

# Volumetric Occupancy Mapping With Probabilistic Depth Completion for Robotic Navigation

Marija Popović<sup>\*,1,2</sup>, Florian Thomas<sup>\*,1</sup>, Sotiris Papatheodorou<sup>1</sup>, Nils Funk<sup>1</sup>,  
Teresa Vidal-Calleja<sup>3</sup>, Stefan Leutenegger<sup>1,4</sup>

**Abstract**—In robotic applications, a key requirement for safe and efficient motion planning is the ability to map obstacle-free space in unknown, cluttered 3D environments. However, commodity-grade RGB-D cameras commonly used for sensing fail to register valid depth values on shiny, glossy, bright, or distant surfaces, leading to missing data in the map. To address this issue, we propose a framework leveraging probabilistic depth completion as an additional input for spatial mapping. We introduce a deep learning architecture providing uncertainty estimates for the depth completion of RGB-D images. Our pipeline exploits the inferred missing depth values and depth uncertainty to complement raw depth images and improve the speed and quality of free space mapping. Evaluations on synthetic data show that our approach maps significantly more correct free space with relatively low error when compared against using raw data alone in different indoor environments; thereby producing more complete maps that can be directly used for robotic navigation tasks. The performance of our framework is validated using real-world data.

## I. INTRODUCTION

In recent years, depth sensors have become a core component in a variety of robotic applications, including scene reconstruction, exploration, and inspection. However, commodity-grade RGB-D cameras, such as Microsoft Kinect and Intel RealSense, suffer from limited range and produce images with noise and missing data in view of surfaces that are too shiny, glossy, bright, or simply too far away. In robotic scenarios, this may lead to inefficient and inaccurate mapping performance when only the raw sensor data is used.

This paper studies the problem of *depth completion* applied in the context of *robotic mapping*. Our goal is to create more complete spatial maps of cluttered 3D environments for robotic navigation purposes. This is achieved by filling in holes found in raw depth images that are used for mapping.

Recently, several deep learning-based approaches for depth completion using RGB-D images have been proposed [1, 2] which effectively use colour information to enhance depth. However, propagating the completed depth into robotic frameworks remains an open challenge. A key

This research is supported by the EPSRC ORCA Robotics Hub (EP/R026173/1), EPSRC grant Aerial ABM (EP/N018494/1), Imperial College London (including President’s Scholarship), and SLAMcore Ltd. It is partially funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under Germany’s Excellence Strategy - EXC 2070 - 390732324. \*Equal contribution. <sup>1</sup>Smart Robotics Lab, Department of Computing, Imperial College London. <sup>2</sup>Cluster of Excellence PhenoRob, Institute of Geodesy and Geoinformation, University of Bonn. <sup>3</sup>Centre for Autonomous Systems, Faculty of Engineering and IT, University of Technology Sydney. <sup>4</sup>Smart Robotics Lab, Department of Informatics, Technical University of Munich. mpopovic@uni-bonn.de.

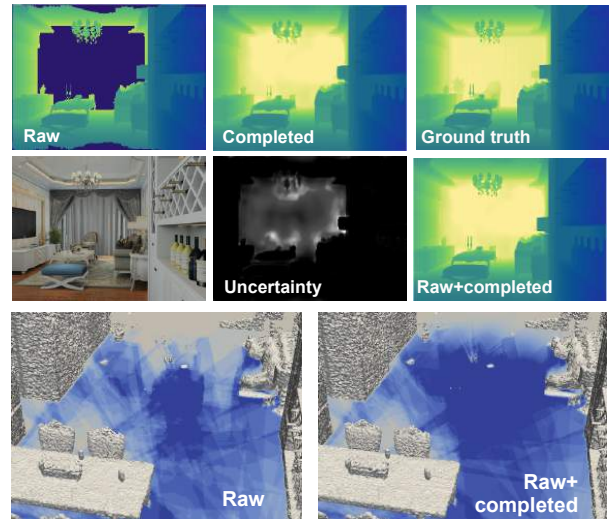


Fig. 1: Overview of our approach for mapping with depth completion. *Top*: Our depth completion network takes raw RGB-D images to predict the completed depth and depth uncertainty, which are used as additional inputs for probabilistic 3D mapping. *Bottom*: By leveraging the network completions to complement the original raw depth data (right), we obtain more complete maps of free space when compared against using the raw depth alone (left). Darker shades of blue indicate areas of lower occupancy probability.

issue is associating the completed areas with reliable measures of depth uncertainty, such that they can be used as an input for probabilistic mapping. Though several works have tackled uncertainty estimation for depth completion, they do not address using this information for 3D reconstruction [3] and largely focus on LiDaR-based sensors in outdoor environments [4–6].

To address this, we propose a new pipeline for mapping with probabilistic depth completion; thus bridging the gap between computer vision algorithms and robotic applications. Inspired by the methods of Huang et al. [2], we introduce a network architecture that jointly predicts both depth and depth uncertainty from RGB-D images by leveraging principles of Bayesian deep learning. Our approach exploits the processed images online as an additional input in the occupancy-based volumetric framework of Vespa et al. [7] and Funk et al. [8]. This procedure enables us to produce maps with more discovered *obstacle-free space* in the environment compared to using the raw images alone, as visualised in [Figure 1](#), which is necessary for robotic navigation tasks in initially unknown environments [8, 9].

The contributions of this work are:

- 1) A new deep learning architecture providing uncertainty

estimates for the probabilistic depth completion of RGB-D images.

- 2) The integration of our network in the volumetric mapping framework of Vespa et al. [7] and Funk et al. [8]. We use the completed depth images with the predicted depth uncertainties in online probabilistic occupancy mapping to obtain more complete free space maps for robotic navigation tasks.
- 3) The extensive evaluation of our framework using synthetic and real-world data showcasing its performance.

We plan to open-source our network implementation for usage and further development by the community.

## II. RELATED WORK

Algorithms for depth estimation and spatial mapping play a key role in many robotic applications and are the subject of a large and growing body of research. In this section, we review previous studies most related to our work.

Traditional methods for depth completion adopt hand-crafted kernels or features to compute the missing values [10, 11]. More recent algorithms [1, 2, 4–6, 12, 13] exploit deep learning for improved performance and generalisation capabilities. Our work focuses on the task of *guided* depth completion, where the goal is to predict the dense depth values at every pixel based on the raw depth and a paired colour image. Uhrig et al. [12] propose a sparse convolution layer which explicitly handles missing data to allow for inputs with varying degrees of sparsity. In a similar problem setup, Ma and Karaman [13] use an encoder-decoder network to combine RGB and depth information within the underlying feature space. Recently, Eldesokey et al. [6] present a network based on normalised convolution layers which supports very sparse depth inputs and also provides confidence measures for the depth predictions. However, the aforementioned studies focus on completing sparse LiDaR-based data in outdoor scenarios and are thus not applicable to the types of degradation obtained with commodity-grade RGB-D cameras, as considered in our work.

For hole-filling with RGB-D cameras, Zhang and Funkhouser [1] exploit the encoder-decoder architecture using dense occlusion boundaries and surface normals predicted from the colour image as secondary features to aid depth completion. Their approach involves an expensive loss optimisation step, making it unsuitable for real-time mapping. Building upon their ideas, Huang et al. [2] introduce a network with a self-attention mechanism and boundary consistency to improve completion accuracy and speed. We propose an extension of their architecture which also predicts the uncertainties in the completed depth.

While significant work has been done on depth completion in the 2D image plane, applying these concepts to 3D mapping in robotics is a relatively unexplored research area. Recently, Teixeira et al. [4] introduced a depth completion algorithm for real-time aerial robotic applications. Similar to us, they obtain probabilistic depth predictions by estimating pixelwise uncertainties. However, they consider LiDaR-based

sensing and do not use the completed images for 3D mapping. Most resembling our work is the approach of Fehr et al. [9], which uses an augmented depth sensor based on sparse inputs for robotic navigation. They show that their system uncovers more free space in unknown environments when compared against using raw depth alone, thereby improving planning performance. Although our work shares the same motivation, a key difference is that, instead of feeding the completed depth directly into a dense mapping framework, we adopt a fully probabilistic strategy based on the depth uncertainties provided by our new modified network.

Uncertainty in depth completion is crucial as it provides a reliability measure for fusing new predicted measurements into the map. One approach is to exploit confidence as a process internal to deep learning [14] to obtain more accurate dense depth outputs, i.e. by leveraging uncertainty as a weight map within the network architecture. An alternative is to treat uncertainty as an auxiliary network output to obtain pixelwise uncertainties [3] or confidence maps [4–6]. We follow the second class of approaches to extract explicit uncertainty values as inputs for mapping. Although, like us, Kendall and Gal [3] learn uncertainty in depth regression problems, to our knowledge, no prior work has applied these ideas in the context of probabilistic robotic mapping.

Another line of work focuses on volumetric scene completion directly in 3D space. For example, Song et al. [15] predict volumetric occupancy and semantic labels from a single-view depth map. Dai et al. [16] complete 3D geometry with per-voxel semantic labels from partial scans. However, as these methods require significant computational processing, they are not viable for real-time, online applications.

## III. APPROACH

In this section, we propose a new approach for tracking and mapping using completed depth images with predicted depth uncertainty. A system overview is depicted in Figure 2. As shown, we process the raw images from a RGB-D camera using a probabilistic depth completion pipeline online to improve the input for Simultaneous Localisation and Mapping (SLAM). Note that, while our approach is applicable for any SLAM scenario, in this paper, we focus on *mapping* only to show improvements for free space mapping in unknown environments. The following sub-sections describe our strategy for probabilistic depth completion before outlining the SLAM framework.

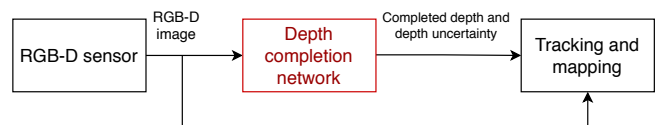


Fig. 2: Overview of our proposed approach. We leverage a network for depth completion with uncertainty to improve the input for probabilistic tracking and mapping.

### A. Network Architecture

Our goal is to complete the depth channel of an RGB-D image and predict the associated pixelwise depth uncertainty,

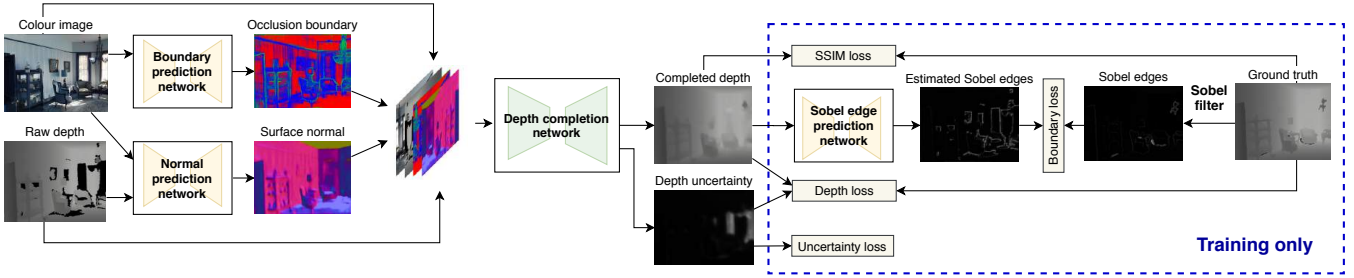


Fig. 3: Our *probabilistic* depth completion system pipeline including the training framework with different training loss components (Section III-B). Given an input RGB-D image, we predict surface normals and boundaries and pass them to the probabilistic depth completion network (Figure 4) to predict depth and associated uncertainty. Black in the depth images indicates missing information.

to be used as input for probabilistic tracking and mapping. To achieve this, we develop a pipeline based on the depth completion network proposed by Huang et al. [2]. An overview of the depth completion sub-system is shown in Figure 3. Figure 4 details our new network architecture.

The main features of the depth completion network are the use of a self-attention mechanism and boundary consistency to produce depth maps with high quality and structure. Following Zhang and Funkhouser [1], we predict surface normals and occlusion boundaries from the raw RGB-D image and use them as additional input features to the network. To estimate surface normals, we employ the hierarchical RGB-D fusion network of Zeng et al. [17], which has state-of-the-art performance. Our boundary estimation network is based on the approach of Zhang et al. [18], using only RGB channels. The normals and boundaries are concatenated with the raw RGB-D image and used for the learning task.

To predict depth uncertainty, we leverage the Bayesian deep learning concepts of Kendall and Gal [3]. Our key idea is to introduce a second decoder on the network output to learn the mapping to the input uncertainty in the completed depth. This is illustrated in the bottom branch of the architecture in Figure 4. We use a SoftPlus activation function (purple) to constrain the output uncertainty to be non-negative. By using a network with two different branches, the encoder of the original network captures the latent features common to the completed output depth and associated uncertainty, before they are processed separately to account for individual information. Moreover, we increase the number of channels per layer to 64 from 48 in the original network to provide a larger latent space for learning in the dual prediction task.

### B. Loss Function

We assume a Gaussian likelihood to model our aleatoric uncertainty [3]. Our loss function for depth completion with uncertainty is then the weighted sum of errors:

$$\begin{aligned}
 L = \frac{1}{N} \sum_{p=1}^N \frac{1}{\sigma(\mathbf{x}_p)^2} |y_p - f(\mathbf{x}_p)|^2 + \log(\sigma(\mathbf{x}_p)^2) \\
 + \lambda_{BC} |y_{Sobel,p} - f_{Sobel}(f(\mathbf{x}_p))| \\
 + \lambda_{SSIM} SSIM(y_p, f(\mathbf{x}_p)), \quad (1)
 \end{aligned}$$

where  $N$  is the number of pixels  $p$  in an image,  $\mathbf{x}_p$  is the input vector of features for the depth completion network

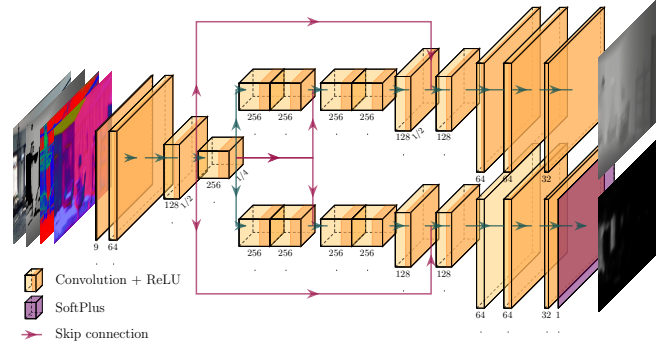


Fig. 4: Our architecture for depth completion with uncertainty. We extend the network of Huang et al. [2] with a second output decoder for uncertainty prediction. Our network takes as input raw RGB-D, surface normals and boundaries (left), and outputs completed depth (top-right) and pixelwise uncertainty (bottom-right).

(raw depth, RGB, estimated normals and boundaries),  $y_p$  is the ground truth depth,  $f(\cdot)$  and  $\sigma(\cdot)$  are the completed depth and associated uncertainty output by the network, respectively, following directly from the negative log likelihood assuming Gaussian uncertainty, and  $\lambda_{BC}$  and  $\lambda_{SSIM}$  are tunable parameters. We only consider pixels in areas *observed*, i.e. non-zero, in the ground truth depth image.

Following Huang et al. [2], we also include a boundary related loss (third term) to enforce boundary consistency and a structural related loss based on the Structural Similarity Index (SSIM) measure [19] to reduce distortion and enhance structural quality (fourth term). The former is computed by training a model to learn the Sobel edges associated with the completed depth supervised by those computed from the ground truth depth. The different components of the loss function are depicted in the dashed blue box in Figure 3.

### C. Mapping

We use an occupancy map to model the environment, as this representation is suitable for integrating noisy sensor measurements and explicitly captures free space for robotic planning applications. Our approach leverages the multi-resolution occupancy mapping (‘MultiresOFusion’) and dense volumetric SLAM framework from Funk et al. [8]. This pipeline is an extension of *supereight* [7] that enables integrating data at multiple octree levels and explicitly maps free space, while performing significantly better than other occupancy mapping frameworks, as shown in [8].

To explain the role of depth uncertainty for mapping in our approach, we briefly overview the ‘MultiresOFusion’ probabilistic inverse sensor model used to fuse new depth measurements into the map. The inverse sensor model, inspired by Loop et al. [20], uses a piecewise linear function in log-odds space. The model produces probabilities expressed in log-odds directly to match the representation of occupancy probabilities in the map. Given a noisy depth measurement  $z$ , we assume its standard deviation is:

$$\sigma(z) = \min(\max(k_\sigma z^2, \sigma_{\min}), \sigma_{\max}), \quad (2)$$

as shown in the right plot in Figure 5, where  $k_\sigma$ ,  $\sigma_{\min}$  and  $\sigma_{\max}$  are constants. This corresponds to a triangulation-based depth camera noise model. The inverse sensor model is used to compute the log-odds occupancy probability given the distance  $d_r$  from a query point to the measured depth  $z_r$  along the ray as shown in the left plot. Log-odd values in front of the surface are clipped at  $l_{\min}$  reached at  $\mu = 3\sigma$  and gradually increases with distance, peaking halfway through the surface thickness  $\tau(z)$ . Surface thickness is computed as:

$$\tau(z) = \min(\max(k_\tau z, \tau_{\min}), \tau_{\max}), \quad (3)$$

where  $k_\tau$ ,  $\tau_{\min}$  and  $\tau_{\max}$  are constants. No voxels beyond  $z_r + \tau$  from the camera are updated. Larger values of  $\sigma$  result in a more gradual increase of occupancy probability. For fusing multiple measurements, we use a clamped accumulation log-odd occupancy as described by Vespa et al. [7].

Our aim is to exploit the completed depth and depth uncertainty provided by our network to complement the raw sensor data captured by this model and improve mapping performance. Specifically, we propose using the network-generated depth uncertainty instead of the measurement uncertainty above for completed depth areas.

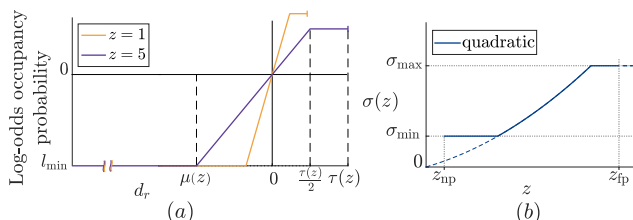


Fig. 5: (a) Inverse sensor model for fusing new data into a map. Occupancy probability as a function of the difference  $d_r$  from a query point to depth measured along a ray. (b) Measurement uncertainty model for mapping with raw depth in *supereight* ‘MultiresOFusion’ [8]. Standard deviation  $\sigma$  given depth measurement  $z$ .

#### IV. EXPERIMENTAL RESULTS

This section presents our experimental results. First, we validate our probabilistic depth completion network via an ablation study. We then evaluate our 3D mapping framework in indoor environments using synthetic and real-world data.

##### A. Training Procedure

We trained our depth completion system end-to-end on Matterport3D, a large-scale RGB-D dataset [21] representative of an indoor exploration task. For training, (complete) ground truth depth was obtained from Zhang and Funkhouser

[1] based on multi-view reconstruction. Unless otherwise specified, 129816 and 36252 images from the dataset were used for training and testing, respectively, with a resolution of  $320 \text{ px} \times 240 \text{ px}$ . We used the Adam optimiser with a weight decay of  $10^{-3}$ , a learning rate of  $10^{-5}$ , and set  $\lambda_{SSIM} = \lambda_{BC} = 1$  in Equation (1). The models were implemented in PyTorch and training was done on a NVIDIA GeForce RTX 2080 Ti GPU with a 3.9GHz AMD Ryzen 9 3900X CPU. On this machine, one forward pass through the pipeline takes  $\sim 0.2 \text{ s}$ .

##### B. Ablation Study

Our first aim is to evaluate our new two-decoder network for depth completion with uncertainty. To this end, an ablation study is conducted to investigate the benefits of separating the two outputs in the proposed architecture and training the model end-to-end for both depth completion and uncertainty. We compare: (i) our proposed architecture with two decoders (Figure 4); (ii) the original architecture of Huang et al. [2] simply extended with a single shared output decoder for depth and uncertainty; (iii) a smaller variant of our architecture, using 48 channels per layer as in [2] instead of 64; and (iv) the same architecture as in (iii), but freezing the weights of the encoder and depth completion decoder parts using the trained model in [2], such that only the uncertainty output features are learned. Apart from (iv), we initialised each model with random weights, then let it train and report the best epoch. We also experimented with initialising (iii) using the pretrained weights from (iv) for training and confirmed that this has no significant impact on the final results, but does speed up convergence.

To evaluate prediction accuracy, we consider the Root Mean Squared Error (RMSE), the Mean Absolute Error (MAE), the percentages of predicted pixels  $\mathbf{p}_{\text{pred}}$  within an interval  $\delta = \frac{|\mathbf{p}_{\text{pred}} - \mathbf{p}_{\text{true}}|}{\mathbf{p}_{\text{true}}}$ , where  $\mathbf{p}_{\text{true}}$  are the corresponding pixel values obtained from the ground truth image and  $\delta \in \{1.05, 1.10, 1.25, 1.25^2, 1.25^3\}$ , and the SSIM. To evaluate the quality of the uncertainty, we measure the Area Under the Sparsification Error curve (AUSE). As explained by Ilg et al. [22], this metric captures the correlation between the estimated uncertainty and prediction error. Our AUSE values are computed based on all pixels in the test set.

Table I summarises our results. With the lowest AUSE, the single-decoder network produces the best uncertainty measure at the cost of lowest prediction accuracy. In contrast, training *only* a second decoder yields low error, but poor uncertainty, as the shared encoder weights are fixed and optimised for the depth completion problem. By increasing the model latent space and training the network end-to-end, our proposed model obtains relatively low AUSE without compromising prediction accuracy with respect to the original implementation. For visual validation, Figure 6 presents example results of our proposed larger 2-decoder model on the Matterport3D test set. This way, we achieve high-quality completed depth with reliable, and, importantly, consistent uncertainty estimates, which we exploit in the next sub-sections to improve probabilistic mapping performance.

Model	RMSE (m) ↓	AUSE ↓	SSIM ↑	MAE (m) ↓	1.05(%) ↑	1.10(%) ↑	1.25(%) ↑	1.25 <sup>2</sup> (%) ↑	1.25 <sup>3</sup> (%) ↑
(i) 2 decoders (64) (ours)	<b>0.3154</b>	0.1849	<b>0.9054</b>	0.1282	<b>85.34</b>	<b>90.01</b>	93.92	96.55	97.85
(ii) 1 decoder (48)	0.3484	<b>0.1771</b>	0.8820	0.1538	80.19	85.08	90.55	95.16	97.22
(iii) 2 decoders (48)	0.3187	0.2115	0.8994	0.1335	85.22	89.61	93.52	96.29	97.68
(iv) 2 decoders (48), depth weights from [2]	0.3166	0.1996	0.9034	<b>0.1268</b>	80.30	88.69	<b>94.20</b>	<b>96.89</b>	<b>98.02</b>

TABLE I: Comparison of our proposed 2-decoder, 64-channel network for depth completion with uncertainty (top) against benchmarks derived from Zhang and Funkhouser [1] on the Matterport3D test set. Our architecture achieves good depth uncertainty measures without compromising depth prediction accuracy. The number of channels per layer in the network is in parentheses.

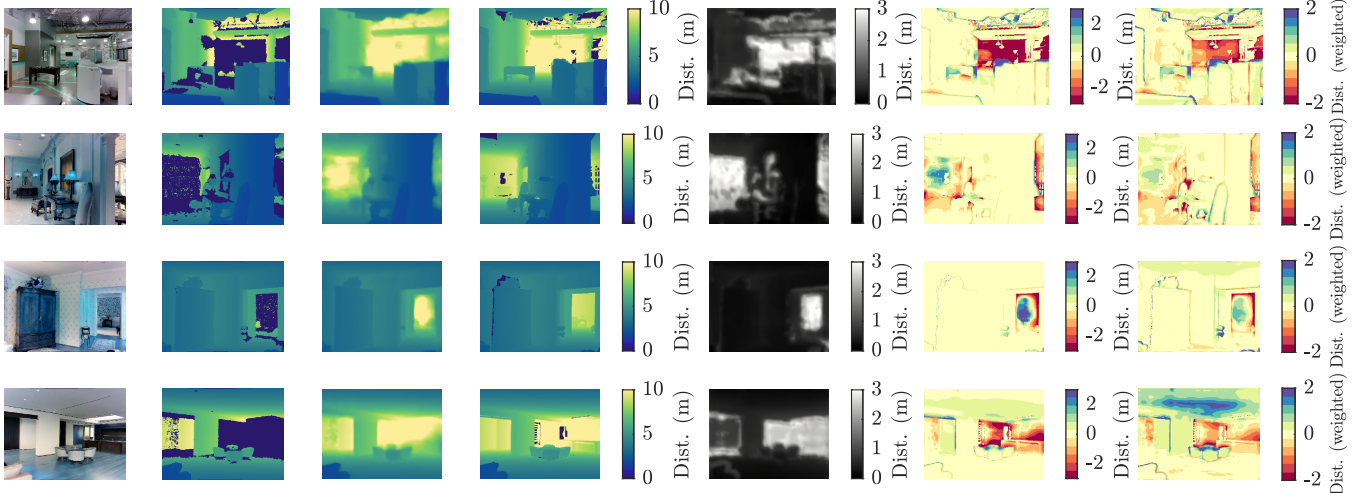


Fig. 6: Examples of our proposed 2-decoder, 64 channel per layer probabilistic depth completion network outputs on images from different sequences from the Matterport3D test set. *Left to right*: RGB, raw depth, completed depth, ground truth depth (obtained from [1]), depth uncertainty (standard deviation), depth error, depth error weighted by the standard deviation. We validate that our network completes the holes in the raw depth images and provides higher uncertainty estimates in areas where data is missing.

### C. Evaluation on Synthetic Data

We perform a quantitative evaluation of our approach for occupancy mapping using trajectory sequences from the synthetic RGB-D dataset InteriorNet [23], in which ground truth depth images and pose data are available. Our aim is to show that mapping using the network predicted depth and uncertainty leads to more complete final maps and a greater volume of free space discovered in the environment, which is a key requirement for safe robotic planning and navigation.

The ground truth depth from InteriorNet contains no measurement noise. To simulate realistic noisy depth images, we degrade the ground truth following the quadratic noise model for the Kinect sensor developed by Nguyen et al. [24]:

$$\sigma_z(z) = 0.0012 + 0.0019(z - 0.4)^2, \quad (4)$$

where  $\sigma_z$  is the standard deviation of lateral noise in metres at a pixel and  $z$  is the corresponding depth measurement in metres. Additionally, a Gaussian filter with a standard deviation of 0.5 px is applied on the depth image to blur the noise between adjacent pixels.

The ground truth depth does not contain occlusion holes or missing depth measurements. To create missing data for depth completion, we set practical sensing limits of 0.8 m–6 m beyond which depth readings are zero. To generate occlusion holes and remove measurements based on the context and structure of the scene within this range, e.g. on textureless/reflective surfaces, we use the generative

adversarial framework of Atapour-Abarghouei et al. [25], which learns to predict depth holes from RGB. This network is trained on our Matterport3D training set, since it contains real RGB/raw depth image pairs. We then train CycleGAN, an unpaired image-to-image translation network [26] to learn the visual domain shift between Matterport3D and InteriorNet using 10000 random images from each dataset, and apply the hole predictor on InteriorNet RGB images translated to the Matterport3D style. This procedure allows us to transfer the learnt structures of real holes to the synthetic dataset and thus generate realistically corrupted depth images.

For depth completion, we use our trained 2-decoder, 64-channel network from Section IV-B, further fine-tuned on the corrupted InteriorNet depth data with 110000 and 28000 images for training and testing, respectively. The images are downsampled to 320 px × 240 px and we apply the same optimisation algorithm as detailed in Section IV-A.

For spatial mapping, we use *supereight* ‘MultiresO-Fusion’ [8] with a voxel resolution of 0.0146 m in a 15 m × 15 m × 15 m volume. We use the inverse sensor model in Figure 5 to integrate raw depth data into the map, setting the constants  $k_\tau$ ,  $\tau_{\min}$ , and  $\tau_{\max}$  as 0.026, 0.06 m, and 0.16 m, respectively. To capture measurement uncertainty, we consider the quadratic uncertainty model in Equation (2) with  $k_\sigma = 0.0016$  m,  $\sigma_{\min} = 0.005$  m, and  $\sigma_{\max} = 0.02$  m resulting in  $\sigma_r = 0.005$  m at  $z_r = 1$  m for the raw depth.

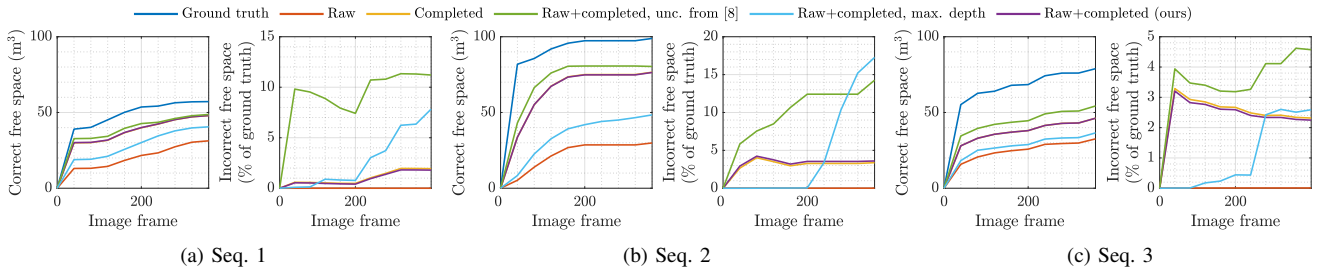


Fig. 7: Comparison of correct and incorrect free space volume discovered during three trajectory sequences from InteriorNet using different depth and uncertainty inputs for mapping. Mapping with probabilistic depth completion leads to more correct free space mapped throughout the image sequence with relatively small errors. Note the different scales on the  $y$  axes for correct and incorrect free space.

We map three trajectories from different InteriorNet scenes not present in the training set<sup>1</sup>. We use 400 images for mapping per sequence, picking large rooms with wide ranges of motion to highlight the advantages of applying depth completion when data is missing. Our experiments compare: (i) the raw depth images with a quadratic depth uncertainty as given in Equation (2) (R); (ii) the completed images with the predicted depth uncertainty from our network (C); and a combination of the two (R+C), where the raw depth is used in known areas, and the invalid depth pixels are completed. As baselines with completion, we study (iii) using the quadratic sensor model in *both* the raw and network completed areas (R+C, unc. from [8]) and (iv) simply filling in the invalid areas with the maximum camera range (8 m) and depth uncertainty  $3\sigma_r$ , conservatively set to this value (R+C, max. depth). Our proposed approach (v) is using the detailed raw depth with the quadratic sensor model in valid areas and depth completion with the network predicted uncertainty to fill in the rest (R+C (ours)). This method is depicted in the top images in Figure 1 and the system diagram in Figure 2. In (ii) and (v), for pixels where the network-generated depth is used, we compare the network depth uncertainty with the quadratic uncertainty model (Equation (2)) when using the network completed depth. If the network depth uncertainty is more than 2 times greater than the quadratic uncertainty *only* free space is integrated for this pixel, otherwise normal integration is performed. This prevents us from creating incorrect surfaces for very uncertain completions while still obtaining usable probabilistic free space estimates.

Our evaluation metrics are the volumes of correct and incorrect mapped free space in the environment with respect to the map generated with ground truth depth at a given image frame. Voxels with occupancy probability  $< 0.04\%$  are considered to be free in the reconstructions; for ground truth, we use a less conservative threshold of  $< 3\%$ . These thresholds can be tuned to balance between extra free space exploration and false-positive free space in a given scenario. To measure accuracy in the final reconstructions, we create meshes using marching cubes, and compute the average distance from the ground truth mesh to an output mesh.

The evaluation results are summarised in Table II. Using depth completion in our mapping pipeline leads to remark-

Seq.	Method	Correct free space (m <sup>3</sup> ) ↑	Incorrect free space (m <sup>3</sup> ) ↓	Mesh accuracy (m) ↓
1	R	31.6748	0.0092	0.0891
	C	48.2372	1.1175	0.1652
	R+C, unc. from [8]	48.9464	6.4669	1.0337
	R+C, max. depth	40.8428	4.5132	0.3379
	R+C (ours)	48.2307	1.0351	0.2008
2	R	32.1136	0.0066	0.0884
	C	77.9038	3.7560	0.2073
	R+C, unc. from [8]	81.5002	15.6306	1.1153
	R+C, max. depth	47.8514	17.4437	0.2942
	R+C (ours)	77.7847	3.9155	0.2939
3	R	34.0504	0.0070	0.0704
	C	47.6763	1.8470	0.0885
	R+C, unc. from [8]	55.6404	3.6456	0.6194
	R+C, max. depth	37.8850	2.0651	0.3551
	R+C (ours)	47.5621	1.7918	0.1725

TABLE II: Evaluation of different depth and uncertainty inputs for mapping using three sequences from InteriorNet. ‘R’, ‘C’ and ‘R+C’ represent raw depth, completed depth, and a combination of the two. Depth completion enables mapping more free space without significantly compromising reconstruction accuracy.

ably more discovered free space volume since, thanks to the predictions, we can capture the free space associated with *all* pixels, instead of only those with valid raw depth measurements. The two ‘R+C’ baselines produce high volumes of incorrect free space with their simple heuristics. In contrast, using our probabilistic network output yields much smaller inaccuracies relative to the gain in correct free space. Our proposed combined approach (‘R+C’) has the benefit of preserving real, detailed raw depth where it is available. Finally, though the reconstruction accuracy with completion is slightly worse when compared to using the raw depth alone (‘R’), it is not drastically degraded with respect to the size of the rooms. We emphasise here that our aim is *not* to achieve higher-quality fine-scale surface reconstructions, but rather create a free space map of the environment suitable for motion planning while preserving its global structure.

Figure 7 depicts the evolution of the mapped free space during the three sequences using different mapping strategies. As a qualitative result, Figure 8 illustrates occupancy map cross-sections at image frame 160 of Seq. 1. The plots in Figure 7 verify that the completion methods consistently map significantly more free space compared to using only the raw depth (orange), even with a small number of images, while free space error is relatively low and grows slowly.

<sup>1</sup>Seq. 1: ‘3FO4K9H4NDAO (7)’; Seq. 2: ‘3FO4JVHRJC4T (7)’; Seq. 3: ‘3FO4JXILITSO (7)’. Trajectory numbers are given in the parentheses.

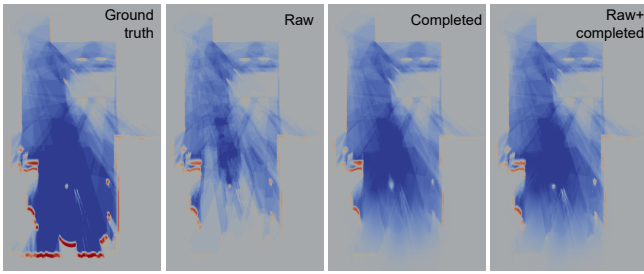


Fig. 8: Comparison of occupancy map cross-sections for Seq. 1 of InteriorNet at image frame 160. Blue to red colours encode occupancy probability (cream is unknown). Depth completion enables mapping more free space (blue) on the side of the room opposite the depth camera (bottom), which is outside the measurement range.

Figure 8 depicts visually the greater proportion of free space (blue areas) achieved using depth completion, especially on the side of the room away from the depth camera (bottom). This portrays the benefit of using our framework in large environments where raw depth coverage is limited.

The bottom images of Figure 1 show the occupancy map cross-sections from Figure 8 overlaid on the output meshes obtained using the ‘R’ and ‘R+C’ approaches. We confirm that the free space in the room is much more complete using our depth completion pipeline. As expected, reconstruction quality remains visually similar; our strategy for integrating free space with highly uncertain completions prevents creating artefacts which may compromise navigation safety.

#### D. Evaluation on Real-World Data

We demonstrate our pipeline for occupancy mapping with probabilistic depth completion using five sequences from the TUM RGB-D dataset [27], which contain trajectory ground truth and RGB-D images captured with a Microsoft Kinect. Note that TUM RGB-D does not include ground truth depth for evaluating accuracy as in Section IV-C. Instead, the aim is to validate qualitatively the benefit of using our trained depth completion system to map free space using real images.

For mapping, we use *supereight* ‘MultiresOFusion’ with a 0.0146 m voxel resolution in a 15 m × 15 m × 15 m volume. The constants  $k_\tau$ ,  $\tau_{\min}$ , and  $\tau_{\max}$  are 0.05, 0.06 m, and 0.16 m, respectively, and the quadratic uncertainty model in Equation (2) uses  $k_\sigma = 0.0025$  m,  $\sigma_{\min} = 0.0098$  m, and  $\sigma_{\max} = 0.0294$  m. These parameters correspond to those used by Funk et al. [8] to evaluate ‘MultiresOFusion’ on real-world data. For depth completion, we use our 2-decoder, 64-channel network from Section IV-B trained on Matterport3D with 320 px × 240 px images. Note that the weights finetuned on InteriorNet in Section IV-C produced similar results.

We compare mapping using: raw depth with the quadratic sensor model in Equation (2) (denoted by ‘R’ in Section IV-C) and our proposed combined approach (‘R+C’), using the network-generated completed depth and depth uncertainty in invalid raw depth areas. Examples of probabilistic depth completion can be seen in Figure 9. As described in Section IV-C, for our proposed method, we integrate only the free space for completed pixels where the network depth uncertainty is more than 2 times greater than the quadratic

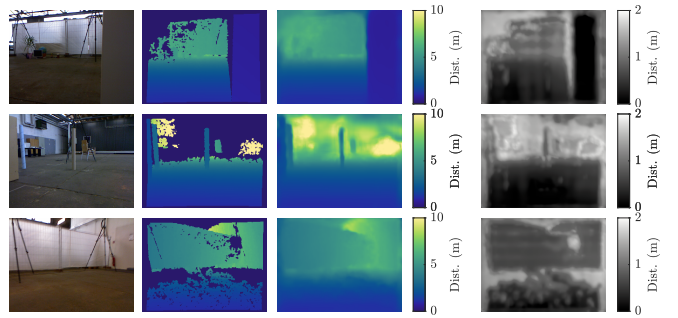


Fig. 9: Examples of our probabilistic depth completion network outputs on various images from TUM RGB-D fr2/pioneer\_slam2. Left to right: RGB, raw depth, completed depth, depth uncertainty (standard deviation). Our network completes holes in the raw depth and provides valid uncertainty estimates for free space mapping.

Sequence	Method	Discovered free space volume (m <sup>3</sup> ) at sequence completion (%)			
		25 (%)	50 (%)	75 (%)	100 (%)
fr1/ 360	R	21.78	60.21	78.41	80.48
	R+C	22.41	62.10	82.68	84.63
fr2/ pioneer_slam	R	108.25	239.18	289.55	306.33
	R+C	170.08	313.39	368.51	382.46
fr2/ pioneer_slam2	R	93.28	159.02	245.67	273.20
	R+C	103.50	202.45	326.58	354.35
fr2/ pioneer_slam3	R	87.40	116.11	212.29	268.37
	R+C	142.21	185.44	291.43	348.24
fr2/ pioneer_360	R	81.21	203.56	264.14	268.29
	R+C	149.52	272.98	339.95	350.31

TABLE III: Comparison of free space volumes (m<sup>3</sup>) discovered during five sequences from TUM RGB-D using different depth and uncertainty inputs for mapping. Using the probabilistic network completed depth to complement raw depth (‘R+C’) yields faster free space mapping compared to using the raw depth alone (‘R’).

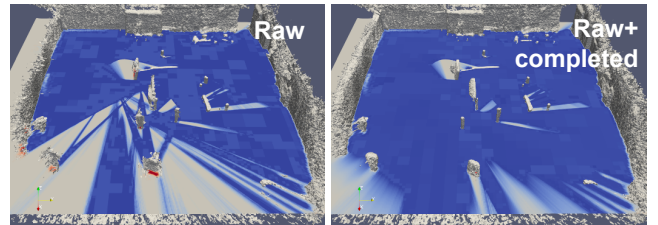


Fig. 10: Comparison of final occupancy map cross-sections and 3D meshes for TUM RGB-D fr2/pioneer\_360. Blue to red colours encode occupancy probability. Depth completion uncovers more free space while the surface reconstruction remains similar.

uncertainty of the raw sensor model. As before, we measure free voxels based on an occupancy probability of < 0.04%.

Table III shows the free space mapped during the sequences. Similar to our InteriorNet experiments, using probabilistic depth completion for mapping consistently yields faster free space discovery when compared against the raw data alone. The occupancy cross-sections in Figure 10 depict visually more free space (blue) at the end of the sequence using completion while the final 3D reconstruction remain similar and do not compromise global navigation safety. These results validate our pipeline in real-world settings.

## V. CONCLUSIONS AND FUTURE WORK

This paper introduced a framework for volumetric mapping using depth completion with uncertainty. A core component of our pipeline is a new network architecture for jointly predicting missing depth and depth uncertainty based on images from commodity-grade RGB-D cameras in cluttered indoor environments. The probabilistic depth is used as an input for mapping to complement the raw depth images, allowing us to obtain more complete free space maps.

We performed an ablation study to validate our network for depth completion with uncertainty. The integrated system for mapping with probabilistic depth was evaluated using synthetic RGB-D data. Our proposed approach using both raw and completed depth was shown to discover correct free space most rapidly without compromising map accuracy and safety in terms of false positive free space. This property is crucial for robotic planning in unknown environments. Further tests validate our approach using real-world images.

One limitation is that our network completions are over-smoothed on depth discontinuities due to the high uncertainties present in these regions. Though our combined approach mitigates this issue by using raw depth data, in more complex environments, one could exploit the edge predictions available from training to preserve sharp boundaries. Another idea is to use recurrent networks to ensure consistency in depth prediction between consecutive images. Finally, we will extend our framework to active mapping problems.

## ACKNOWLEDGEMENT

We would like to thank Xiao Gu for his advice on simulating realistic depth cameras and Dr. Željko Popović and Alexis Laignelet for helping with our experimental setup.

## REFERENCES

- [1] Y. Zhang and T. Funkhouser, "Deep Depth Completion of a Single RGB-D Image," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2018, pp. 175–185.
- [2] Y. K. Huang, T. H. Wu, Y. C. Liu, and W. H. Hsu, "Indoor Depth Completion with Boundary Consistency and Self-Attention," in *International Conference on Computer Vision*. IEEE, 2019, pp. 1070–1078.
- [3] A. Kendall and Y. Gal, "What uncertainties do we need in Bayesian deep learning for computer vision?" in *Advances in Neural Information Processing Systems*. Curran Associates, 2017, pp. 5575–5585.
- [4] L. Teixeira, M. R. Oswald, M. Pollefeys, and M. Chli, "Aerial Single-View Depth Completion With Image-Guided Uncertainty Estimation," *IEEE Robotics and Automation Letters*, vol. 5, no. 2, pp. 1055–1062, 2020.
- [5] A. Eldesokey, M. Felsberg, K. Holmquist, and M. Persson, "Uncertainty-Aware CNNs for Depth Completion: Uncertainty from Beginning to End," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2020.
- [6] A. Eldesokey, M. Felsberg, and F. S. Khan, "Confidence Propagation through CNNs for Guided Sparse Depth Regression," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 42, no. 10, pp. 2423–2436, 2020.
- [7] E. Vespa, N. Nikolov, M. Grimm, L. Nardi, P. H. J. Kelly, and S. Leutenegger, "Efficient Octree-Based Volumetric SLAM Supporting Signed-Distance and Occupancy Mapping," *IEEE Robotics and Automation Letters*, vol. 3, no. 2, pp. 1144–1151, 2018.
- [8] N. Funk, J. Tarrío, S. Papatheodorou, M. Popović, P. F. Alcantarilla, and S. Leutenegger, "Multi-resolution 3D mapping with explicit free space representation for fast and accurate mobile robot motion planning," *arXiv*, 2020.
- [9] M. Fehr, T. Taubner, Y. Liu, R. Siegwart, and C. Cadena, "Predicting Unobserved Space for Planning via Depth Map Augmentation," in *International Conference on Advanced Robotics*, 2019, pp. 30–36.
- [10] D. Ferstl, C. Reinbacher, R. Ranftl, M. Ruether, and H. Bischof, "Image Guided Depth Upsampling Using Anisotropic Total Generalized Variation," in *IEEE International Conference on Computer Vision*, 2013, pp. 993–1000.
- [11] H. Xue, S. Zhang, and D. Cai, "Depth Image Inpainting: Improving Low Rank Matrix Completion With Low Gradient Regularization," *IEEE Transactions on Image Processing*, vol. 26, no. 9, pp. 4311–4320, 2017.
- [12] J. Uhrig, N. Schneider, L. Schneider, U. Franke, T. Brox, and A. Geiger, "Sparsity Invariant CNNs," in *International Conference on 3D Vision*, 2017, pp. 11–20.
- [13] F. Ma and S. Karaman, "Sparse-to-Dense: Depth Prediction from Sparse Depth Samples and a Single Image," 2018.
- [14] J. Qiu, Z. Cui, Y. Zhang, X. Zhang, S. Liu, B. Zeng, and M. Pollefeys, "DeepLiDAR: Deep Surface Normal Guided Depth Prediction for Outdoor Scene From Sparse LiDAR Data and Single Color Image," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2019.
- [15] S. Song, F. Yu, A. Zeng, A. X. Chang, M. Savva, and T. Funkhouser, "Semantic Scene Completion from a Single Depth Image," *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [16] A. Dai, D. Ritchie, M. Bokeloh, S. Reed, J. Sturm, and M. Nießner, "Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans," in *IEEE Conference on Computer Vision and Pattern Recognition*, 2018.
- [17] J. Zeng, Y. Tong, Y. Huang, Q. Yan, W. Sun, J. Chen, and Y. Wang, "Deep Surface Normal Estimation with Hierarchical RGB-D Fusion," in *IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, 2019.
- [18] Y. Zhang, S. Song, E. Yumer, M. Savva, J.-Y. Lee, H. Jin, and T. Funkhouser, "Physically-based rendering for indoor scene understanding using convolutional neural networks," *IEEE Conference on Computer Vision and Pattern Recognition*, 2017.
- [19] Zhou Wang, A. C. Bovik, H. R. Sheikh, and E. P. Simoncelli, "Image quality assessment: from error visibility to structural similarity," *IEEE Transactions on Image Processing*, vol. 13, no. 4, pp. 600–612, 2004.
- [20] C. Loop, Q. Cai, S. Orts-Escolano, and P. A. Chou, "A Closed-Form Bayesian Fusion Equation Using Occupancy Probabilities," in *International Conference on 3D Vision*, 2016, pp. 380–388.
- [21] A. Chang, A. Dai, T. Funkhouser, M. Halber, M. Niessner, M. Savva, S. Song, A. Zeng, and Y. Zhang, "Matterport3d: Learning from rgb-d data in indoor environments," *International Conference on 3D Vision*, 2017.
- [22] E. Ilg, O. Cicek, S. Galesso, A. Klein, O. Makansi, F. Hutter, and T. Brox, "Uncertainty estimates and multi-hypotheses networks for optical flow," in *European Conference on Computer Vision*, 2018, pp. 652–667.
- [23] W. Li, S. Saeedi, J. McCormac, R. Clark, D. Tzoumanikas, Q. Ye, Y. Huang, R. Tang, and S. Leutenegger, "InteriorNet: Mega-scale Multi-sensor Photo-realistic Indoor Scenes Dataset," in *British Machine Vision Conference*, 2018.
- [24] C. V. Nguyen, S. Izadi, and D. Lovell, "Modeling Kinect Sensor Noise for Improved 3D Reconstruction and Tracking," in *International Conference on 3D Imaging, Modeling, Processing, Visualization Transmission*, 2012, pp. 524–530.
- [25] A. Atapour-Abarghouei, S. Akcay, G. Payen de La Garanderie, and T. P. Breckon, "Generative adversarial framework for depth filling via Wasserstein metric, cosine transform and domain transfer," *Pattern Recognition*, vol. 91, pp. 232 – 244, 2019.
- [26] J.-Y. Zhu, T. Park, P. Isola, and A. A. Efros, "Unpaired Image-to-Image Translation using Cycle-Consistent Adversarial Networks," in *IEEE International Conference on Computer Vision*, 2017.
- [27] J. Sturm, N. Engelhard, F. Endres, W. Burgard, and D. Cremers, "A Benchmark for the Evaluation of RGB-D SLAM Systems," in *IEEE/RSJ International Conference on Intelligent Robots and Systems*. IEEE, 2012, pp. 573–580.
- [28] J. C. Chow and D. D. Lichti, "Photogrammetric bundle adjustment with self-calibration of the PrimeSense 3D camera technology: Microsoft Kinect," *IEEE Access*, vol. 1, pp. 465–474, 2013.