# Neural Rendering based Urban Scene Reconstruction for Autonomous Driving

*Shihao Shen* [1] *, Louis Kerofsky* [1] *, Varun Ravi Kumar* [1] *and Senthil Yogamani* [2]

[1] *Qualcomm Technologies, Inc., San Diego, California, U.S.*

[2] *Automated Driving, QT Technologies Ireland Limited.*

## ABSTRACT

*Dense 3D reconstruction has many applications in automated driving including automated annotation validation, multimodal data augmentation, providing ground truth annotations for systems lacking LiDAR, as well as enhancing auto-labeling accuracy. LiDAR provides highly accurate but sparse depth, whereas camera images enable estimation of dense depth but noisy particularly at long ranges. In this paper, we harness the strengths of both sensors and propose a multimodal 3D scene reconstruction using a framework combining neural implicit surfaces and radiance fields. In particular, our method estimates dense and accurate 3D structures and creates an implicit map representation based on signed distance fields, which can be further rendered into RGB images, and depth maps. A mesh can be extracted from the learned signed distance field and culled based on occlusion. Dynamic objects are efficiently filtered on the fly during sampling using 3D object detection models. We demonstrate qualitative and quantitative results on challenging automotive scenes.*
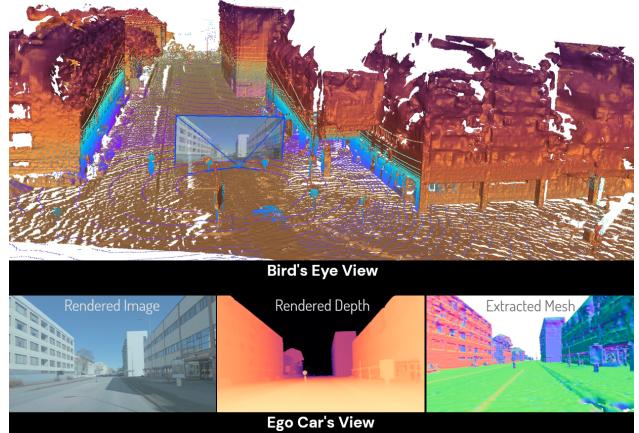
Figure 1: We demonstrate dense and accurate 3D structure being represented by an implicit model, which can be further rendered into RGB images and depth maps or reconstructed into a high-quality mesh.

## I INTRODUCTION

Advancement in the field of computer vision has enabled the rapid development of perception systems for autonomous vehicles (AV) in recent years. Deep learning in particular has accelerated this progress by achieving rapid advancement in various perception tasks including object detection [1–3], semantic segmentation [4–6], depth prediction [7–13], adverse weather detection [14–17], moving object detection [18–20], SLAM [21–23], multi-task learning [24–26] and sensor fusion [27–29]. However, dense scene construction of urban scenes is relatively less mature due to a variety of challenges. Firstly, it is challenging to reconstruct a spatially consistent dense structure map over time due to odometry errors. Secondly, there are many moving objects in urban scenes which hinders the 3D reconstruction process. Finally, urban scenes have fine 3D structures like curb and poles which are challenging to estimate.

3D scene reconstruction, refers to the creation of three-dimensional models from available data modalities, e.g. a set of images. Traditionally, 3D scenes are represented by point clouds, voxels, or meshes that are explicit and discrete. Recent Neural Radiance Fields (NeRF) models are able to represent a continuous scene implicitly through Multi-Layer Perceptrons (MLPs) that learn the scene geometry and appearance simultaneously. Since NeRF was first introduced by Mildenhall et al. [30], it has been widely explored and adapted into variants for a variety of applications, including city reconstruction [31], image processing [32],

generative AI [33] to point out a few. This paper aims to tackle the application of reconstructing unconstrained urban scenes from monocular recordings.

Reconstructing 3D urban environments from sensor data is an important problem with applications in autonomous vehicles, augmented reality, city planning and more. Traditional approaches rely on fusing data from LiDAR and structure from motion techniques applied to camera images. However, these explicit representations have limitations such as sparsity, noise, or difficulty scaling to large scenes. Neural implicit functions provide an alternative representation that is compact, smooth and can easily scale to model complex scenes. Recent works have shown promising results using neural radiance fields and occupancy networks to reconstruct indoor and small outdoor scenes from images [34–36]. However, applying these techniques to reconstruct large-scale urban environments from vehicle sensors poses new challenges. In this paper, we introduce a method to urban scene reconstruction using a framework combining neural implicit surfaces and radiance fields that addresses these challenges. Example use cases of the method include online 3D environmental model or offline extraction of 3D instances for multimodal data augmentation. In particular, we plan to deploy our trained model to aid our automated labelling pipeline.
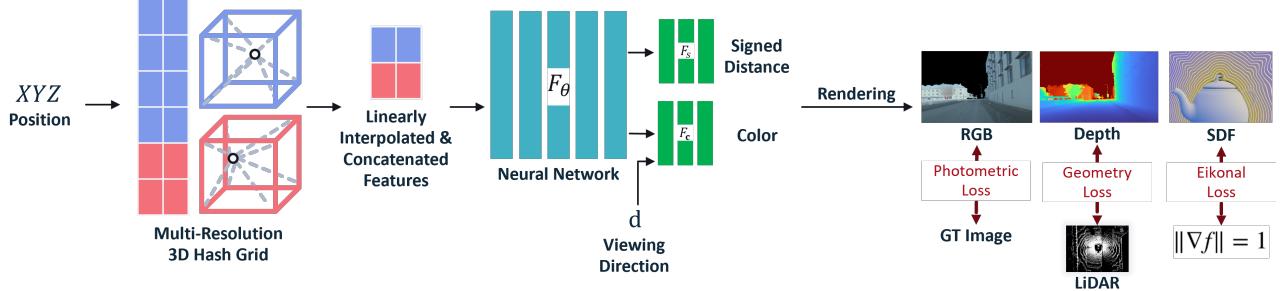
Figure 2: Overview of our foreground model. It follows the same multiresolution hash grid design as in [37] and predicts signed distance rather than density. The other MLP head predicts view-dependent color by taking in the viewing direction. Three major supervisions are photometric loss to supervise the reconstructed scene appearance, Eikonal loss to regularize the learned SDF, and geometry loss to supervise the reconstructed scene geometry.

## II  METHODOLOGY

In this section, we describe the 3D scene reconstruction framework, the challenges we have faced, as well as the corresponding solutions.

### II. A  Background

Neural Radiance Fields (NeRF) [30] are composed of multiple MLPs that implicitly represent a target scene. A NeRF takes in a 3D position x and viewing direction d, and outputs volume density $\sigma$ and color $c$. The viewing direction input enables NeRF to learn view-dependent color $c$ and hence to reproduce the reflective appearance of certain materials in the scene. On the other hand, the volume density output $\sigma$ is made viewing direction independent because intuitively volume density aims to reproduces the scene geometry, which is constant irrespective of viewing direction.

To reproduce the appearance of the scene, or in other words, to render a pixel in an image, the renderer shoots a ray from the camera center o through the pixel to infinity, denoted as $r(t) = o + td$, where d is the direction from center camera to the pixel. Then it randomly samples points along the ray with distances $\{t_i\}_{i=0}^N$ from the o. Those points $r(t_i)$ and viewing direction d are passed to the MLPs which produce $c_i$ and $\sigma_i$ for every point. Alpha compositing is used to get the final color of this pixel, with alphas equal to the predicted densities:

$$c_{\text{out}} = \sum_{i=1}^N T_i(1 - e^{-\delta_i \sigma_i})c_i$$

$$\text{where } T_i = \exp\left(-\sum_{j=0}^{i-1} \delta_j \sigma_j\right) \text{ and } \delta_i = t_i - t_{i-1} \quad (1)$$

Similarly, to reproduce the geometry of the scene, or in other words, to render a pixel in a depth map, we simply replace $c_i$ in Equation 1 with depth $t_i$:

$$d_{\text{out}} = \sum_{i=1}^N T_i(1 - e^{-\delta_i \sigma_i})t_i \quad (2)$$

Simple coordinates $x = (X, Y, Z)$ lack the expressiveness for high-frequency details, so x and d are encoded into higher-dimensional vectors of sine and cosine functions, known as sinusoidal positional encoding [30] $\gamma_{PE}$.

However, naively sampling points along a ray becomes intractable in outdoor scenes because the spatial coordinate could theoretically be infinite (e.g., the sky). Therefore, NeRF++ [34] decomposes the scene into foreground and background, with foreground enclosed by a volume of unit sphere and the background enclosed by a volume of an inverted sphere. Instead of inputting x into the background model, the input becomes $(x', 1/r)$ where $x'$ is x projected onto the unit sphere and $r$ is the distance from x to the sphere's center, which effectively maps unbounded spatial inputs to bounded ones, avoids numeral unstable problems, and hence facilitates the representation of background scene.

After achieving high-quality representation, speed becomes a concern. Instant-NGP [37] improves both training and rendering by using spatial data structures to store neural features which can be subsequently interpolated into feature vectors per spatial coordinate. It also adopts a multiresolution hash table for encodings, which significantly decreases the capacity required of the prediction MLPs. We refer to [37] for details.

Although radiance field representations like a NeRF have shown great performance in novel view synthesis, they have been demonstrated to be unstable and ambiguous toward accurate geometry [38, 39] because the underlying volume density is non-smooth and prone to artifacts. Intuitively, there's no constraint enforced on the volume density in empty space and hence it becomes difficult to extract watertight geometry from it. Therefore, prior work embraces neural implicit surface to achieve accurate geometry. Both VolSDF [40] and NeuS [38] replace the output density $\sigma$ in NeRF by signed distance and then analytically convert the predicted signed distance to density to be used in the same way as volume rendering in NeRF. Because the network is now essentially a signed distance function (SDF), we can readily enforce the SDF prior in supervision, known as the Eikonal loss $|\|\nabla f\| - 1|$. As such, the network represents the scene geometry with regularization, which is critical to our scene reconstruction use case.

### II. B  Scene Decomposition

Inspired by prior work [34, 41, 42], we decompose the scene into foreground and background and fit two different models to each of them. We define the foreground to be the street scene, including buildings, trees, road elements, etc., and define the background to be landscapes or the sky, regions that can't be reached by the ego vehicle in the current data sequence. We adapt our solution from StreetSurf [41] as well as use part of it as the backbone

of our solution.

Figure 2 illustrates the foreground model. Specifically, we follow the same design in [37], which stores feature embeddings in a multiresolution 3D hash grid and hence allows a much smaller and faster decoder than the standalone MLP in NeRF [30]. Our network is an implicit surface model because besides the view-dependent color, it predicts the signed distance $s$ from the input position to the nearest surface, instead of predicting volume density as in NeRF. This helps extract watertight geometry from the learned model. We follow the analytical solution in NeuS [40] to convert the signed distance $s$ to density $\sigma$ in or to enable volume rendering for both color and depth:

$$\sigma_i = \alpha \Phi_\beta \left( -s_i \right) \tag{3}$$

where $\Phi_\beta$ is the cumulative distribution function (CDF) of the Laplace distribution with zero mean and $\beta$ scale:

$$\Phi_\beta(x) = \begin{cases} \frac{1}{2} exp \left( \frac{x}{\beta} \right) & \text{if } x \leq 0 \\ 1 - \frac{1}{2} exp \left( -\frac{x}{\beta} \right) & \text{if } x > 0 \end{cases} \tag{4}$$

Both $\alpha$ and $\beta$ are learnable parameters. Under this conversion, the density becomes 0 when $s_i < 0$ (i.e., outside the surface) and becomes $\alpha$ when $s_i \geq 0$ (i.e., inside or on the surface), with the sharpness of drop of the CDF controlled by $1/\beta$. After approximating the density, we use Equation 1 to render RGB color and Equation 2 to render depth.

Our background model takes in the spatial coordinate x and applies the inverse sphere parameterization technique in NeRF++ [34] to x, so that the re-parameterized input becomes bounded even for landscapes and the sky at nearly infinite distances.

The final rendering becomes the composite of foreground and background, equivalent to breaking the integral in Equation 1 (or similarly Equation 2) into two parts:

$$c_{\text{out}} = \sum_{i=1}^{N^{fg}} T_i (1 - e^{-\delta_i \sigma_i}) c_i + \sum_{i=1}^{N^{bg}} T_i (1 - e^{-\delta_i \sigma_i}) c_i \tag{5}$$

where $t_{N^{fg}}$ is both the farthest sampled depth in the foreground model and the closest sampled depth in the background model, and $t_{N^{bg}} = \infty$.

## II. C   Dynamic Object Filtering

We aim at 3D scene reconstruction to aid automated ground truth generation purposes and we focus on the static scene only. We filter independently moving objects in the scene based on annotations obtained from our 3D object detection (3DOD) pipeline which detects vehicles, pedestrians and two wheeled vehicles. Given the detections from the 3DOD pipeline, we can readily avoid any point being sampled by checking if the point is inside any bounding box during sampling and ray marching. To ensure that we use only the detections corresponding to dynamic objects, rather than transient objects that remain static in the current scene (e.g., parked vehicles, construction sites, etc.), we leverage two complementary signals: (1) the relative speed of the bounding box and (2) the difference in the absolute position of the bounding box across frames. The relative speed per bounding box is obtained

by the 3DOD pipeline, which tells whether the annotated object is moving or not relative to the ego vehicle. The difference in the absolute position is obtained by transforming bounding boxes' positions to the world frame using ego vehicle's positions. If the relative speed is zero when the ego is moving, or the relative speed is nonzero when the vehicle is not moving, or the difference in the absolute position is nonzero, we keep the annotation for filtering; otherwise, we regard it as static objects. We demonstrate the effectiveness of dynamic object filtering in Section III.

When we filter out dynamic objects during scene reconstruction, pose refinement becomes a well-posed problem as it would not be in dynamic environments where motion cues are not consistent between ego vehicle and dynamic objects [43]. Therefore, we set the input pose [R|t] to be a learnable parameter that is trained together with grid features and MLP weights. Due to difficulty of optimizing in SO(3), we parameterize the rotation component into unit quaternion $q$ and regress $\Delta q$.

## II. D   Large-Scale Support

The size of the scene covered by each sequence varies, depending on the speed of the ego vehicle. However, the model capacity is predetermined. Therefore, we adopt a simple but effective 'divide and conquer' approach to larger-scale input sequences. With the same rationale as in Block-NeRF [44], we divide the input sequence into subsequences with a pre-fixed number of frames. Then we train our model on each subsequence in parallel as the reconstruction of each is independent. All subsequences still share the same world coordinate so that during rendering, we can readily merge all subsequences based on ego vehicle's position and orientation.

## II. E   Supervision

In addition to the flow, Figure 2 also shows three major supervision signals (only for the foreground model for simplicity). To reconstruct appearance of the scene, we calculate the photometric loss $L_C$ between the full rendered image (i.e., after compositing both foreground and background) and input image. To reconstruct geometry of the scene, we first regularize the learned SDF field by the Eikonal loss $L_S$ (as introduced in Section II. A), and then make use of the LiDAR measurement to calculate the geometry loss $L_D$. Point clouds are first projected onto the camera
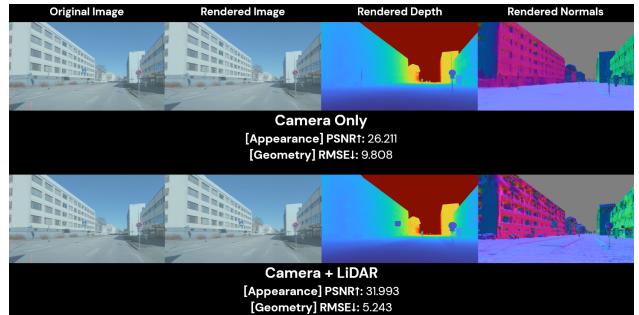


Figure 3: Qualitative and quantitative results that demonstrate benefits of combining LiDAR and camera. PSNR↑: 26.211 without LiDAR and 31.993 with LiDAR. RMSE↓: 9.808 without LiDAR and 5.243 with Li-DAR.

frame and compared against the rendered depth at valid coordinates. The final loss becomes $L = L_C + L_S + L_D$.

## III   EXPERIMENTAL RESULTS

We train and evaluate our method on the challenging internal automotive dataset. We train both foreground and background models jointly for 10,000 iterations with 8196 rays per batch. We use Adam [45] with a learning rate of $1 \times 10^{-2}$. Without large-scale support, we train one sequence on an NVIDIA Tesla V100, with large-scale support we parallelize training using multiple cards.

We first demonstrate the reconstruction result in both bird's eye view (BEV) and ego car's view in Figure 1. The complete video demonstration has been presented at the Electronic Imaging Autonomous Vehicles and Machines conference. In BEV, we extract the foreground mesh from the learned SDF field by the marching cubes algorithm [46], and further apply occlusion culling to it due to the noise in unobserved or occluded regions, such as the rear of the buildings. In the video, we overlay input image at each time and point clouds on top of the mesh. In the ego car's view, we show the rendered image, the rendered depth, and the extracted mesh colored by surface normals.

Figure 3 illustrates qualitative and quantitative results of combining LiDAR and camera in neural rendering for urban scene reconstruction. Peak Signal-to-Noise Ratio (PSNR) and Root Mean Square Error (RMSE) are used to evaluate appearance reconstruction and geometry reconstruction respectively. Both metrics are calculated and averaged over the entire 300-frame sequence. In the camera-only baseline, we follow [41], using off-the-shelf models [47] to estimate depth and surface normal as pseudo ground truth in depth supervision. Quantitatively, we see there's a 23% increase in PSNR and 46% decrease in RMSE when we make use of LiDAR in supervision. Qualitatively, we see the camera-only baseline has missed one of the road signs in all three rendered modalities. Since LiDAR sweeps provide much more reliable depth even on tiny objects, the model is able to capture finer details. On the other hand, off-the-shelf monocular depth estimators still take camera data as input and hence the camera-only model can't recover accurate geometry. Note that the windows on the building are also distinguishable after using LiDAR, because monocular depth estimator can't perceive such detailed geometry cues. We have observed that the surface normals rendered by the camera+LiDAR model are much noisier than the camera-only one, especially on the building. This is because unlike per-pixel depth estimated by the monocular depth estimator, point clouds from LiDAR are much sparser and hence not every rendered depth pixel can be supervised by $L_D$ in Section II. E, which causes am-
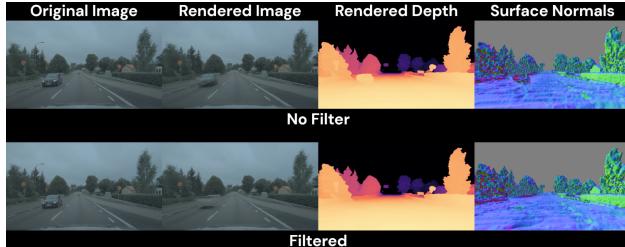


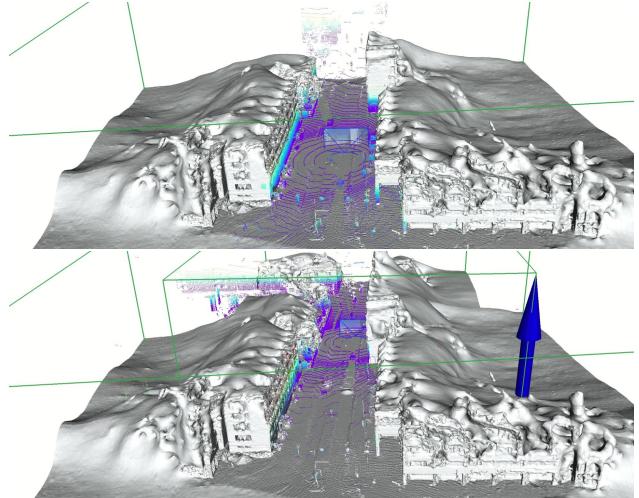Figure 4: Qualitative results of dynamic object filtering



Figure 5: Large-scale support demonstration in BEV. Occlusion culling is not applied to the mesh for simplicity. Green boxes denote allocated spatial size for each subsequence. Top: Extracted mesh of the first subsequence. Bottom: Extracted mesh of the next subsequence and merged into the first subsequence.

biguity in the surface. Further regularization losses can be applied to smoothen the surface in regions with undefined depth.

Figure 4 illustrates the effectiveness of our dynamic object filtering implementation based on 3DOD detections. Without filter, the dynamic object appears as ghosting artifacts in the rendered image. And it's observable in both rendered depth and surface normals, which is challenging for accurate 3D scene reconstruction. By avoiding sampling inside the 3D bounding box during ray marching, the dynamic object becomes completely "invisible" in rendered depth and surface normals. Artifacts also disappear in render image, except that the shadow of the vehicle still remains on the ground because it's not enclosed in the annotation. Some recent works solve this by having an additional prediction head to predict a shadow ratio, and we consider it as a topic to be explored in the future.

Finally we demonstrate the scalability of our solution to larger scenes. As detailed in Section II. D, the sequence shown in Figure 5 has been divided into two subsequences and two models are trained on each subsequence independently. Since they share the same coordinate, all types of renderings can be merged seamlessly. Trained models for each subsequence are queried sequentially according to the input position. Meshes are extracted independently from learned SDF fields and then merged together as shown in the bottom figure.

## IV   Conclusion

In this paper, we introduce our solution to large-scale urban scene reconstruction based on neural rendering. We demonstrate the method's efficacy in reliable neural scene representation by combining neural radiance fields and neural implicit surface. By leveraging LiDAR measurements, fine details in the scene can be captured and reconstructed accurately. Furthermore, we show the method is immune to disturbance of dynamic objects by leverage 3D object detections. We also prove the scalability of the method by reconstructing arbitrarily large environments through

divide and conquer. In future work, we aim to to integrate our solution into our offline perception automated driving stack.

# REFERENCES

[1] S. Mohapatra, S. Yogamani, H. Gotzig, S. Milz, and P. Mader, "Bevdetnet: bird's eye view lidar point cloud based real-time 3d object detection for autonomous driving," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 2809–2815, IEEE, 2021.

[2] H. Rashed, E. Mohamed, G. Sistu, V. R. Kumar, C. Eising, A. El-Sallab, and S. Yogamani, "Fisheyeyolo: Object detection on fisheye cameras for autonomous driving," in *Proceedings of the Machine Learning for Autonomous Driving NeurIPS 2020 Virtual Workshop, Virtual*, vol. 11, 2020.

[3] H. Rashed, E. Mohamed, G. Sistu, V. R. Kumar, C. Eising, A. El-Sallab, and S. Yogamani, "Generalized object detection on fisheye cameras for autonomous driving: Dataset, representations and baseline," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2272–2280, 2021.

[4] S. Chennupati, G. Sistu, S. Yogamani, and S. Rawashdeh, "Auxnet: Auxiliary tasks enhanced semantic segmentation for automated driving," *arXiv preprint arXiv:1901.05808*, 2019.

[5] H. Rashed, A. El Sallab, S. Yogamani, and M. ElHelw, "Motion and depth augmented semantic segmentation for autonomous navigation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 0–0, 2019.

[6] A. Briot, P. Viswanath, and S. Yogamani, "Analysis of efficient cnn design techniques for semantic segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pp. 663–672, 2018.

[7] V. Ravi Kumar, "Cnn based depth map prediction on raw monocular fisheye camera images trained with sparse ground truth."

[8] V. R. Kumar, M. Klingner, S. Yogamani, M. Bach, S. Milz, T. Fingscheidt, and P. Mäder, "Svdistnet: Self-supervised near-field distance estimation on surround view fisheye cameras," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 8, pp. 10252–10261, 2021.

[9] 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 *2018 21st International Conference on Intelligent Transportation Systems (ITSC)*, pp. 2853–2858, 2018.

[10] V. Ravi Kumar, S. A. Hiremath, M. Bach, S. Milz, C. Witt, C. Pinard, S. Yogamani, and P. Mäder, "Fisheyedistancenet: Self-supervised scale-aware distance estimation using monocular fisheye camera for autonomous driving," in *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 574–581, 2020.

[11] V. R. Kumar, M. Klingner, S. Yogamani, S. Milz, T. Fingscheidt, and P. Mader, "Syndistnet: Self-supervised monocular fisheye camera distance estimation synergized with semantic segmentation for autonomous driving," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 61–71, 2021.

[12] 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*, 2021.

[13] A. R. Sekkat, Y. Dupuis, V. R. Kumar, H. Rashed, S. Yogamani, P. Vasseur, and P. Honeine, "Synwoodscape: Synthetic surround-view fisheye camera dataset for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 7, no. 3, pp. 8502–8509, 2022.

[14] M. M. Dhananjaya, V. R. Kumar, and S. Yogamani, "Weather and light level classification for autonomous driving: Dataset, baseline and active learning," in *2021 IEEE International Intelligent Transportation Systems Conference (ITSC)*, pp. 2816–2821, IEEE, 2021.

[15] M. Uricár, J. Ulicny, G. Sistu, H. Rashed, P. Krizek, D. Hurych, A. Vobecky, and S. Yogamani, "Desoiling dataset: Restoring soiled areas on automotive fisheye cameras," in *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*, pp. 0–0, 2019.

[16] A. Das, P. Křížek, G. Sistu, F. Bürger, S. Madasamy, M. Uřičář, V. R. Kumar, and S. Yogamani, "Tiledsoilingnet: Tile-level soiling detection on automotive surround-view cameras using coverage metric," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–6, IEEE, 2020.

[17] S. Shen, L. Kerofsky, and S. Yogamani, "Optical flow for autonomous driving: Applications, challenges and improvements," *Electronic Imaging*, vol. 35, pp. 1–8, 2023.

[18] M. Siam, H. Mahgoub, M. Zahran, S. Yogamani, M. Jagersand, and A. El-Sallab, "Modnet: Moving object detection network with motion and appearance for autonomous driving," *arXiv preprint arXiv:1709.04821*, 2017.

[19] E. Mohamed, M. Ewaisha, M. Siam, H. Rashed, S. Yogamani, W. Hamdy, M. El-Dakdouky, and A. El-Sallab, "Monocular instance motion segmentation for autonomous driving: Kitti instancemotseg dataset and multi-task baseline," in *2021 IEEE Intelligent Vehicles Symposium (IV)*, pp. 114–121, IEEE, 2021.

[20] M. Ramzy, H. Rashed, A. E. Sallab, and S. Yogamani, "Rst-modnet: Real-time spatio-temporal moving object detection for autonomous driving," *arXiv preprint arXiv:1912.00438*, 2019.

[21] N. Tripathi and S. Yogamani, "Trained trajectory based automated parking system using visual slam on surround view cameras," *arXiv preprint arXiv:2001.02161*, 2020.

[22] L. Yahiaoui, J. Horgan, B. Deegan, S. Yogamani, C. Hughes, and P. Denny, "Overview and empirical analysis of isp parameter tuning for visual perception in autonomous driving," *Journal of Imaging*, vol. 5, no. 10, p. 78, 2019.

[23] V. R. Kumar, C. Eising, C. Witt, and S. Yogamani, "Surround-view fisheye camera perception for automated driving: Overview, survey & challenges," *IEEE Transactions on Intelligent Transportation Systems*, 2023.

[24] I. Leang, G. Sistu, F. Bürger, A. Bursuc, and S. Yogamani, "Dynamic task weighting methods for multi-task networks in autonomous driving systems," in *2020 IEEE 23rd International Conference on Intelligent Transportation Systems (ITSC)*, pp. 1–8, IEEE, 2020.

[25] V. R. Kumar, S. Yogamani, H. Rashed, G. Sitsu, C. Witt, I. Leang, S. Milz, and P. Mäder, "Omnidet: Surround view cameras based multi-task visual perception network for autonomous driving," *IEEE Robotics and Automation Letters*, vol. 6, no. 2, pp. 2830–2837, 2021.

[26] H. Rashed, S. Yogamani, A. El-Sallab, P. Krizek, and M. El-Helw, "Optical flow augmented semantic segmentation networks for automated driving," *arXiv preprint arXiv:1901.07355*, 2019.

[27] K. Dasgupta, A. Das, S. Das, U. Bhattacharya, and S. Yogamani, "Spatio-contextual deep network-based multimodal pedestrian detection for autonomous driving," *IEEE transactions on intelligent*

*transportation systems*, vol. 23, no. 9, pp. 15940–15950, 2022.

[28] M. Uricár, D. Hurych, P. Krizek, *et al.*, "Challenges in designing datasets and validation for autonomous driving," in *Proceedings of the International Conference on Computer Vision Theory and Applications*, 2019.

[29] C. Eising, J. Horgan, and S. Yogamani, "Near-field perception for low-speed vehicle automation using surround-view fisheye cameras," *IEEE Transactions on Intelligent Transportation Systems*, vol. 23, no. 9, pp. 13976–13993, 2021.

[30] B. Mildenhall, P. P. Srinivasan, M. Tancik, J. T. Barron, R. Ramamoorthi, and R. Ng, "Nerf: Representing scenes as neural radiance fields for view synthesis," *Communications of the ACM*, vol. 65, no. 1, pp. 99–106, 2021.

[31] H. Turki, D. Ramanan, and M. Satyanarayanan, "Mega-nerf: Scalable construction of large-scale nerfs for virtual fly-throughs," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 12922–12931, 2022.

[32] C. Wang, X. Wu, Y.-C. Guo, S.-H. Zhang, Y.-W. Tai, and S.-M. Hu, "Nerf-sr: High quality neural radiance fields using supersampling," in *Proceedings of the 30th ACM International Conference on Multimedia*, pp. 6445–6454, 2022.

[33] B. Poole, A. Jain, J. T. Barron, and B. Mildenhall, "Dreamfusion: Text-to-3d using 2d diffusion," *arXiv preprint arXiv:2209.14988*, 2022.

[34] K. Zhang, G. Riegler, N. Snavely, and V. Koltun, "Nerf++: Analyzing and improving neural radiance fields," *arXiv preprint arXiv:2010.07492*, 2020.

[35] R. Martin-Brualla, N. Radwan, M. S. Sajjadi, J. T. Barron, A. Dosovitskiy, and D. Duckworth, "Nerf in the wild: Neural radiance fields for unconstrained photo collections," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 7210–7219, 2021.

[36] J. T. Barron, B. Mildenhall, D. Verbin, P. P. Srinivasan, and P. Hedman, "Mip-nerf 360: Unbounded anti-aliased neural radiance fields," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 5470–5479, 2022.

[37] T. Müller, A. Evans, C. Schied, and A. Keller, "Instant neural graphics primitives with a multiresolution hash encoding," *ACM Transactions on Graphics (ToG)*, vol. 41, no. 4, pp. 1–15, 2022.

[38] P. Wang, L. Liu, Y. Liu, C. Theobalt, T. Komura, and W. Wang, "Neus: Learning neural implicit surfaces by volume rendering for multi-view reconstruction," *Advances in Neural Information Processing Systems*, vol. 34, pp. 27171–27183, 2021.

[39] M. Oechsle, S. Peng, and A. Geiger, "Unisurf: Unifying neural implicit surfaces and radiance fields for multi-view reconstruction," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 5589–5599, 2021.

[40] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman, "Volume rendering of neural implicit surfaces," *Advances in Neural Information Processing Systems*, vol. 34, pp. 4805–4815, 2021.

[41] J. Guo, N. Deng, X. Li, Y. Bai, B. Shi, C. Wang, C. Ding, D. Wang, and Y. Li, "Streetsurf: Extending multi-view implicit surface reconstruction to street views," *arXiv preprint arXiv:2306.04988*, 2023.

[42] J. Ost, F. Mannan, N. Thuerey, J. Knodt, and F. Heide, "Neural scene graphs for dynamic scenes," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2856–2865, 2021.

[43] S. Shen, Y. Cai, W. Wang, and S. Scherer, "Dytanvo: Joint refinement of visual odometry and motion segmentation in dynamic environments," in *2023 IEEE International Conference on Robotics and Automation (ICRA)*, pp. 4048–4055, IEEE, 2023.

[44] M. Tancik, V. Casser, X. Yan, S. Pradhan, B. Mildenhall, P. P. Srinivasan, J. T. Barron, and H. Kretzschmar, "Block-nerf: Scalable large scene neural view synthesis," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 8248–8258, 2022.

[45] D. P. Kingma and J. Ba, "Adam: A method for stochastic optimization," *arXiv preprint arXiv:1412.6980*, 2014.

[46] W. E. Lorensen and H. E. Cline, "Marching cubes: A high resolution 3d surface construction algorithm," in *Seminal graphics: pioneering efforts that shaped the field*, pp. 347–353, 1998.

[47] A. Eftekhar, A. Sax, J. Malik, and A. Zamir, "Omnidata: A scalable pipeline for making multi-task mid-level vision datasets from 3d scans," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 10786–10796, 2021.

## AUTHORS BIOGRAPHY

*Shihao Shen* is an machine learning engineer at Qualcomm. He received B.S. in Electrical Engineering from the University of California San Diego and M.S. in Robotic Systems Development from Carnegie Mellon University. His main research focus is machine learning with applications in geometry vision, localization, and 3D reconstruction.

*Louis Kerofsky* is a researcher in video compression, video processing and display. He received M.S. and Ph.D. degrees in Mathematics from the University of Illinois, Urbana-Champaign (UIUC). He has over 20 years of experience in research and algorithm development and standardization of video compression. He has served as an expert in the ITU and ISO video compression standards committees. He is an author of over 40 publications which have over 5000 citations. He is an inventor on over 140 issued US patents. He is a senior member of IEEE, member of Society for Information Display, and member of Association for Computing Machinery.

*Varun Ravi Kumar* holds a staff engineer position in Qualcomm and leads the multi-modal perception team. He received a Ph.D. degree in Artificial Intelligence from TU Ilmenau in 2021 and an M.Sc. degree in 2017 from TU Chemnitz, Germany. Ph.D. thesis builds a first-ever 6-task multi-task learning near-field perception system that constitutes the necessary modules for a Level3 autonomous stack using surround-view fisheye cameras. He has 8+ years of experience in research focusing on designing self-supervised perception algorithms using neural networks for self-driving cars. He is an author of 27 publications with 900 citations and 80+ filed patents.

*Senthil Yogamani* holds an engineering director position at Qualcomm and leads the data-centric AI for autonomous driving department. He has over 18 years of experience in computer vision and machine learning including 15 years of experience in industrial automotive systems. He is an author of 120 publications which have 6500 citations and 150+ filed patents. He serves on the editorial board of various leading IEEE automotive conferences including ITSC and IV.