	✗	✗	0.142	0.537	2.316	0.179	0.878	0.971	0.985
	✓	✗	✗	✓	✗	0.133	0.491	2.264	0.168	0.868	0.976	0.988
	✓	✓	✗	✓	✗	0.121	0.429	2.128	0.155	0.875	0.980	0.990
	✓	✓	✓	✗	✗	0.105	0.396	1.976	0.143	0.878	0.982	0.992
	✓	✓	✓	✓	✗	0.076	0.368	1.714	0.127	0.891	0.988	0.994
SynDistNet(ResNet50) [11]	✓	✗	✗	✗	✗	0.138	0.540	2.279	0.177	0.880	0.973	0.986
	✓	✗	✗	✓	✗	0.127	0.485	2.204	0.166	0.881	0.975	0.989
	✓	✓	✗	✓	✗	0.115	0.413	2.028	0.148	0.876	0.983	0.992
	✓	✓	✓	✗	✗	0.102	0.387	1.856	0.135	0.884	0.985	0.994
	✓	✓	✓	✓	✗	0.068	0.352	1.668	0.121	0.895	0.990	0.996
SVDistNet (SAN10-patch)	✓	✓	✗	✗	✗	0.128	0.469	2.153	0.164	0.875	0.974	0.986
	✓	✓	✗	✓	✗	0.114	0.413	2.022	0.149	0.878	0.982	0.990
	✓	✓	✗	✓	✓	0.101	0.378	1.861	0.133	0.884	0.984	0.991
	✓	✓	✓	✗	✗	0.094	0.345	1.789	0.128	0.887	0.985	0.992
	✓	✓	✓	✗	✓	0.082	0.316	1.682	0.119	0.890	0.987	0.993
	✓	✓	✓	✓	✗	0.074	0.343	1.641	0.112	0.892	0.985	0.994
	✓	✓	✓	✓	✓	0.057	0.315	1.532	0.108	0.896	0.986	0.996
SVDistNet (SAN10-pair)	✓	✓	✗	✓	✓	0.121	0.457	2.115	0.152	0.879	0.979	0.985
	✓	✓	✓	✗	✓	0.103	0.385	1.882	0.141	0.882	0.983	0.990
	✓	✓	✓	✓	✓	0.081	0.365	1.710	0.128	0.890	0.985	0.994
SVDistNet (SAN19-patch)	✓	✓	✗	✗	✗	0.121	0.437	2.127	0.153	0.878	0.976	0.989
	✓	✓	✗	✓	✗	0.109	0.408	2.006	0.145	0.880	0.982	0.992
	✓	✓	✗	✓	✓	0.098	0.372	1.849	0.138	0.884	0.983	0.991
	✓	✓	✓	✗	✗	0.091	0.351	1.773	0.129	0.886	0.986	0.993
	✓	✓	✓	✗	✓	0.070	0.305	1.669	0.108	0.893	0.986	0.994
	✓	✓	✓	✓	✗	0.067	0.296	1.578	0.106	0.895	0.985	0.994
SVDistNet (SAN19-pair)	✓	✓	✗	✓	✓	0.116	0.461	2.097	0.154	0.881	0.982	0.988
	✓	✓	✓	✗	✓	0.096	0.371	1.846	0.147	0.884	0.985	0.991
	✓	✓	✓	✓	✓	0.074	0.331	1.624	0.101	0.891	0.986	0.994

TABLE II: **Ablation study on the effect of our contributions** up to our final SVDistNet model on the Fisheye Woodscape dataset [3]. We cap the distance estimates to 40m. From our distance estimation baseline [9], we incrementally add up the robust loss (RL), self-attention layers encoder heads (Self-Attn), semantic guidance in the decoder (SEM), dynamic object masking (Mask), and camera geometry tensor (CGT). We showcase our improvements for various network architectures and, in particular, show the superiority of our SVDistNet model over the SynDistNet model as well as the positive effect of using camera geometry tensor C_t .

estimation in a self-supervised manner. In this case, we can remove the dynamic object mask and still see improvements.

By adding the camera geometry tensor to this setting, we can improve the metric results even further and potentially outperform the best setting in SynDistNet [11]. Using these critical features, we create an experiment similar to SynDistNet’s final setting, which combines all of the listed features except the camera geometry tensor. We could attain better results than SynDistNet’s ResNet50 results. It is important to note that ResNet50 is comparable to the SAN19-pair/patch encoder. However, in terms of computational complexity, we were able to outperform ResNet50 with the SAN10-patch encoder (cf. Section V).

By introducing our camera geometry tensor, we completed our SVDistNet model for the surround view-camera framework. We also performed a few vital evaluations using the pairwise self-attention encoder head. We were unable to achieve the same level of accuracy as the patchwise self-attention encoder. The patchwise self-attention module outperforms convolution, while the pairwise self-attention module equals or outperforms the convolutional equivalents. We use a higher-

Cams	CGT	Abs Rel	Sq Rel	RMSE	RMSE _{log}	$\delta < 1.25$	$\delta < 1.25^2$	$\delta < 1.25^3$
		lower is better				higher is better		
Front	✗	0.074	0.343	1.641	0.112	0.892	0.985	0.994
	✓	0.057	0.315	1.532	0.108	0.896	0.987	0.996
Rear	✗	0.089	0.358	1.657	0.131	0.888	0.981	0.988
	✓	0.065	0.337	1.579	0.123	0.891	0.983	0.992
Left	✗	0.102	0.398	1.874	0.126	0.886	0.980	0.983
	✓	0.091	0.382	1.781	0.114	0.889	0.985	0.990
Right	✗	0.105	0.406	1.889	0.135	0.882	0.980	0.981
	✓	0.093	0.391	1.796	0.120	0.887	0.983	0.986

TABLE III: **Ablation study on multiple cameras** concerning the usage of Camera Geometry Tensor using WoodScape [3].

performing SAN19-patch encoder network to perform the same vital ablation of our contributions. Finally, with all of our additions, we outperform all previous works on the Woodscape dataset [9], [10], [11].

C. Generalization Across Different Cameras

The generalization across different cameras from the surround-view setup is depicted in Table III using the camera geometry tensor as described in Section IV-A. Each camera’s

Method	Abs _{rel}	Sq _{rel}	RMSE	RMSE _{log}	δ_1	δ_2	δ_3	
	lower is better				higher is better			
Original [46]	Monodepth2 [5]	0.115	0.903	4.863	0.193	0.877	0.959	0.981
	PackNet-SfM [13]	0.111	0.829	4.788	0.199	0.864	0.954	0.980
	FisheyeDistanceNet [9]	0.117	0.867	4.739	0.190	0.869	0.960	0.982
	UnRectDepthNet [10]	0.107	0.721	4.564	0.178	0.894	0.971	0.986
	SynDistNet [11]	0.109	0.718	4.516	0.180	0.896	0.973	0.986
	Shu <i>et al.</i> [47]	0.104	0.729	4.481	0.179	0.893	0.965	0.984
	SVDistNet	0.102	0.706	4.459	0.172	0.908	0.974	0.986
	Struct2Depth* [7]	0.109	0.825	4.750	0.187	0.874	0.958	0.983
	GLNet* [25]	0.099	0.796	4.743	0.186	0.884	0.955	0.979
	Shu* <i>et al.</i> [47]	0.088	0.712	4.137	0.169	0.915	0.965	0.982
SVDistNet*	0.086	0.701	4.118	0.170	0.919	0.976	0.985	
Improved [48]	Monodepth2 [5]	0.090	0.545	3.942	0.137	0.914	0.983	0.995
	PackNet-SfM [13]	0.078	0.420	3.485	0.121	0.931	0.986	0.996
	UnRectDepthNet [10]	0.081	0.414	3.412	0.117	0.926	0.987	0.996
	SynDistNet [11]	0.076	0.412	3.406	0.115	0.931	0.988	0.996
	SVDistNet	0.071	0.405	3.345	0.106	0.934	0.988	0.996
	SVDistNet*	0.059	0.392	3.206	0.097	0.935	0.989	0.995

TABLE IV: **Evaluation of the KITTI Eigen split** in comparison to other self-supervised monocular depth estimation methods. We limit depths to 80m in accordance with best practices. We also evaluate using the *Original* depth maps generated from raw point clouds as proposed by [46] and *Improved* annotated depth maps as introduced by [48]. M indicates that sequences are trained on using the monocular approach. The method* indicates the online refinement technique [7], in which the model is trained during inference.

metrics significantly improve as sequences from different cameras aid in generalization during the model’s training phase. For example, the front camera’s distance estimates profit as the side cameras steer the network to generalize close and overlapping objects. Because of the use of C_t , our network does not overfit to a specific camera intrinsic. It adapts to any changes detected by a family of unseen cameras deployed with a pre-calibrated camera setup. This leads to better estimates and generalization of new cameras and the ability to train on images from different cameras.

D. State-of-the-Art Comparison on KITTI

To facilitate comparisons with previous methods, we also train our distance estimation method on the KITTI Eigen split [49], [46] in the classical depth estimation setting. First, we train and evaluate on depth maps generated from LiDAR point clouds, where Table IV shows that with our contributions, textitwe outperform all previous methods.

We obtain a significant improvement using the online refinement method from [7], and our results remain superior to previous methods. Also, when training and evaluating the improved KITTI labels for depth estimation, we show a significant improvement. Simultaneously, our network architecture’s complexity is still comparable to the one of a ResNet18 in terms of computational complexity. We also use the general camera tensor C_t as described in Section IV-A in our model, wherein instead of the angle of incidence maps, we employ the maps generated using Eq. 11 for pinhole cameras.

VII. CONCLUSION

In this paper, we solved the problem of multi-camera distance estimation for surround-view fisheye cameras. We introduced a novel camera model adaptation mechanism wherein camera parameters are transformed into a tensor and used within the CNN model. The specific camera model parameters

are used during training and inference. Using this technique, we demonstrate training of a single distance estimation model for twelve different cameras with different extrinsic and intrinsic parameters and achieve the improved results as training a specialized model for each camera variant. We achieve state-of-the-art results on the fisheye WoodScape dataset and the KITTI dataset. We intend to learn and refine camera parameters’ within the training framework in future work.

REFERENCES

- [1] M. Heimberger, J. Horgan, C. Hughes, J. McDonald, and S. Yogamani, “Computer vision in automated parking systems: Design, implementation and challenges,” in *Image and Vision Computing*. Elsevier, 2017.
- [2] V. Ravi Kumar, S. Milz, C. Witt, M. Simon, K. Amende, J. Petzold, S. Yogamani, and T. Pech, “Monocular fisheye camera depth estimation using sparse lidar supervision,” in *ITSC*, 2018.
- [3] S. Yogamani, C. Hughes, J. Horgan, G. Sistu, P. Varley *et al.*, “Wood-scape: A Multi-Task, Multi-Camera Fisheye Dataset for Autonomous Driving,” in *ICCV*, 2019.
- [4] T. Zhou, M. Brown, N. Snavely, and D. G. Lowe, “Unsupervised learning of depth and ego-motion from video,” in *CVPR*, 2017.
- [5] C. Godard, O. Mac Aodha, M. Firman, and G. J. Brostow, “Digging Into Self-Supervised Monocular Depth Estimation,” in *ICCV*, 2019.
- [6] V. Ravi Kumar, S. Yogamani, S. Milz, and P. Mäder, “FisheyeDistanceNet++: Self-Supervised Fisheye Distance Estimation with Self-Attention, Robust Loss Function and Camera View Generalization,” in *Electronic Imaging*. Society for Imaging Science and Technology, 2021.
- [7] V. Casser, S. Pirk, R. Mahjourian, and A. Angelova, “Depth Prediction Without the Sensors: Leveraging Structure for Unsupervised Learning from Monocular Videos,” in *AAAI*, 2019.
- [8] M. Klingner, J.-A. Termöhlen, J. Mikolajczyk, and T. Fingscheidt, “Self-Supervised Monocular Depth Estimation: Solving the Dynamic Object Problem by Semantic Guidance,” in *ECCV*, 2020.
- [9] V. Ravi Kumar, S. A. Hiremath, M. Bach, S. Milz, C. Witt *et al.*, “Fisheyedistancenet: Self-supervised scale-aware distance estimation using monocular fisheye camera for autonomous driving,” in *ICRA*, 2020.
- [10] V. Ravi Kumar, S. Yogamani, M. Bach, C. Witt, S. Milz *et al.*, “UnRect-DepthNet: Self-Supervised Monocular Depth Estimation using a Generic Framework for Handling Common Camera Distortion Models,” in *IROS*, 2020.
- [11] V. Ravi Kumar, M. Klingner, S. Yogamani, S. Milz, T. Fingscheidt *et al.*, “SynDistNet: Self-Supervised Monocular Fisheye Camera Distance Estimation Synergized with Semantic Segmentation for Autonomous Driving,” in *WACV*, 2021.
- [12] C. Godard, O. Mac Aodha, and G. J. Brostow, “Unsupervised Monocular Depth Estimation With Left-Right Consistency,” in *CVPR*, 2017.
- [13] V. Guizilini, R. Ambrus, S. Pillai, and A. Gaidon, “3D Packing for Self-Supervised Monocular Depth Estimation,” in *CVPR*, 2020.
- [14] F. Tosi, F. Aleotti, M. Poggi, and S. Mattocchia, “Learning Monocular Depth Estimation Infusing Traditional Stereo Knowledge,” in *CVPR*, 2019.
- [15] R. Wang, S. M. Pizer, and J.-M. Frahm, “Recurrent Neural Network for (Un-)Supervised Learning of Monocular Video Visual Odometry and Depth,” in *CVPR*, 2019.
- [16] H. Zhang, C. Shen, Y. Li, Y. Cao, Y. Liu *et al.*, “Exploiting Temporal Consistency for Real-Time Video Depth Estimation,” in *ICCV*, 2019.
- [17] V. Casser, S. Pirk, R. Mahjourian, and A. Angelova, “Unsupervised Monocular Depth and Ego-Motion Learning With Structure and Semantics,” in *CVPR-Workshops*, 2019.
- [18] A. Kendall, Y. Gal, and R. Cipolla, “Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics,” in *CVPR*, 2018.
- [19] L. Wang, J. Zhang, O. Wang, Z. Lin, and H. Lu, “SDC-Depth: Semantic Divide-and-Conquer Network for Monocular DepthEstimation,” in *CVPR*, 2020.
- [20] M. Klingner, A. Bär, and T. Fingscheidt, “Improved Noise and Attack Robustness for Semantic Segmentation by Using Multi-Task Training with Self-Supervised Depth Estimation,” in *CVPR - Workshops*, 2020.
- [21] M. Klingner and T. Fingscheidt, “Online Performance Prediction of Perception DNNs by Multi-Task Learning with Depth Estimation,” *IEEE Transactions on Intelligent Transportation Systems*, 2021.

- [22] S. Zhao, H. Fu, M. Gong, and D. Tao, "Geometry-Aware Symmetric Domain Adaptation for Monocular Depth Estimation," in *CVPR*, 2019.
- [23] P. Liu, M. Lyu, I. King, and J. Xu, "SelfFlow: Self-Supervised Learning of Optical Flow," in *CVPR*, 2019.
- [24] L. Liu, G. Zhai, W. Ye, and Y. Liu, "Unsupervised Learning of Scene Flow Estimation Fusing With Local Rigidity," in *IJCAI*, 2019.
- [25] Y. Chen, C. Schmid, and C. Sminchisescu, "Self-Supervised Learning With Geometric Constraints in Monocular Video Connecting Flow, Depth, and Camera," in *ICCV*, 2019.
- [26] V. Nekrasov, T. Dharmasiri, A. Spek, T. Drummond, C. Shen *et al.*, "Real-time joint semantic segmentation and depth estimation using asymmetric annotations," in *ICRA*, 2019.
- [27] V.-C. Miclea and S. Nedeveschi, "Real-time semantic segmentation-based stereo reconstruction," in *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [28] V. Guizilini, R. Hou, J. Li, R. Ambrus, and A. Gaidon, "Semantically-Guided Representation Learning for Self-Supervised Monocular Depth," in *ICLR*, 2020.
- [29] P.-Y. Chen, A. H. Liu, Y.-C. Liu, and Y.-C. F. Wang, "Towards Scene Understanding: Unsupervised Monocular Depth Estimation With Semantic-Aware Representation," in *CVPR*, 2019.
- [30] S. Zhu, G. Brazil, and X. Liu, "The Edge of Depth: Explicit Constraints between Segmentation and Depth," in *CVPR*, 2020.
- [31] G. Yang, H. Zhao, J. Shi, Z. Deng, and J. Jia, "SegStereo: Exploiting Semantic Information for Disparity Estimation," in *ECCV*, 2018.
- [32] A. Gordon, H. Li, R. Jonschkowski, and A. Angelova, "Depth from Videos in the Wild: Unsupervised Monocular Depth Learning from Unknown Cameras," in *ICCV*, 2019.
- [33] J. M. Facil, B. Ummenhofer, H. Zhou, L. Montesano, T. Brox *et al.*, "CAM-Convs: Camera-Aware Multi-Scale Convolutions for Single-View Depth," in *CVPR*, 2019.
- [34] F.-E. Wang, Y.-H. Yeh, M. Sun, W.-C. Chiu, and Y.-H. Tsai, "BiFuse: Monocular 360° Depth Estimation via Bi-Projection Fusion," in *CVPR*, 2020.
- [35] L. Jin, Y. Xu, J. Zheng, J. Zhang, R. Tang *et al.*, "Geometric Structure Based and Regularized Depth Estimation From 360° Indoor Imagery," in *CVPR*, 2020.
- [36] J. T. Barron, "A General and Adaptive Robust Loss Function," in *CVPR*, 2019.
- [37] H. Zhao, J. Jia, and V. Koltun, "Exploring self-attention for image recognition," in *CVPR*, 2020.
- [38] J. P. Barreto, "Unifying image plane liftings for central catadioptric and dioptric cameras," in *Imaging Beyond the Pinhole Camera*. Springer, 2006.
- [39] B. Khomutenko, G. Garcia, and P. Martinet, "An Enhanced Unified Camera Model," in *RA-L*, 2016.
- [40] V. Usenko, N. Demmel, and D. Cremers, "The double sphere camera model," in *3DV*, 2018.
- [41] R. Liu, J. Lehman, P. Molino, F. P. Such, E. Frank *et al.*, "An Intriguing Failing of Convolutional Neural Networks and the Coordconv solution," in *NIPS*, 2018.
- [42] H. Su, V. Jampani, D. Sun, O. Gallo, E. Learned-Miller *et al.*, "Pixel-Adaptive Convolutional Neural Networks," in *CVPR*, 2019.
- [43] L. Liu, H. Jiang, P. He, W. Chen, X. Liu, J. Gao, and J. Han, "On the variance of the adaptive learning rate and beyond," in *ICLR*, 2020.
- [44] M. Zhang, J. Lucas, J. Ba, and G. E. Hinton, "Lookahead Optimizer: k steps forward, 1 step back," in *NIPS*, 2019.
- [45] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," in *arXiv preprint arXiv:1412.6980*, 2014.
- [46] D. Eigen, C. Puhrsch, and R. Fergus, "Depth Map Prediction from a Single Image Using a Multi-Scale Deep Network," in *NIPS*, 2014.
- [47] C. Shu, K. Yu, Z. Duan, and K. Yang, "Feature-metric Loss for Self-supervised Learning of Depth and Egomotion," in *ECCV*, 2020.
- [48] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox *et al.*, "Sparsity Invariant CNNs," in *3DV*, 2017.
- [49] A. Geiger, P. Lenz, C. Stiller, and R. Urtasun, "Vision Meets Robotics: The KITTI Dataset," in *The Int. Journal of Robotics Research*, 2013.



Varun Ravi Kumar received a B.E. degree in 2015 & an M.Sc. degree in 2017 from TU Chemnitz, Germany. He is currently a Ph.D. student in Deep Learning for autonomous driving affiliated to TU Ilmenau and is currently working at Valeo. His research is mainly focused on the design of self-supervised perception algorithms using neural networks for self-driving cars. His diverse expertise lies in geometric, semantic tasks, 2D & 3D object detection, point-cloud processing & multi-task modeling.



Marvin Klingner received a B.Sc. degree in 2016 & an M.Sc. degree in 2018 in physics from Georg-August-Universität Göttingen, Germany. He is currently a Ph.D. student in the Faculty of Electrical Engineering, Information Technology, Physics at Technische Universität Braunschweig, Germany. His research interests in computer vision for autonomous driving include self-supervised 3D geometry perception with neural networks, multi-task learning approaches, the robustness of neural networks to domain shifts & adversarial attacks & performance monitoring neural networks during online deployment. He is the recipient of the Dr. Berliner - Dr. Ungewitter award of the Faculty of Physics at Georg-August-Universität Göttingen, Germany, in 2018, was given the CVPR SIAID 2020 Workshop Best Paper Award.



Senthil Yogamani is an Artificial Intelligence architect & holds a director-level technical leader position at Valeo Ireland. He leads the research & design of AI algorithms for various modules of autonomous driving systems. He has over 14 years of experience in computer vision & machine learning, including 12 years of experience in industrial, automotive systems. He is an author of 100+ publications with 1600+ citations & 100+ inventions with 70 filed patent families. He serves on the editorial board of various leading IEEE automotive conferences.



Markus Bach received the Dipl.-Phys. degree in 2013 from TU Dresden, Germany. He has over 5 years of research experience in the field of theoretical particle physics. In 2019 he joined Valeo as a software engineer in Driving Assistance Research & is working on diverse topics: SLAM & localization based on various inputs (camera, Lidar, Radar), computer vision, intrinsic & extrinsic sensor calibration, assurance of AI-based perception functions for automated driving.



Stefan Milz received his Ph.D. degree in Physics from the Technical University of Munich. He has a strong history in professional software development & automotive. He is Managing Director of Spleenlab.ai, a self-founded machine learning company focusing on safety-critical computer vision applications (UAV, Automated Driving, Air-Taxis) deploying SLAM, sensor-fusion, perception functions into the real world. Besides, he is also a research fellow at the TU-Ilmenau. Stefan Milz is the author & co-author of more than 60 patents & publications.



Tim Fingscheidt received the Dipl.-Ing. degree in electrical engineering in 1993 & the Ph.D. degree in 1998 from RWTH Aachen University, Germany. He joined AT&T Labs, Florham Park, NJ, USA, in 1998, & Siemens AG (Mobile Devices), Munich, Germany, in 1999. With Siemens Corporate Technology, Munich, Germany, he led the speech technology development activities (2005–2006). Since 2006, he has been a Full Professor with the Institute for Communications Technology, Technische Universität Braunschweig, Germany.



Patrick Mäder is a Professor at the Technische Universität Ilmenau, Germany & heading the endowed chair on Software Engineering for Critical Systems. His research interests include software engineering focusing on requirements traceability, requirements engineering, object-oriented analysis & design, & development of safety-critical systems. He received a Diploma degree in industrial engineering & a Ph.D. degree (Distinction) in computer science from the Technische Universität Ilmenau in 2003 & 2009, respectively.