# Spatial Acoustic Projection for 3D Imaging Sonar Reconstruction
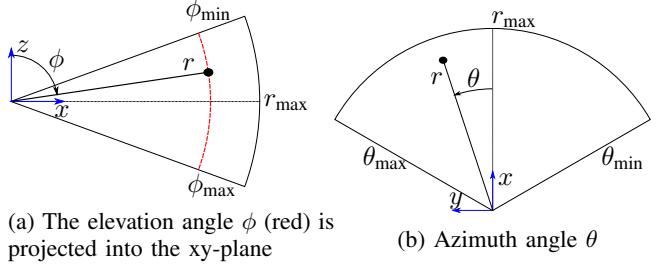
Sascha Arnold[1,2] and Bilal Wehbe[1]

***Abstract*— In this work we present a novel method for reconstructing 3D surfaces using a multi-beam imaging sonar. We integrate the intensities measured by the sonar from different viewpoints for fixed cell positions in a 3D grid. For each cell we integrate a feature vector that holds the mean intensity for a discretized range of viewpoints. Based on the feature vectors and independent sparse range measurements that act as ground truth information, we train convolutional neural networks that allow us to predict the signed distance and direction to the nearest surface for each cell. The predicted signed distances can be projected into a truncated signed distance field (TSDF) along the predicted directions. Utilizing the marching cubes algorithm, a polygon mesh can be rendered from the TSDF. Our method allows a dense 3D reconstruction from a limited set of viewpoints and was evaluated on three real-world datasets.**

(a) The elevation angle $\phi$ (red) is projected into the xy-plane

(b) Azimuth angle $\theta$

Fig. 1: Projection of a 3D point (black dot) into the 2D sonar image $I(r, \theta)$. The elevation angle $\phi$ is lost during acquisition as the 3D point could lie anywhere on the red-dashed arc.

## I. INTRODUCTION

Imaging sonars are a key sensor modality for under water vehicles, in particular in higher ranges when optical cameras are limited due to turbid water. The reconstruction of 3D information and mapping of the underwater environment is particularly interesting for autonomous underwater vehicles (AUVs) in order to fulfill inspection, exploration and mapping tasks [1].

Imaging sonars have a wide vertical opening angle (elevation) allowing to get returns from a subsection of the scene. While the horizontal angle (azimuth) and the range can be measured by the sonar, the elevation angle is lost during acquisition (Fig. 1a). The challenge of 3D imaging sonar reconstruction therefore is to recover the elevation angles to the surfaces in the scene.

Previous works on 3D reconstruction with imaging sonars can be roughly grouped into feature based methods, generative models and volumetric methods. Feature based methods utilize feature points in the sonar image, and match the corresponding features between images taken from different viewpoints. Based on the feature point correspondences, nonlinear optimization or filtering can be applied to create a sparse 3D reconstruction [12], [13], [18], [14]. The drawback of those methods is that the reconstruction is sparse, as feature points are often difficult to extract in sonar images and correspondences can be reliably found only at nearby viewpoints. Generative models utilize the measured intensities and a known starting position in order to derive the slopes of the corresponding surfaces in the scene [3], [6], [8], [19], [22], [24]. While allowing a locally dense 3D reconstruction, these methods rely on the estimate of object edges and knowledge of the reflective properties of the surfaces. Volumetric methods discretize the environment into a voxel grid and determine for each voxel their contribution to the sonar image. Space carving is one example of volumetric methods that utilizes only the free space information before the first high-intensity return in order to carve out the free space, while the remaining voxels are considered occupied [4], [5], [16], [9]. Other methods utilize occupancy grid mapping [21], [20] or more recently albedo-based methods [10], [23]. These methods rely on a high variety of viewpoints in order to achieve a dense 3D reconstruction.

This work presents a novel volumetric acoustic projection method that can achieve a dense 3D reconstruction with a limited set of viewpoints. As in other volumetric methods we discretize the environment into a voxel grid, with the difference that every cell stores a feature vector of intensities. Our approach aims to create dense 3D reconstructions with a limited set of viewpoints, for instance if the environment was traversed only once as in a typical survey mission (imaging sonar facing forward and slightly downward). We model the problem of reconstructing the positions and orientations of surfaces in the scene by defining features with fixed 3D positions that integrate intensities measured from different sensor viewpoints over time. Based on these features, we utilize neural networks in order to first classify outliers and second predict the signed distance and direction to the nearest surface. In order to train the neural networks we rely on the sparse range measurements from a Doppler Velocity Log (DVL) sensor, which act as ground truth information. The predicted signed distances are projected along their predicted direction into a truncated signed distance field (TSDF). TSDFs have recently become a common implicit surface representation for 3D reconstruction applications [7], [17]. Based on the TSDF, a polygon mesh can be rendered using the marching cubes algorithm [15]. The resulting mesh

[1]German Research Center for Artificial Intelligence, Bremen, Germany
[2]Kraken Robotics, Bremen, Germany
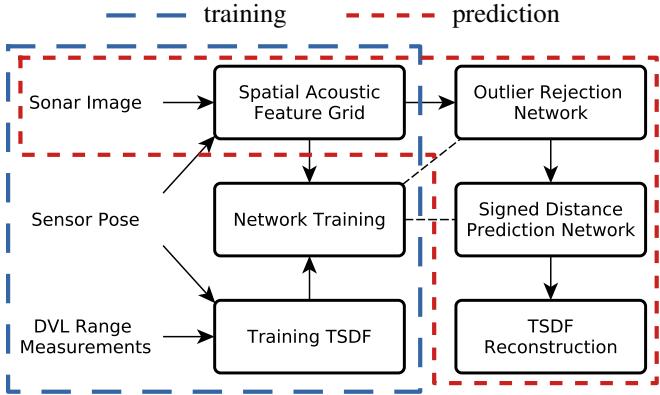sarnold@ieee.org, bilal.wehbe@dfki.de

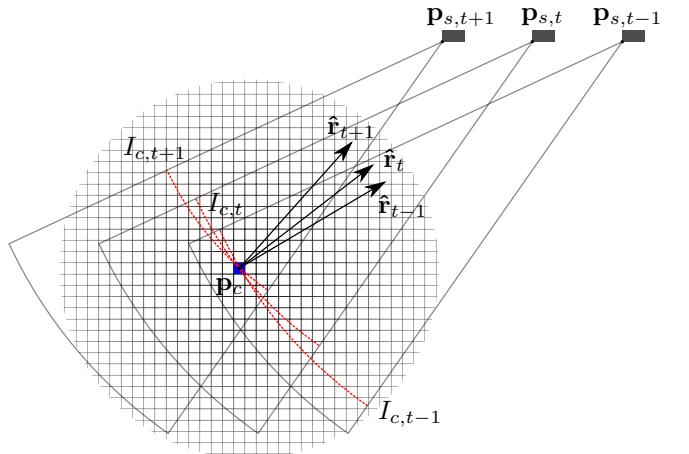Fig. 2: Spatial acoustic reconstruction architecture.



Fig. 3: Projection of measured intensities into a cell $c$ of the 3D grid from three different sensor viewpoints. The red-dashed arcs visualize the intensity measurement along the elevation angle for each viewpoint.

is used to evaluate our approach against the ground truth information.

Our method relies on the estimated poses provided by an inertial navigation system (INS) in order to relate the different sensor viewpoints. Additionally we utilize the four range measurements provided by a DVL sensor as independent measurements of the surfaces in the scene.

The main contributions of the paper are as follows:

1) Novel volumetric acoustic projection method that can predict the signed distance and direction to the nearest surface for each voxel.
2) Definition of a feature vector for acoustic intensities from different viewpoints.
3) Show that sparse DVL range measurements can be utilized to train models which allow a dense 3D reconstruction.

In the following section we describe the various aspects of our approach. Section III presents the results using three real-world datasets. Finally we give conclusion remarks in section IV.

## II. SPATIAL ACOUSTIC PROJECTION

This section details the various aspects of our 3D imaging sonar reconstruction approach. Firstly, a spatial acoustic grid is constructed using the measured intensities by the sonar from different viewpoints. Using the range measurement form a DVL, a TSDF is constructed which acts as a ground truth. Two neural networks are trained in order to relate the feature vectors from the spatial grid with a ground truth signed distance and direction to the nearest surface. The trained networks are then deployed to predict the signed distances and directions creating a TSDF that is used to reconstruct a 3D mesh. A block diagram of the overall architecture is illustrated in Fig. 2.

### A. Imaging sonar sensor

An imaging sonar is an active acoustic sensor which emits a pulse of sound and measures the intensity $I$ of the reflected pulse by the scene. The 2D image of a multi-beam sonar $I(r, \theta)$ is described by the range $r$, which is determined by the speed of sound in water and by the azimuth angle $\theta$, which is determined by a 1D array of transducers. The elevation angle $\phi$ is lost during the projection from the 3D world into the 2D image. Therefore, the measured intensity $I(r, \theta)$ includes all reflections along the elevation arc ($\phi_{\min}$, $\phi_{\max}$) defined by the vertical opening angle of the imaging sonar. Fig. 1 illustrates the projection of a point in Euclidean coordinates into the sonar image.

### B. Spatial Acoustic Feature Grid

We discretize the 3D Euclidean space into a 3D volumetric grid with a fixed cell resolution along each axis. Each grid cell stores the measured intensity $I$ for each angle between the vector $\hat{\mathbf{r}}_t$ and each axis vector. $\hat{\mathbf{r}}_t$ points from the grid cell to the sonar for a measurement at time $t$:

$$\hat{\mathbf{r}}_t = \frac{(\mathbf{p}_{s,t} - \mathbf{p}_c)}{\|(\mathbf{p}_{s,t} - \mathbf{p}_c)\|} \tag{1}$$

where $\mathbf{p}_{s,t}$ is the position of the sonar at time $t$ and $\mathbf{p}_c$ is the center of the cell $c$.

Given $\hat{\mathbf{r}}_t$, the angles $\gamma_i$ to each unit axis $\hat{\mathbf{e}}_i$ in the range $[0, \pi/2]$ for $i \in [x, y, z]$ are defined as:

$$\gamma_{i,t} = \arccos(|\hat{\mathbf{e}}_i \cdot \hat{\mathbf{r}}_t|) \tag{2}$$

with

$$\|\hat{\mathbf{r}}_t\| = \|\hat{\mathbf{e}}_i\| = 1 \ \forall i, t \tag{3}$$

We limit the angular range to $[0, \pi/2]$ since we assume the same intensity characteristic above and below the surface.

Fig. 3 illustrates the projection of three intensity measurements into the same cell from different sensor viewpoints. The intensity $I_{c,t}$ for a measurement at time $t$ at the position of a cell $c$ can be expressed in the sonar image $I_t(r, \theta)$ as:

$$I_{c,t} = I_t(\|\mathbf{p}_l\|, \arctan(y_l/x_l)) \tag{4}$$

where $\mathbf{p}_l$ is the cell center $\mathbf{p}_c$ expressed in the sonar frame:

$$\mathbf{p}_l = \begin{bmatrix} x_l \\ y_l \\ z_l \end{bmatrix} = \mathbf{C}_{c,t}^s \mathbf{p}_c \tag{5}$$
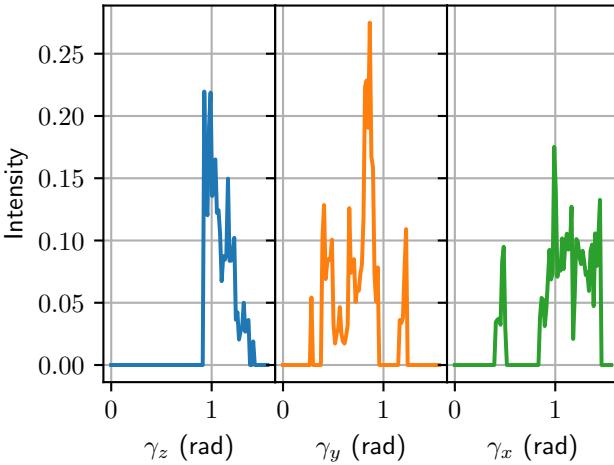
Fig. 4: Feature vector $\mathbf{X}$ split into its three components $[\mathbf{x}_z(\gamma_z), \mathbf{x}_y(\gamma_y), \mathbf{x}_x(\gamma_x)]^T$.



Fig. 5: Architecture of the outlier rejection network.

and $\mathbf{C}^s_{c,t}$ is the coordinate transformation from cell to sonar frame at time $t$.

Based on (2) and (4) we define the feature vectors $\mathbf{x}_i$ as:

$$\mathbf{x}_i(\gamma_{i,t}) = \frac{1}{(T-1)} \sum_{t=0}^{T} I_{c,t}, \; \forall i \in [x,y,z] \quad (6)$$

where the range of the angles $\gamma_i$ is discretized to a fixed resolution $n$. A feature vector $\mathbf{x}_i$ therefore is of size $n \times 1$. Each measured intensity $I_{c,t}$ at the cell position $\mathbf{p}_c$ is integrated in each feature vector $\mathbf{x}_i$. The three feature vectors $\mathbf{x}_i$ are arranged into the feature vector $\mathbf{X}$ as follows:

$$\mathbf{X} = [\mathbf{x}_z, \mathbf{x}_y, \mathbf{x}_x]^T \quad (7)$$

$\mathbf{X}$ therefore is of size $3n \times 1$ and will be the input of the neural networks. Each cell $c$ in the 3D grid has a feature vector $\mathbf{X}_c$. Fig. 4 shows an example of $\mathbf{X}$ split into its three components $\mathbf{x}_i$.

### C. Training TSDF

We build a TSDF based on the range measurements of a DVL sensor in order to determine a sparse ground truth of the true shape of the environment that can be utilized to train the neural networks. The range measurements are projected into the TSDF by tracing along the measurement from the direction of the sensor origin around the margin defined by a truncation threshold $\tau$, similar to [17]. By that each visited cell of the TSDF holds the truncated signed distance $d_c$ to the nearest surface. In order to account for the opening angle of each DVL beam we trace along 9 rays, one in the center and 8 equidistantly placed along the outside cone of the DVL beam defined by the opening angle.

To determine the gradients for each cell we apply a 3D Sobel filter with a kernel size of $3 \times 3 \times 3$. This gives us the distance gradients $\mathbf{n}_c$ for each cell with a defined neighborhood. Since we are interested in values close to zero and a compact representation as model output, we represent the direction as a delta on the unit sphere with respect to the
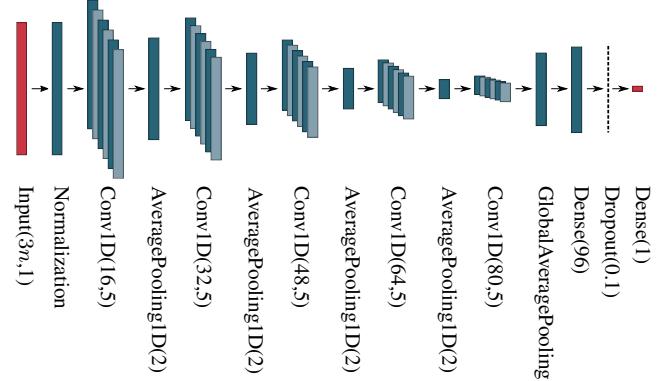
unit vector $\hat{\mathbf{e}}_z$. (8) computes the logarithm map that maps $\mathbf{n}_c$ to the tangent plane $\log_{\hat{\mathbf{e}}_z}(\mathbf{n}_c)$ determined by $\hat{\mathbf{e}}_z$ [11]:

$$\delta\mathbf{n}_c = \log_{\hat{\mathbf{e}}_z}\left(\frac{\mathbf{n}_c}{\|\mathbf{n}_c\|}\right) \in \mathbb{R}^2 \quad (8)$$

The gradient $\hat{\mathbf{n}}_c$ with unit length can be recovered using the exponential map:

$$\hat{\mathbf{n}}_c = \exp_{\hat{\mathbf{e}}_z}(\delta\mathbf{n}_c) \in \mathbb{S}^2 \quad (9)$$

With the tuple $[d_c, \delta\mathbf{n}_c]^T$ the signed distance and direction to the nearest surface can be described for each cell $c$.

### D. Neural network architectures

We are using two cascaded neural networks, the first one is a classifier and the second is a regressor. The outlier rejection network predicts for each feature vector if the corresponding cell is inside of the truncation threshold to the nearest surface. The signed distance prediction network predicts for positively classified feature vectors the signed distance to the nearest surface and the distance gradient.

The outlier rejection network is a classifier that predicts for each feature vector $\mathbf{X}_c$ if the corresponding cell $c$ is within the truncation threshold $\tau$ to the nearest surface. The boolean training samples $\mathbf{y}_b$ are defined based on the signed distance $d_c$:

$$\mathbf{y}_b = [|d_c| < \tau] \in [0,1] \quad (10)$$

The outlier classification is modeled as a convolutional neural network (CNN) as shown in Fig. 5 with 5 convolutional layers and 2 fully connected layers at the end. As activation function the rectified linear unit is used in all except the last layer, which uses the sigmoid activation function. We use a global average pooling layer instead of a flattening in order to reduce the number of trainable parameters. With a feature vector resolution of $n = 100$ the network consists of $59393$ trainable parameters.

The signed distance prediction network is a regressor that predicts for each feature vector $\mathbf{X}_c$ of the corresponding cell $c$ the signed distance to the nearest surface and the distance gradient:

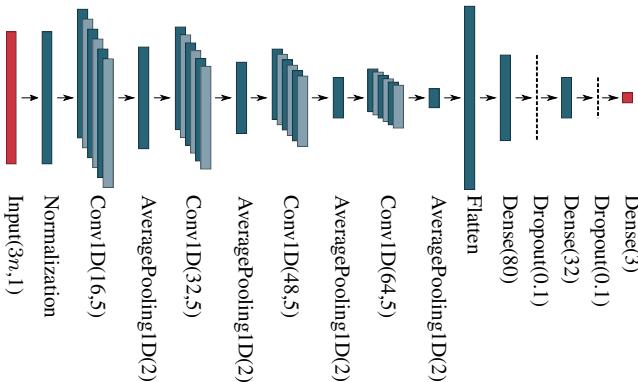$$\mathbf{y}_d = [d_c, \delta\mathbf{n}_c]^T \in \mathbb{R}^3 \quad (11)$$

Fig. 6: Architecture of the signed distance prediction network.



(a) Relative bathymetry height  (b) Measurements used for training (blue) and evaluation (orange)

Fig. 7: Accumulated DVL range measurements of the square trajectory dataset.

The signed distance prediction is modeled as a CNN as shown in Fig. 6. It has 4 convolutional layers, a flatten layer and 3 fully connected layers at the end. As activation function the rectified linear unit is used in all convolutional and dense layers. With a feature vector resolution of $n = 100$ the network consists of $105411$ trainable parameters.

The CNN architectures were selected in an progressive fashion, starting from a simple architecture and adding layers and/or cells as long as the performance improved.

*E. TSDF Reconstruction*

As seen in the overview in Fig. 2 the TSDF reconstruction is the last step of our approach that utilizes the outputs of the neural networks in order to reconstruct a polygon mesh based on the measurements of the imaging sonar.
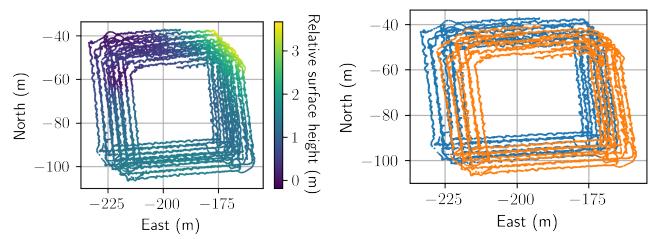
For each cell in the spatial acoustic feature grid the trained models can predict, using the feature vector $\mathbf{X}_c$ as input, if a cell is inside of the truncation threshold to the nearest surface $\tilde{\mathbf{y}}_b$ and if this is true the signed distance and distance gradient $\tilde{\mathbf{y}}_d$.

Based on the predictions $\tilde{\mathbf{y}}_d$ of the signed distance prediction network, the TSDF is built by tracing along the predicted gradient $\tilde{\hat{n}}_c$ with the predicted distance $\tilde{d}_c$ in the range of $[\tilde{d}_c - \tau, \tilde{d}_c + \tau]$ using the center of cell $c$ as origin. This process is repeated for all cells. The TSDF and the spatial acoustic feature grid share the same cell resolution.

Using the well known marching cubes algorithm [15] a polygon mesh can be extracted from the TSDF.

## III. RESULTS

For the evaluation of our method we utilize three datasets that have been collected in the open ocean using the *FlatFish* AUV [1]. The AUV is equipped with a Tritech Gemini 720i Multibeam Imaging Sonar and an INS described in [2]. The estimated pose of the AUV during the experiments is affiliated with an growing error throughout the mission [2]. Since in our experiments we are only integrating a limited time window of measurements we are disregarding the pose error and assume the pose to be known. From the three datasets one was split into a training and evaluation part, while the other two datasets were used only for evaluation.

In all experiments a truncation threshold of $\tau = 1.0$ and a feature vector resolution of $n = 100$ was selected.

*A. Network Training*

Fig. 7 shows the sparse DVL range measurements of the dataset that was used to train the networks defined in II-D. During the mission the vehicle followed repeated square trajectories with an edge length of 50 meter. As can be seen in Fig. 7a the trajectories have a displacement, which allows to cover a wider area of the bathymetry with the four DVL range measurements. The networks have been trained using the measurements of the first half while the measurements of the second half have been used for the evaluation of the resulting surface reconstruction (Fig. 7b).

The networks are trained on the associated data for each grid cell $c$ in the spatial acoustic feature grid (section II-B) and the training TSDF (section II-C). Both 3D grids therefore have the same resolution during the training step. The outlier rejection network is trained with the feature vector $\mathbf{X}_c$ as input and the boolean class $\mathbf{y}_{b,c}$ as output. The signed distance prediction network is trained with the same input and the vector $\mathbf{y}_{d,c}$ as output.

The training was performed using the Adam optimizer with a learning rate of $0.0001$, a validation split of $0.2$ and a grid cell resolution of $0.1m$. As loss the logarithm of the hyperbolic cosine was used for the signed distance prediction network and the binary cross-entropy for the outlier rejection network. The training was performed until the validation loss was stable for several epochs. The outlier rejection network reached an accuracy of $0.957$ after $85$ epochs and the signed distance prediction network a mean squared error (MSE) of $0.0267$ after $204$ epochs. The selected CNN architectures were outperforming networks utilizing only fully connected layers with a similar amount of trainable parameters (Accuracy: $0.952$, MSE: $0.0415$).

*B. Evaluation*

For the evaluation of the 3D reconstruction the second half of the square trajectory dataset, a dataset with a lower altitude to the seafloor and a dataset with man-made structures on the seafloor were used. All reconstructions are preformed based on the models trained with the first half of the square trajectory dataset.
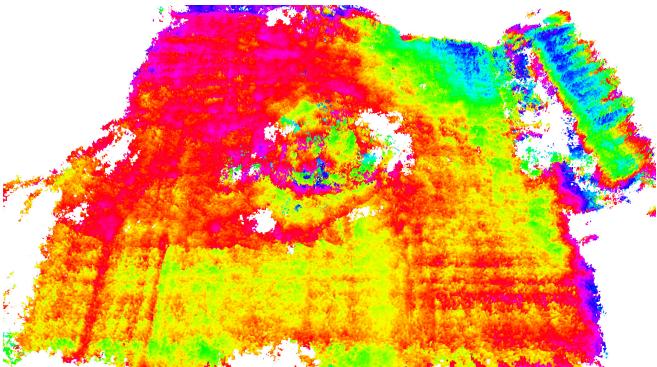
Fig. 8: Second half of the square trajectory dataset: Polygon mesh reconstruction of the seafloor with a HSV color mapping repeating every 5 meters along the z axis.
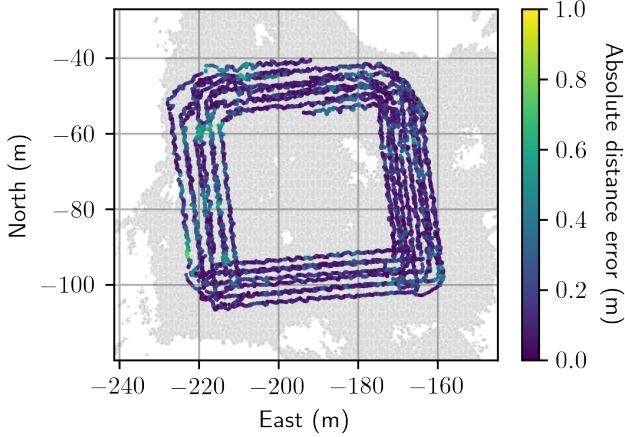


Fig. 9: Square trajectory dataset: Absolute distance error in meter between the DVL range measurements and the reconstructed mesh. The outline of the mesh is illustrated in grey.

The result of the 3D reconstruction described in section II-E is a polygon mesh of the surfaces in the scene based only on the imaging sonar intensity measurements from the viewpoints defined by the trajectory of the AUV.

*1) Square trajectory dataset:* The reconstructed mesh of the second half of the square trajectory dataset is visualized in Fig. 8. It shows the 3D reconstruction of the seafloor with a HSV color mapping repeating every 5 meters along the z axis. The mesh was reconstructed with a grid cell resolution of $0.2m$ and based only on the sonar measurements during the traverse of one square trajectory. In areas that have been covered from a variety of sonar viewpoints, the reconstruction shows to be consistent with the sparse DVL range measurements. For some cells, that have been covered by the sonar only close to a high azimuth angle and from very limited viewpoints, the models fail to classify the cells correctly, as can be seen in the center and the upper right part of the mesh.

In order to evaluate the accuracy of the reconstructed mesh the closest distance between the sparse DVL range measure-



Fig. 10: Low altitude dataset: Polygon mesh reconstruction of the seafloor with a HSV color mapping repeating every 5 meters along the z axis.
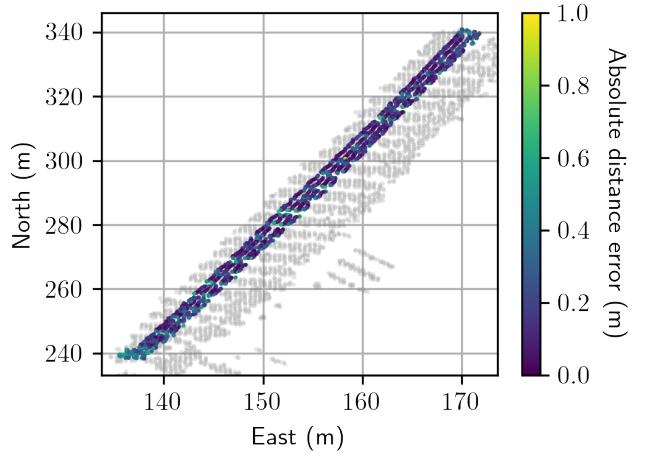


Fig. 11: Low altitude dataset: Absolute distance error in meter between the DVL range measurements and the reconstructed mesh. The outline of the mesh is illustrated in grey.

ments and the mesh was computed. Fig. 9 shows the absolute distance error between the DVL range measurements and the reconstructed mesh with a MSE of $0.033m^2$. The outline of the mesh is visible in the background.

*2) Low altitude dataset:* In a second dataset the vehicle followed a straight line for $100m$ in a lower altitude (2 to $3m$) to the seafloor. The polygon mesh was reconstructed based on a cell resolution of $0.25m$. Fig. 10 shows the mesh of the seafloor sensed by the imaging sonar with a HSV color mapping repeating every 5 meters along the z axis.

Fig. 11 shows the absolute distance error between the DVL range measurements and the reconstructed mesh with a MSE of $0.067m^2$. The mesh shows to be consistent with the ground truth, even though the models have been trained with higher altitude and on a different dataset. The mesh also shows regular gaps that are likely related to the pitch motion of the AUV.

*3) Man-made structures dataset:* In order to evaluate if the model can also be applied to datasets, containing (vertical) structures that the model has not seen before, the model was applied to a third dataset containing mockups of a pipeline and a subsea isolation valve (SSIV). The mockups are shown in Fig. 12.

In the reconstructed polygon mesh (Fig. 13) the locations of the pipeline and the SSIV can be identified. The model is able to reconstruct some of the vertical surfaces of the $2.8m$ high SSIV mockup. However the shadows casted by the SSIV in the sonar image create gaps in the dense

(a) SSIV         (b) Pipeline elements

Fig. 12: Subsea mockup structures. Images: Jan Albiez, SENAI CIMATEC
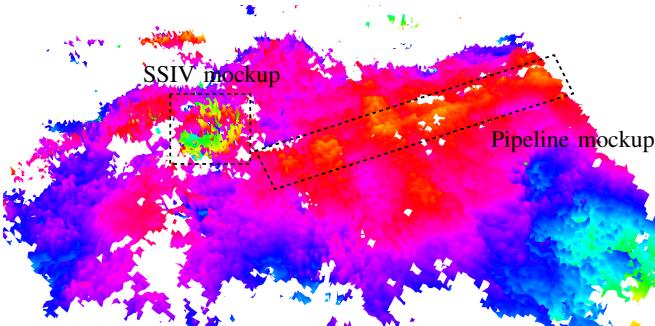


Fig. 13: Man-made structures dataset: Polygon mesh reconstruction of the seafloor sensed by the imaging sonar with a HSV color mapping repeating every 5 meters along the z axis.
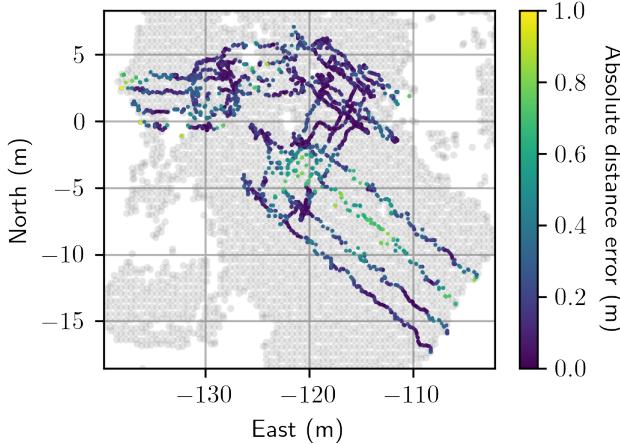


Fig. 14: Man-made structures dataset: Absolute distance error in meter between the DVL range measurements and the reconstructed mesh. The outline of the mesh is illustrated in grey.

reconstruction close to the SSIV.

Fig. 14 shows the absolute distance error between the DVL range measurements and the reconstructed mesh with a MSE of $0.12m^2$. In this dataset the mesh shows inconsistencies with the ground truth especially in the area of the pipeline.

TABLE I: Accuracy of the reconstructed meshes with respect to the cell resolution.

| Cell resolution | MSE to ground truth |
|---|---|
| $0.2m$ | $0.033m^2$ |
| $0.25m$ | $0.038m^2$ |
| $0.4m$ | $0.051m^2$ |
| $0.6m$ | $0.064m^2$ |

### C. Different cell resolutions

Since the models predict for each cell individually, they can be applied to different resolutions of the spatial acoustic feature grid while still being able to reconstruct the same consistent mesh. If the resolution is decreased of course the accuracy of the reconstruction suffers. Table I gives an overview on different cell resolutions and the MSE of the reconstructed meshes for the square trajectory dataset. While the memory footprint grows cubically with the selected cell resolution, the MSE shows to be closer to linear improvement.

## IV. CONCLUSION

In this work we have presented a novel approach for the reconstruction of 3D surfaces using an imaging sonar sensor. We defined a feature vector which can be utilized to train models that can predict the signed distance and direction to the nearest surface in the scene. The feature vector however could be replaced by another structure integrating the measured intensities from different viewpoints, for instance the discretized surface of a unit sphere. This would on the other hand also require a graph convolutional network architecture in order to model the neighborhood correctly. One limitation of our current approach is that the predicted surface direction is represented as a delta on the unit sphere with respect to a fixed reference axis (unit z-axis). While for datasets that are focused on the bathymetry this seems sufficient, the reference could be selected uniquely for each cell using additional constraints. Another limitation is that we disregard shadows and multi-path reflections present in the sonar image. These sections could be masked or modeled in order to improve the reconstruction result. It is part of future work to compare the results of our work with existing solutions, like space carving. We demonstrated the effectiveness of our approach on three real world datasets. While the sonar sensor in the datasets mainly imaged the seafloor, we could also show that man-made structures with vertical surfaces could be partially reconstructed by the models, even so they have not been present during the training step.

## ACKNOWLEDGMENT

## REFERENCES

[1] Jan Albiez, Sylvain Joyeux, Christopher Gaudig, Jens Hilljegerdes, Sven Kroffke, Christian Schoo, Sascha Arnold, Geovane Mimoso, Pedro Alcantara, Rafael Saback, et al. Flatfish-a compact subsea-resident inspection auv. In *OCEANS'15 MTS/IEEE Washington*, pages 1–8. IEEE, 2015.

[2] Sascha Arnold and Lashika Medagoda. Robust model-aided inertial localization for autonomous underwater vehicles. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*, pages 4889–4896. IEEE, 2018.

[3] Murat D Aykin and Shahriar Negahdaripour. Forward-look 2-d sonar image formation and 3-d reconstruction. In *2013 OCEANS-San Diego*, pages 1–10. IEEE, 2013.

[4] Murat D Aykin and Shahriar Negahdaripour. On 3-d target reconstruction from multiple 2-d forward-scan sonar views. In *OCEANS 2015-Genova*, pages 1–10. IEEE, 2015.

[5] Murat D Aykin and Shahriar Negahdaripour. Three-dimensional target reconstruction from multiple 2-d forward-scan sonar views by space carving. *IEEE Journal of Oceanic Engineering*, 42(3):574–589, 2016.

[6] Murat D Aykin and Shahriar S Negahdaripour. Modeling 2-d lens-based forward-scan sonar imagery for targets with diffuse reflectance. *IEEE Journal of Oceanic Engineering*, 41(3):569–582, 2016.

[7] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312, 1996.

[8] Robert DeBortoli, Fuxin Li, and Geoffrey A Hollinger. Elevatenet: A convolutional neural network for estimating the missing dimension in 2d underwater sonar images. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8040–8047. IEEE, 2019.

[9] Thomas Guerneve and Yvan Petillot. Underwater 3d reconstruction using blueview imaging sonar. In *OCEANS 2015-Genova*, pages 1–7. IEEE, 2015.

[10] Thomas Guerneve, Kartic Subr, and Yvan Petillot. Three-dimensional reconstruction of underwater objects using wide-aperture imaging sonar. *Journal of Field Robotics*, 35(6):890–905, 2018.

[11] Christoph Hertzberg, René Wagner, Udo Frese, and Lutz Schröder. Integrating generic sensor fusion algorithms with sound state representations through encapsulation of manifolds. *Information Fusion*, 14(1):57–77, 2013.

[12] Tiffany A Huang and Michael Kaess. Towards acoustic structure from motion for imaging sonar. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 758–765. IEEE, 2015.

[13] Tiffany A Huang and Michael Kaess. Incremental data association for acoustic structure from motion. In *2016 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1334–1341. IEEE, 2016.

[14] Jie Li, Michael Kaess, Ryan M Eustice, and Matthew Johnson-Roberson. Pose-graph slam using forward-looking sonar. *IEEE Robotics and Automation Letters*, 3(3):2330–2337, 2018.

[15] William E Lorensen and Harvey E Cline. Marching cubes: A high resolution 3d surface construction algorithm. *ACM siggraph computer graphics*, 21(4):163–169, 1987.

[16] Shahriar Negahdaripour, Victor M Milenkovic, Nikan Salarieh, and Mahsa Mirzargar. Refining 3-d object models constructed from multiple fs sonar images by space carving. In *OCEANS 2017-Anchorage*, pages 1–9. IEEE, 2017.

[17] Richard A Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J Davison, Pushmeet Kohi, Jamie Shotton, Steve Hodges, and Andrew Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *2011 10th IEEE international symposium on mixed and augmented reality*, pages 127–136. IEEE, 2011.

[18] Jinkun Wang, Tixiao Shan, and Brendan Englot. Underwater terrain reconstruction from forward-looking sonar imagery. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 3471–3477. IEEE, 2019.

[19] Yusheng Wang, Yonghoon Ji, Dingyu Liu, Hiroshi Tsuchiya, Atsushi Yamashita, and Hajime Asama. Elevation angle estimation in 2d acoustic images using pseudo front view. *IEEE Robotics and Automation Letters*, 6(2):1535–1542, 2021.

[20] Yusheng Wang, Yonghoon Ji, Hanwool Woo, Yusuke Tamura, Atsushi Yamashita, and Hajime Asama. Three-dimensional underwater environment reconstruction with graph optimization using acoustic camera. In *2019 IEEE/SICE International Symposium on System Integration (SII)*, pages 28–33. IEEE, 2019.

[21] Yusheng Wang, Yonghoon Ji, Hanwool Woo, Yusuke Tamura, Atsushi Yamashita, and Asama Hajime. 3d occupancy mapping framework based on acoustic camera in underwater environment. *IFAC-PapersOnLine*, 51(22):324–330, 2018.

[22] Eric Westman, Ioannis Gkioulekas, and Michael Kaess. A theory of fermat paths for 3d imaging sonar reconstruction. In *2020 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 5082–5088. IEEE, 2020.

[23] Eric Westman, Ioannis Gkioulekas, and Michael Kaess. A volumetric albedo framework for 3d imaging sonar reconstruction. In *2020 IEEE International Conference on Robotics and Automation (ICRA)*, pages 9645–9651. IEEE, 2020.

[24] Eric Westman and Michael Kaess. Wide aperture imaging sonar reconstruction using generative models. In *2019 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 8067–8074. IEEE, 2019.