# AiSDF: Structure-aware Neural Signed Distance Fields in Indoor Scenes

Jaehoon Jang[1*], Inha Lee[1*], Minje Kim[1] and Kyungdon Joo[2†]

*Abstract*—Indoor scenes we are living in are visually homogenous or textureless, while they inherently have structural forms and provide enough structural priors for 3D scene reconstruction. Motivated by this fact, we propose a structure-aware online signed distance fields (SDF) reconstruction framework in indoor scenes, especially under the Atlanta world (AW) assumption. Thus, we dub this incremental SDF reconstruction for AW as AiSDF. Within the online framework, we infer the underlying Atlanta structure of a given scene and then estimate planar surfel regions supporting the Atlanta structure. This Atlanta-aware surfel representation provides an explicit planar map for a given scene. In addition, based on these Atlanta planar surfel regions, we adaptively sample and constrain the structural regularity in the SDF reconstruction, which enables us to improve the reconstruction quality by maintaining a high-level structure while enhancing the details of a given scene. We evaluate the proposed AiSDF on the ScanNet and ReplicaCAD datasets, where we demonstrate that the proposed framework is capable of reconstructing fine details of objects implicitly, as well as structures explicitly in room-scale scenes.

*Index Terms*—Deep learning for visual perception, mapping, incremental learning

## I. INTRODUCTION

VARIOUS 3D scene representations, such as explicit geometric primitive and implicit functions, have been actively studied in computer vision and robotics [2], [3], [4]. As one of the implicit representations, signed distance fields (SDF) inherently encode the surface information as the signed distance between the position in space and the closest surface, where the zero-level set corresponds to the surface. By virtue of this characteristic, many vision tasks, such as rendering [5], and path planning [6], use SDF as a medium, especially neural implicit reconstruction based on SDF has gained a lot of

[1]J. Jang, I. Lee and M. Kim are with the Artificial Intelligence Graduate School, UNIST, Ulsan, South Korea. (e-mail: erick1997@unist.ac.kr; epsilon8854@unist.ac.kr; minje617@unist.ac.kr)

[2]Kyungdon Joo is with the Artificial Intelligence Graduate School and the Department of Computer Science and Engineering, UNIST, Ulsan, South Korea. (e-mail: kyungdon@unist.ac.kr)

*Equal contribution (alphabetical order by last name)

†Corresponding author

Project website: https://vision3d-lab.github.io/AiSDF/

Digital Object Identifier (DOI): see top of this page.

attention [7], [8]. In those tasks, inferring accurate SDF with low latency is important.

Recently, Ortiz *et al.* [9] presented an incremental SDF estimation framework (iSDF in short) that reconstructs the SDF of room-scale indoor environments using continual learning in real-time. Given a stream of posed depth images, iSDF focuses on a neural SDF-based mapping module within a SLAM framework. By employing a compact MLP network and sparse sampling, they show that online SDF reconstruction is feasible using continual learning. However, iSDF misses several properties in *indoor scenes* that can improve reconstruction quality while maintaining efficiency.

From the layout or floorplan of the rooms to various objects such as furniture, most objects, including the scene itself, have structural forms in indoor scenes. Concretely, while they are homogeneous or textureless from a visual perspective, they consist of a set of orthogonal or parallel planes (planar segments) from a geometric viewpoint. These structural characteristics of indoor scenes can be represented by a few dominant directions; structural assumptions, such as the Manhattan world (MW) [10], or Atlanta world (AW) [11] assumptions, have been explored in the literature [12], [13], [14], [15]. For example, the MW assumption, represented by three orthogonal directions, describes a given scene with a strictly aligned shape, like a cuboid shape. In the case of the AW assumption, it can cover more general indoor scenes, such as non-orthogonal walls, using vertical and a set of horizontal directions. These structural assumptions have been exploited as prior information on indoor scenes.

Motivated by this fact, we propose a structure-aware online SDF reconstruction framework, AiSDF, in indoor scenes under the AW assumption (see Fig. 1). To this end, we continually estimate the underlying Atlanta structure of a given scene inside the online SDF reconstruction framework. This structural understanding provides several advantages within our AiSDF framework. 1) Based on the estimated dominant Atlanta directions as a priori, we can efficiently extract planar regions following the AW assumption in the form of surface elements (surfels). This Atlanta-aware surfel representation provides an explicit planar map for a given scene. 2) We can exploit the structural regularity as a constraint. Concretely, we can adaptively sample points according to the Atlanta-aware surfels, which enables us to enforce the additional structural constraint and focus on complex regions. We seamlessly integrate the structural understanding into the online SDF reconstruction framework. We demonstrate our AiSDF framework on the ScanNet [1] and ReplicaCAD [16] datasets. AiSDF shows that the overall details of the scene and
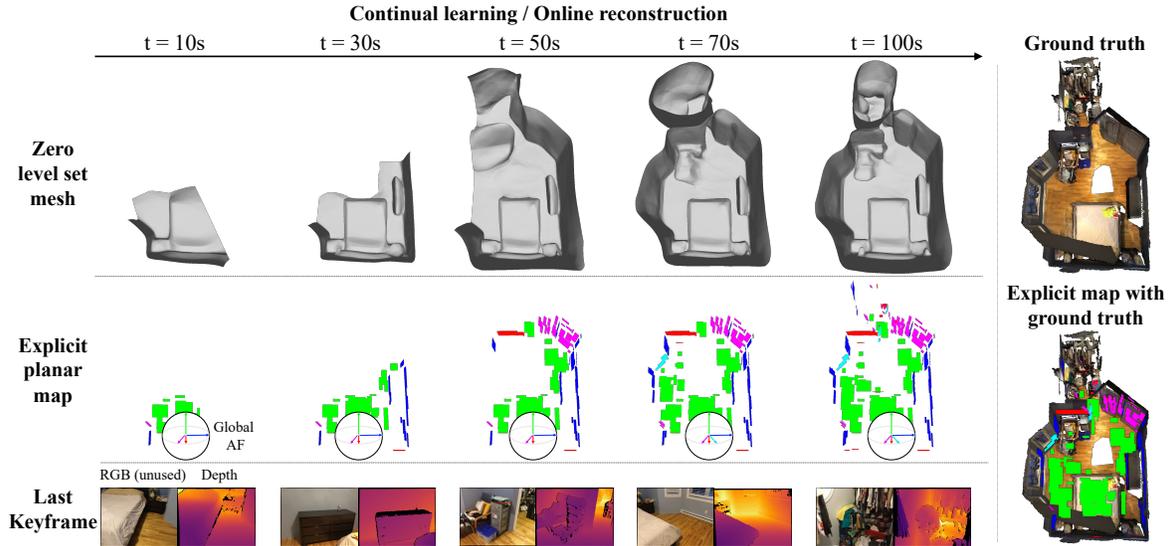
Fig. 1: **The proposed `AiSDF` on the ScanNet dataset [1]**. *Top*: Our framework represents the scene as a signed distance fields (SDF) by considering the structure of the scene in a continual manner. *Middle*: In addition, we estimate the underlying Atlanta structure (global Atlanta frame) and extract a 3D explicit planar map in the form of Atlanta-aware surfels. We colorize each surfel with the associated Atlanta direction. *Bottom*: We visualize the last keyframe used by `AiSDF` (RGB image is unused in practice).

the structure in the form of a 3D planar map are recovered better than comparison methods [17], [9]. In summary, our contributions are:

- We propose a new structure-aware online neural SDF, `AiSDF` that reconstructs a given indoor scene under the AW assumption with an online process.
- Based on the structural understanding, we introduce an Atlanta-aware surfel representation, which approximates a given indoor scene to a set of rectangular surfels.
- By utilizing the Atlanta-aware surfels, we effectively sample points considering the structure of the scene. In addition, we perform a structure-aware tight bound computation for self-supervised learning.
- In addition to obtaining a neural implicit map, `AiSDF` extracts a low-memory explicit planar map that can make it easier for robots to access the structure information of the scene.

## II. RELATED WORK

**Structural assumptions**. Thanks to their simplicity, represented by a few dominant directions, and their applicability in structured environments, various structural assumptions have been studied in robotics and computer vision [18], [19], [15] (please refer to [18] for a detailed review). The Manhattan world (MW) assumption [10], represented by three orthogonal directions, can approximate cuboid shape scenes, such as an indoor room. Beyond the MW assumption, the Atlanta world (AW) assumption [11] is defined by a vertical direction and a set of horizontal directions orthogonal to the vertical one, which can describe most indoor scenes, including non-orthogonal walls.

These two structural assumptions have been broadly used in various vision applications, such as scene understanding [20], [21], camera calibration [22], visual odometry [14], SLAM [13], [15], and so on. In particular, Joo *et al*. [15] propose a linear SLAM framework for

structured environments. They estimate the underlying Atlanta structure and explicitly use planar features supporting the Atlanta structure as measurements. In this work, we assume a given indoor scene follows the AW assumption.

**Neural scene reconstruction**. Recently, research on neural scene reconstruction using implicit representations, such as occupancy, neural radiance fields (NeRF), and SDF has been actively conducted [23], [24], [25], [26]. Among various implicit representations, SDF has gained much attention in that it can implicitly encode surface information using continuous values [27], [28], [29], [30]. In addition, SDF can make significant synergy with volumetric rendering and produce a high-quality reconstruction [31], [32], [33].

Neural scene reconstruction also can be utilized in conjunction with online processes, such as SLAM [7], [8]. Within the traditional real-time SLAM pipeline (*i.e.*, front-end tracking and back-end mapping), iMAP [7] takes an RGB-D image as input and exploits an MLP to implicitly represent a 3D volumetric map in the form of volume density. Recently, Ortiz *et al*. [9] propose iSDF, a continual SDF reconstruction framework given a stream of posed depth images. Inspired by iMAP, iSDF adopts a keyframe selection module to achieve real-time performance as well as train the model in a continual manner. Although this recent progress in online implicit 3D scene reconstruction is impressive and impactful, they less attention to improving the reconstruction quality. In other words, they focus on the efficiency of 3D neural scene representation.

Several research works [34], [35], [36] pay attention to enhancing the quality of neural scene reconstructions by combining geometric priors (*e.g*., depth and normal).[1] Guo *et al*. [36] propose a neural 3D scene reconstruction method with a strict structural assumption of indoor scenes. This ManhattanSDF method seamlessly combines

---

[1]Note that this line of research works focuses on enhancing reconstruction quality by offline process regardless of computational efficiency.
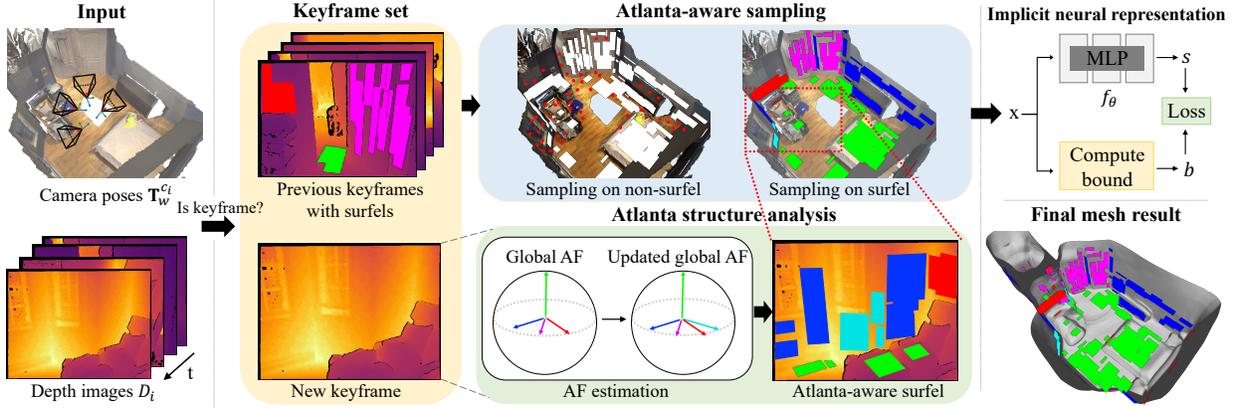
Fig. 2: **Overview of `AiSDF`**. Given a stream of posed depth images, AiSDF first selects the keyframe and adds it to the keyframe set for continual learning. We update the global Atlanta frame (AF) by extracting the dominant directions from a new keyframe and then generate surfels that represent the planar regions supported by the updated global AF. From a set of keyframes with Atlanta-aware surfels, we sample the 3D points considering the structure of the scene. Finally, sampled point x is queried to MLP that outputs signed distance value s, and we optimize the network in a self-supervised manner by measuring the loss between s and bound b. Note that we intentionally present intermediate steps of continual learning to show the process of extracting the new Atlanta direction and surfels supported by updated global AF. In Atlanta-aware sampling (blue box), we use the ground truth mesh to visualize the sampling effectively. The final mesh result indicates the reconstructed mesh by AiSDF using all keyframes.

the planar regions following the MW assumption with the learning of implicit neural representations and improves the reconstruction quality, especially in textureless and homogeneous regions. While the above methods show improved reconstruction quality, they cannot be utilized in the online process. Therefore, we seamlessly integrate an online SDF reconstruction framework with a general structural assumption, the AW assumption.

**Sampling strategy**. Unlike offline scene reconstruction approaches [33], [27], [36] that fully utilize input data regardless of computational time and memory consumption, sampling strategy becomes essential for online scene reconstruction to ensure efficiency. In the 2D image domain, iMAP [7] proposes loss-guided active sampling, which samples more points in the complex (high-frequency) regions based on loss calculated from the image grid. In the depth domain, iSDF [9] randomly samples a small number of pixel coordinates ($\leq 200$) from each selected keyframe for efficiency. Such sparse sampling strategies for scene reconstructions require an appropriate sampling rate based on the structural complexity of the scene. However, the accurate decision on the complexity of scenes in 2D or 2.5D domains is still a difficult challenge. To alleviate this issue, we exploit planar surfel regions supported by the AW assumption, which allows us to use Atlanta structure-aware sampling.

## III. ATLANTA-AWARE iSDF FRAMEWORK

In this work, we propose a new structure-aware online SDF reconstruction framework in indoor scenes (see Fig. 2). Unlike the previous iSDF [9] that focuses on reconstructing SDF itself, we exploit the structural regularity of indoor environments, especially the AW assumption [11]. We dub this structure-aware online SDF estimation `AiSDF` in short.

`AiSDF` takes as input a stream of posed depth images $\{D_i\}$ and camera poses $\{\mathbf{T}_w^{c_i}\}$, and aims to learn a neural network $f(\mathbf{x})$ based on the structural understanding of the AW, where $f(\mathbf{x})$ estimates SDF value s at a 3D point $\mathbf{x} \in \mathbb{R}^3$. Specifically,

for a consecutively selected keyframe $\mathcal{K}_j = (D_j, \mathbf{T}_w^{c_j})$, `AiSDF` infers the underlying Atlanta structure (*i.e.*, Atlanta frame) within a continual framework. Based on this structural understanding, we then estimate planar regions that support the estimated AF in the form of surfel representation. According to surfel regions, we perform Atlanta-aware sampling to force the structural regularity and focus on complex regions in SDF reconstruction adaptively. Thus, the proposed `AiSDF` framework improves the SDF reconstruction quality at the structure level as well as generates an explicit 3D planar map composed of Atlanta-aware surfels.

### A. Structural assumption: Atlanta frame

The AW assumption [11] can approximate a given indoor scene into a set of orthogonal and parallel planes, where planar walls are orthogonal to floors but do not have to be orthogonal to each other. We can formally define a set of dominant directions satisfying the AW assumption, which consists of a vertical dominant direction $\mathbf{v}_v$ and a set of $M$ horizontal dominant directions $\mathbf{v}_{h_m}$, where $\mathbf{v}_v \perp \mathbf{v}_{h_m}$. We call this direction set $\mathcal{V} = \{\mathbf{v}_v, \mathbf{v}_{h_1}, \mathbf{v}_{h_2}, \cdots, \mathbf{v}_{h_M}\}$ the *Atlanta frame (AF)* or *Atlanta directions*. In this work, we assume that a given indoor scene follows the AW assumption and use the AF parametrization [19][2] that represents Atlanta directions using the rotation matrix $\mathbf{R}$ and a set of 1D angles $\{\alpha_m\}$. We denote the estimated AF for the given $j$-th keyframe at the camera coordinate as local AF $\mathcal{V}_L^j$ and the unified AF for the observed keyframes at the world coordinate as global AF $\mathcal{V}_G$.

### B. Estimation of underlying Atlanta frame

For a consecutively selected keyframe, we estimate the underlying Atlanta structure (*i.e.*, global AF) in a continual

[2] AF parametrization [19] uses the rotation matrix $\mathbf{R} = [\mathbf{r}_1, \mathbf{r}_2, \mathbf{r}_3] \in SO(3)$ to represent the vertical direction and the first horizontal direction (*i.e.*, $\mathbf{v}_v = \mathbf{r}_1$ and $\mathbf{v}_{h_1} = \mathbf{r}_2$, where $\mathbf{v}_{h_1}$ acts as a reference location). Then, each $\mathbf{v}_{h_m}$ can be defined as a 1D angle parameter $\alpha_m$ by rotating $\mathbf{v}_{h_1}$ by $\alpha_m$ around $\mathbf{v}_v$.
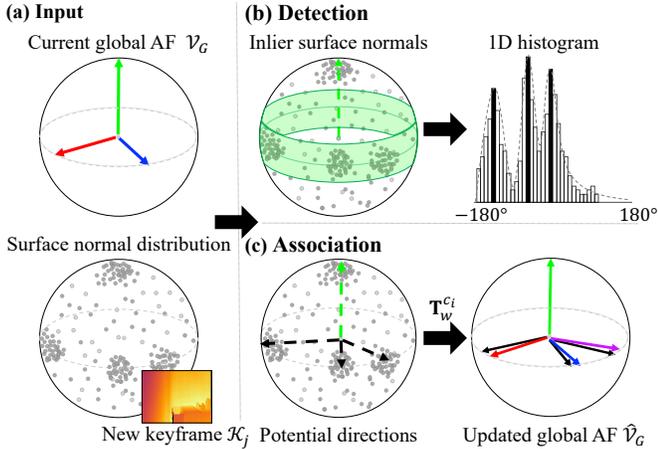
Fig. 3: **Illustration of underlying AF estimation**. (a) Given global AF $\mathcal{V}_G$ (solid arrows) and surface normal distribution of new keyframe $\mathcal{K}_j$, the Atlanta structure analysis proceeds in two steps. (b) First, we estimate potential dominant horizontal directions (black bars) from a 1D histogram of inlier surface normals. (c) The new dominant direction (purple arrows) is extracted by associating potential directions with the global AF in the world coordinate.

manner, as shown in Fig. 3. Unlike the previous work [15] that detects the AF at each frame and then naively associates them, we fully exploit the property of the AW assumption that dominant horizontal directions are orthogonal to the vertical direction. That is, the potential horizontal direction does exist on the horizon defined by the vertical direction. Based on this property, we reduce search space for the potential horizontal directions to 1D space in the range of $[-180°, 180°]$. In this 1D domain, we detect dominant horizontal directions and then associate them with the current global AF.

Specifically, given the current global AF $\mathcal{V}_G$ and a new posed keyframe $\mathcal{K}_j$ with its surface normal distribution[3], we first transform the vertical direction $\mathbf{v}_v \in \mathcal{V}_G$ into camera coordinate by $\mathbf{T}_w^{c_j -1}$. We then set an inlier horizon region *w.r.t.* the vertical direction and identify the inlier surface normals within the inlier horizon region (see the green band region in Fig. 3(b)). These inlier surface normals describe the distribution of potential dominant horizontal directions of the new keyframe. Thus, we construct a 1D histogram with 361 buckets, spaced at $1°$ intervals in the range $[-180°, 180°]$ for the inlier surface normals, and detect the dominant directions from peak points of smoothed histogram via Gaussian filtering. By doing so, we can efficiently estimate potential dominant horizontal directions under a known vertical direction. Finally, we transform the potential directions to the world coordinate by $\mathbf{T}_w^{c_j}$ and then associate them with the current global AF to determine the local AF $\mathcal{V}_L^j$ and update the global AF $\hat{\mathcal{V}}_G$ (see Fig. 3(c)).

We initialize $\mathcal{V}_G$ as the dominant Manhattan directions (*e.g.*, [37]) and set the angle thresholds for inlier horizon region and association as $20°$ to detect the dominant one and be robust against noise.

---

[3]Given a depth image and intrinsic parameter, we can directly compute the corresponding surface normal map using x-axis and y-axis gradient values and their cross-product.
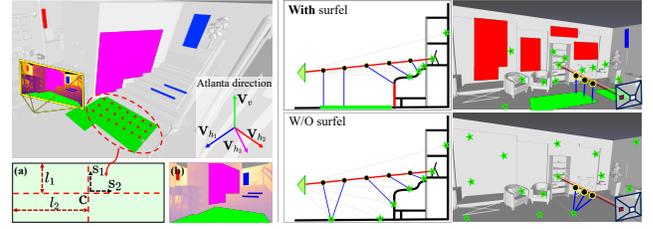


Fig. 4: **Illustration of the Atlanta-aware surfels and surfel-aware bound computation**. *Left*: (a) Atlanta-aware surfel representation. (b) 2D surfel mask $\mathbf{M}_{\mathfrak{s}}$ overlaid on the depth image. *Right*: Selected points to compute bound values with surfels (top) and without surfels (bottom). With Atlanta-aware surfels, AiSDF can compute more tight bound values while covering a denser and wider area.

### C. Atlanta-aware surfel representation

Given a keyframe $\mathcal{K}_j$ with the estimated local AF $\mathcal{V}_L^j$, we extract dominant planar regions supporting the AF in the form of surface elements, *Atlanta-aware rectangular surfels*, as shown in Fig. 4(a). We represent an Atlanta-aware rectangular surfel $\mathfrak{s}$ with its center point $\mathbf{c}$, surface normal $\mathbf{n}$, two axes $\mathbf{s}_1$ and $\mathbf{s}_2$, and the corresponding lengths $l_1$ and $l_2$. Through Atlanta directions, we can constrain the surface normal of surfels and define axes of rectangular surfels using other Atlanta directions orthogonal to the surface normal. This process allows us to explicitly extract a set of rectangular surfels without a learning-based module (*e.g.*, wall segmentation module [36]) and force the structural regularity inside the online SDF framework.

We first detect an Atlanta-aware plane using plane RANSAC as in [15]. Let $\mathbf{v}$ be one of the given Atlanta directions and $\mathcal{X}_{\mathbf{v}}$ be a set of 3D points having similar normals with $\mathbf{v}$. With randomly sampled 1-point in $\mathcal{X}_{\mathbf{v}}$ and $\mathbf{v}$, we can compute a potential Atlanta-aware plane and perform plane RANSAC using the Euclidean distance between the potential plane and $\mathcal{X}_{\mathbf{v}}$. We denote the detected Atlanta-aware planes by RANSAC as $\pi_{\mathbf{v}}$ and its inlier points as $\hat{\mathcal{X}}_{\mathbf{v}} \subset \mathcal{X}_{\mathbf{v}}$. We then extract Atlanta-aware rectangular surfel $\mathfrak{s}$ on the detected plane, where its surface normal and two axes are directly computed from $\mathbf{v}$ and two orthogonal Atlanta directions (*e.g.*, $\mathbf{n} = \mathbf{v}_v$, $\mathbf{s}_1 = \mathbf{v}_{h_1}$, and $\mathbf{s}_2 = \mathbf{v}_v \times \mathbf{v}_{h_1}$). Specifically, we project $\hat{\mathcal{X}}_{\mathbf{v}}$ on $\pi_{\mathbf{v}}$, of which 2D space is defined by two axes $\mathbf{s}_1$ and $\mathbf{s}_2$. On this 2D projection domain, we randomly sample 2-point and construct an axis-aligned rectangular surfel. We generate a set of candidate rectangular surfels by random sampling and select dominant surfels in terms of their area and occupancy by points, where the occupancy is measured by a ratio of associated points on the surfel against $\hat{\mathcal{X}}_{\mathbf{v}}$. For each direction $\mathbf{v} \in \mathcal{V}_L^j$, we extract Atlanta-aware surfels $\{\mathfrak{s}\}$ in 3D space and 2D surfel mask $\mathbf{M}_{\mathfrak{s}}$ in the depth image domain.

### D. Atlanta-aware sampling

Given a set of keyframes $\{\mathcal{K}_j\}$ with the estimated structure-aware surfels $\{\mathfrak{s}\}$, we perform structure-aware sampling that effectively considers planar surfel regions (*i.e.*, 2D surfel mask on depth image domain). According to the surfel mask $\mathbf{M}_{\mathfrak{s}}$, we separately sample points on the surfel and non-surfel regions ($N_{\mathfrak{s}}$ number of surfel points $\mathcal{X}_{\mathfrak{s}} = \{\mathbf{x}_{\mathfrak{s}}\}$ and $N_{\mathbf{n}}$ number of non-surfel points $\mathcal{X}_{\mathbf{n}} = \{\mathbf{x}_{\mathbf{n}}\}$). By performing this structure-aware

sampling, we intend to sample more non-surfel points (*i.e.*, relatively dense sampling) on the complex non-surfel regions while fewer surfel points on homogeneous surfel regions, as illustrated in the blue-colored box of Fig. 2.

Concretely, for the non-surfel regions, we randomly select a set of pixels and then sample a bunch of 3D points along the corresponding ray for each sample pixel as in [9]. For the surfel regions, we uniformly sample 3D points on the 3D surfel region directly (see the left side of Fig. 4), where the number of sampled points for each surfel is the same as the number of sampled points for a pixel along the ray for the non-surfel region.

We then compute the bound and approximated gradient depending on where the points are sampled. In the case of bound computation, we use zero bound value for $\mathcal{X}_\mathfrak{s}$ and the closest distance to surface points and surfel points for $\mathcal{X}_\mathfrak{n}$, where we consider the sign by computing the difference of depth between the depth from a sensor and the point sample along the ray:

$$b(\mathbf{x}, \mathcal{P}) = \begin{cases} \text{sgn}(D[u,v] - d)\min_{\mathbf{p} \in \mathcal{P}} |\mathbf{x} - \mathbf{p}| & \mathbf{x} \in \mathcal{X}_\mathfrak{n} \\ 0 & \mathbf{x} \in \mathcal{X}_\mathfrak{s} \end{cases} \quad (1)$$

where $\text{sgn}(\cdot)$ is the sign function, $\mathcal{P}$ is a set of surface points and surfel points, and $D[u,v]$ and $d$ are depth values of pixel $[u,v]$ and query point, respectively. The approximated gradient is computed with the closest surface point for $\mathbf{x} \in \mathcal{X}_\mathfrak{n}$, while it is replaced with the Atlanta direction for $\mathcal{X}_\mathfrak{s}$:

$$g(\mathbf{x}, \mathcal{P}) = \begin{cases} \text{sgn}(D[u,v] - d)(\mathbf{x} - \arg\min_{\mathbf{p} \in \mathcal{P}} |\mathbf{x} - \mathbf{p}|) & \mathbf{x} \in \mathcal{X}_\mathfrak{n} \\ \mathbf{v} & \mathbf{x} \in \mathcal{X}_\mathfrak{s} \end{cases} \quad (2)$$

where $\mathbf{v}$ is an Atlanta direction. By directly sampling points from surfels, we can cover a denser and wider area without increasing the total number of samples. This leads to a more accurate approximation of the SDF and gradient, as shown on the right side of Fig. 4. As a result, using these approximated values according to the type of sampled points, we can enforce `AiSDF` to learn more tight SDF reconstruction in a self-supervised manner.

### E. Loss for training `AiSDF`

Similar to [9], the loss function for training `AiSDF` basically consists of three terms: SDF, gradient, and Eikonal loss terms to satisfy the geometric properties of SDF. We adaptively utilize each loss term according to the identity of sampled 3D points (*i.e.*, surfel points $\mathcal{X}_\mathfrak{s}$ or non-surfel points $\mathcal{X}_\mathfrak{n}$). In the case of non-surfel points as query points, we optimize the network using three loss functions used in iSDF (see details in [9]) with approximated bound and gradient computed by Eqs. (1) and (2). In this work, we introduce a surfel-guided loss function for surfel points, which allows us to enforce structural regularity.

**Surfel loss**. Given sampled surfel points $\mathbf{x}_\mathfrak{s} \in \mathcal{X}_\mathfrak{s}$, we strongly constrain planar regions. To represent its geometry, the predicted SDF value is forced to have zero value:

$$\mathcal{L}_{\text{sdf}}^\mathfrak{s}(f(\mathbf{x}_\mathfrak{s}; \theta)) = |(f(\mathbf{x}_\mathfrak{s}; \theta)|. \quad (3)$$

Also, since the normal of $\mathbf{x}_\mathfrak{s}$ is aligned with Atlanta direction $\mathbf{v}$, we can apply the gradient loss to align the two vectors:

$$\mathcal{L}_{\text{grad}}^\mathfrak{s}(\nabla_{\mathbf{x}_\mathfrak{s}} f(\mathbf{x}_\mathfrak{s}; \theta), \mathbf{v}) = 1 - \frac{\nabla_{\mathbf{x}_\mathfrak{s}} f(\mathbf{x}_\mathfrak{s}; \theta) \cdot \mathbf{v}}{\|\nabla_{\mathbf{x}_\mathfrak{s}} f(\mathbf{x}_\mathfrak{s}; \theta)\| \|\mathbf{v}\|}, \quad (4)$$

where $\nabla_{\mathbf{x}_\mathfrak{s}} f(\mathbf{x}_\mathfrak{s}; \theta)$ denotes the gradient of the predicted SDF. In particular, we can constrain the normal in a common direction regardless of keyframes due to using global AF, which is one of our contributions using structural regularity.

Following [38], the Eikonal loss enforces the SDF to have a unit norm gradient to produce a high-fidelity surface:

$$\mathcal{L}_{\text{eik}}^\mathfrak{s}(f(\mathbf{x}_\mathfrak{s}; \theta)) = |\|\nabla_{\mathbf{x}_\mathfrak{s}} f(\mathbf{x}_\mathfrak{s}; \theta)\| - 1|. \quad (5)$$

Our entire surfel loss is as follows:

$$\mathcal{L}_{\text{surfel}} = \lambda_{\text{sdf}}^\mathfrak{s} \mathcal{L}_{\text{sdf}}^\mathfrak{s} + \lambda_{\text{grad}}^\mathfrak{s} \mathcal{L}_{\text{grad}}^\mathfrak{s} + \lambda_{\text{eik}}^\mathfrak{s} \mathcal{L}_{\text{eik}}^\mathfrak{s}, \quad (6)$$

where $\lambda_{\text{sdf}}^\mathfrak{s}$, $\lambda_{\text{grad}}^\mathfrak{s}$ and $\lambda_{\text{eik}}^\mathfrak{s}$ are the weight factors of the SDF loss, gradient loss and eikonal loss for surfel points. We set $\{\lambda_{\text{sdf}}^\mathfrak{s}, \lambda_{\text{grad}}^\mathfrak{s}, \lambda_{\text{eik}}^\mathfrak{s}\} = \{1, 0.4, 0.2\}$.

**Total loss**. Finally, the total loss of `AiSDF` is as follows:

$$l(\theta) = \mathcal{L}_{\text{iSDF}} + \mathcal{L}_{\text{surfel}}, \quad (7)$$

where $\mathcal{L}_{\text{iSDF}}$ indicates the loss term composed of SDF, gradient, and Eikonal loss for non-surfel sampled 3D points. We set the weight factors of each loss term of $\mathcal{L}_{\text{iSDF}}$ to be the same as iSDF.

## IV. EVALUATION

In this section, we demonstrate the proposed `AiSDF` framework on synthetic and real-world datasets. In Sec. IV-A, we provide the details of the experiment setting. In Sec. IV-B and Sec. IV-C, we show the qualitative and quantitative results to give a better understanding of reconstruction in various scenarios. In addition, we analyze the proposed `AiSDF` via an ablation study in Sec. IV-D. It should be noted that additional experiments are available in the supplementary video.

### A. Experiment setting

**Implementation details**. Following [9], we model `AiSDF` as a single MLP composed of four hidden layers. Specifically, both sampled surfel points $\mathcal{X}_\mathfrak{s}$ and non-surfel points $\mathcal{X}_\mathfrak{n}$ are transformed by positional embedding to make the network learn high-frequency regions. Then, embedded features are passed to four hidden layers that each composed of a linear layer and a softplus activation function. The output of the network is the SDF value and produces a mesh from this value by running a marching cube algorithm. `AiSDF` is implemented by PyTorch [39] and trained on a single RTX 3090 GPU. To optimize our `AiSDF`, we employ an AdamW optimizer [40] with a learning rate $1.3 \times 10^{-3}$ and a weight decay $1.2 \times 10^{-2}$. In addition, we utilize the same keyframes and iteration process as iSDF [9] for a fair comparison.

**Dataset**. We validate the proposed `AiSDF` on two datasets. 1) The ScanNet dataset [1] captured by an RGB-D camera contains $1,513$ scans of real-world indoor scenes with camera parameters, semantic labels, and surface reconstructions.
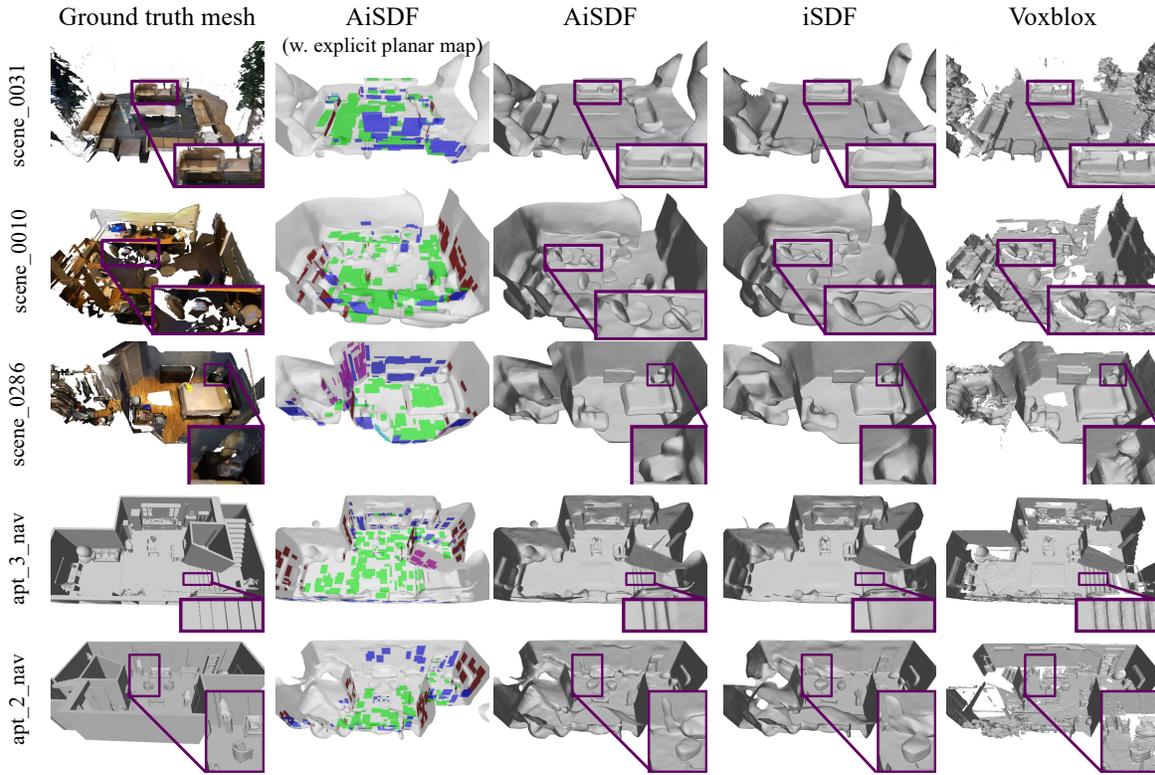
Fig. 5: **Qualitative evaluation on the ScanNet and ReplicaCAD datasets**. The purple boxes are close-up views to show the details of each scene. The second column presents the half-transparent explicit planar map composed of surfels overlaid on the mesh. Here, the green color denotes the surfels supported by the vertical Atlanta direction, and the other colors represent the surfels by the other horizontal Atlanta directions. Note that the size of the surfels may look smaller in the ReplicaCAD due to the different scales of scenes.

Among these scans, we use six sequences satisfying the Manhattan structure and three sequences following the Atlanta structure. 2) The ReplicaCAD dataset [16] is a photo-realistic 3D indoor scene reconstruction dataset, which is a recreated version of the 'FRL apartment' from the Replica dataset [41]. This dataset provides six scenes with variations of placement of large furniture, as well as small objects, for object rearrangement tasks. In our work, we use the same sequences used by iSDF [9] for a fair comparison. Specifically, we utilize two rooms, where each room structure has three sequences according to different trajectories with their own purpose, such as navigation, object reconstruction, and manipulation (six sequences in total).

**Comparison method**. We compare `AiSDF` with two methods. 1) Voxblox [17] is an algorithm for incrementally propagating wavefronts from updated TSDF voxels to create a voxel grid with the Euclidean signed distance. To conserve hardware resources, Voxblox uses voxels of large size and raycast grouping to accelerate integration. These techniques make it possible to create a map in real-time. We set the voxel size to $5.5cm$. 2) iSDF [9], which is our baseline, is a real-time SDF reconstruction method from a stream of posed depth images of room-scale environments.

**Metric**. For quantitative evaluation, we use the same three metrics as in iSDF [9]: SDF error, collision cost error, and gradient cosine distance.
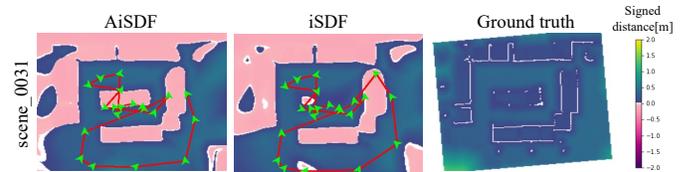


Fig. 6: **Visualization of SDF at a constant height**. The SDF is extracted at the end of the sequence.

### B. Qualitative evaluation

**Mesh reconstruction quality**. Figure 5 shows the mesh reconstruction results of `AiSDF` and comparison methods. Overall, Voxblox captures the details of even small objects, but it produces uneven surfaces. On the other hand, iSDF generates even surfaces, whereas has difficulty capturing the details of small objects. Unlike comparison methods, the proposed `AiSDF` reconstructs flat and complete surfaces while maintaining a certain level of detail thanks to Atlanta-aware sampling. For example, `AiSDF` reconstructs the complete shape of the sofa and the detail of the sofa armrest in $scene\_0031$ of ScanNet. For more complex scenes, such as $scene\_0010$, we can observe that the shape and details of two nearby chairs are clearly distinguished by `AiSDF`. In the $apt\_3\_nav$ of ReplicaCAD, `AiSDF` reconstructs the structure of stairs clearly compared to iSDF, which shows a collapsed structure. Note that in the case of Voxblox, they show the clear structure of stairs since they use most of the input frames for reconstruction, unlike keyframe-based `AiSDF` and iSDF (*i.e.*, the viewpoint covered by a few keyframes is

TABLE I: **Quantitative evaluation**. We highlight the best and second-best by **bold** and underline, respectively.

| Dataset | Scene | Voxblox [17] | | | iSDF [9] | | | AiSDF (w/o explicit map) | | | AiSDF (w. explicit map) | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | SDF | Collision | Gradient | SDF | Collision | Gradient | SDF | Collision | Gradient | SDF ↓ | Collision ↓ | Gradient ↓ |
| ScanNet | *scene_0004* | 6.352 | 0.048 | 0.215 | 5.355 | 0.037 | 0.134 | <u>4.317</u> | <u>0.029</u> | <u>0.125</u> | **4.258** | **0.028** | **0.125** |
| | *scene_0005* | 6.416 | 0.050 | 0.120 | 3.580 | 0.029 | 0.131 | <u>3.470</u> | <u>0.028</u> | <u>0.127</u> | **3.402** | **0.028** | **0.124** |
| | *scene_0009* | 5.608 | 0.041 | 0.112 | 4.515 | 0.034 | 0.161 | <u>3.726</u> | <u>0.028</u> | <u>0.149</u> | **3.593** | **0.027** | **0.145** |
| | *scene_0010* | 4.850 | 0.055 | 0.347 | 3.822 | 0.032 | 0.154 | <u>3.494</u> | <u>0.029</u> | <u>0.149</u> | **3.285** | **0.027** | **0.147** |
| | *scene_0030* | 5.500 | 0.042 | 0.942 | 3.188 | 0.026 | 0.099 | <u>2.910</u> | <u>0.023</u> | <u>0.095</u> | **2.734** | **0.022** | **0.092** |
| | *scene_0031* | 11.512 | 0.104 | 0.317 | 4.891 | 0.036 | **0.126** | <u>4.721</u> | <u>0.034</u> | 0.133 | **4.646** | **0.034** | <u>0.132</u> |
| ReplicaCAD | *apt_2_nav* | 4.672 | 0.038 | 0.158 | 4.349 | 0.036 | 0.156 | <u>3.585</u> | <u>0.029</u> | <u>0.150</u> | **3.317** | **0.027** | **0.144** |
| | *apt_2_mnp* | 9.744 | 0.073 | 0.396 | 7.707 | 0.051 | 0.339 | <u>7.464</u> | <u>0.050</u> | <u>0.327</u> | **7.123** | **0.047** | **0.298** |
| | *apt_2_obj* | 8.178 | 0.071 | 0.237 | 4.604 | 0.034 | <u>0.144</u> | <u>4.294</u> | <u>0.031</u> | 0.146 | **4.119** | **0.030** | **0.143** |
| | *apt_3_nav* | 4.626 | 0.036 | 0.128 | 3.713 | 0.029 | 0.119 | <u>3.288</u> | <u>0.026</u> | <u>0.116</u> | **2.812** | **0.022** | **0.107** |
| | *apt_3_mnp* | 8.040 | 0.076 | 0.232 | 6.336 | 0.044 | <u>0.192</u> | <u>5.979</u> | <u>0.042</u> | 0.194 | **5.818** | **0.040** | **0.187** |
| | *apt_3_obj* | 5.802 | 0.073 | 0.282 | <u>4.368</u> | 0.032 | **0.133** | 4.406 | <u>0.032</u> | 0.138 | **4.294** | **0.031** | <u>0.134</u> |

limited). In addition, We can also produce the slice map of the reconstructed SDF, as shown in Fig. 6.

**Explicit planar map**. In addition to mesh reconstruction, AiSDF can generate an explicit 3D planar map composed of Atlanta-aware surfels (see the second column in Fig. 5). In particular, explicit maps mainly contain planes of scene structure (*e.g.*, walls and floors), while also involving relatively small planes (*e.g.*, stairs and objects). Specifically, even in *scene_0010* containing many objects, and *scene_0286* with Atlanta structure, AiSDF robustly extracts various planes supported by Atlanta directions. Furthermore, these explicit maps give additional geometric cues to reconstructed mesh. For example, in *apt_3_nav* of ReplicaCAD, AiSDF reconstructs the structure of stairs at the mesh level by implicitly utilizing the extracted surfels in surfel loss. We believe that explicit planar maps can also serve as an auxiliary resource for various downstream (*e.g.*, footstep planners for humanoids [42] and Roomplan scans [43]). It should be noted that additional qualitative results are available in the supplementary video.

## C. Quantitative evaluation

Table I shows the quantitative results. Following [9], for six scenes providing ground truth SDF, we measure three metrics for 200k sampled points on all rays computed by randomly sampled pixel coordinates from the frames. In addition, we consider explicit planar maps together to measure the performance of AiSDF. To this end, we compute explicit SDF values by calculating the closest distance from each of the 200k sampled points to surfel maps. Then, we compare it to the implicit SDF value and choose the smaller value to measure SDF and Collision metric with ground truth values. For the sampled points using explicit SDF values, we replace implicit gradient vectors with explicit gradient vectors, obtained by computing the difference between the points and the corresponding surfel points to measure gradient metric.

Overall, AiSDF outperforms comparison methods on the ScanNet and ReplicaCAD datasets. In particular, when we consider the estimated explicit planar map of AiSDF together, it shows better performance. Concretely, in the case of Voxblox, since it uses large voxel size for real-time, its performance is much worse than iSDF and AiSDF. Regarding iSDF, it is difficult to predict an accurate SDF because of the noisy depth of homogeneous regions. In addition, iSDF

TABLE II: **Ablation study of Atlanta-aware sampling**.

| Surfel mask | Surfel loss | scene_0010 | | | apt_2_nav | | |
|---|---|---|---|---|---|---|---|
| | | SDF | Collision | Gradient | SDF | Collision | Gradient |
| × | × | 3.822 | 0.032 | 0.154 | 4.349 | 0.036 | 0.156 |
| ✓ | × | 3.667 | 0.030 | 0.154 | 4.042 | 0.033 | 0.152 |
| × | ✓ | 3.760 | 0.031 | 0.151 | 3.959 | 0.032 | 0.155 |
| ✓ | ✓ | **3.494** | **0.029** | **0.149** | **3.585** | **0.029** | **0.150** |

TABLE III: **Runtime of AiSDF**. For the sequences of ScanNet used in the quantitative evaluation, we compute the average time for each module and then round up (unit: ms).

| AF estimation | Extract surfels | Sampling | Compute bound | forward | backward | **Total** |
|---|---|---|---|---|---|---|
| 2 | 54 | 5 | 7 | 6 | 14 | **88** |

may have difficulty assigning the proper number of samples to complex regions, whereas AiSDF samples more points on complex regions by utilizing the Atlanta-aware surfels. On the other hand, AiSDF obtains more accurate results by using Atlanta-aware surfels that are robust to noise and allow us adaptive sampling according to surfels or non-surfels.

## D. Analysis

**Ablation study**. We conduct the ablation study to demonstrate the effectiveness of Atlanta-aware sampling: surfel mask and surfel loss (see Table II). Surfel mask $\mathbf{M_s}$ provides a criterion for dividing a scene into complex and planar regions. Thus, training the model with surface masks allows us to sample more points on complex areas, resulting in improved performance. In addition, we observe that regularizing the network with surfel loss leads to performance improvement. Specifically, by utilizing surfel points $\mathcal{X_s}$, we can directly constrain the planar regions and obtain a more accurate approximated bound and gradient. Consequently, AiSDF using surface mask and surfel loss together shows the best performance.

**Runtime and memory**. Table III shows the running time for each module in AiSDF. Unfortunately, AiSDF does not work in real-time, but it is sufficient to reconstruct the scene in an online process ($\leq 100ms$) when we consider this process that corresponds to the back-end part and takes only keyframes, not each frame. In terms of memory, AiSDF reconstructs the scene with 1MB of network parameters as in iSDF, while effectively representing the scene with only hundreds of KB of the explicit planar map.

**Limitation**. AiSDF proposes a novel way to combine the structural regularity of Atlanta structures with the implicit SDF reconstruction framework working online, which is our

main contribution. Contrary to this novelty, `AiSDF` has several limitations as a pioneer work. To maintain reasonable runtime, `AiSDF` independently extracts Atlanta-aware surfels for each keyframe. Thus, we cannot provide a complete and unified explicit planar map of a given indoor scene. In addition, the current `AiSDF` does not fully exploit the estimated Atlanta structures in the implicit neural representation; we currently enforce surfel loss on surfel points to learn SDF. From this point of view, we believe that encoding Atlanta-aware surfels itself to SDF could be an interesting research direction.

## V. Conclusion

Under the AW assumption, we have proposed the novel `AiSDF`, a structure-aware online SDF reconstruction framework of a given indoor scene. To fully exploit the inherent property of indoor scenes (*i.e.*, structural regularity), we estimate the underlying Atlanta structure in the form of dominant Atlanta directions within a continual framework. This structural understanding allows us to extract Atlanta-aware surfels, which explicitly play a role as a 3D planar map. Moreover, Atlanta-aware surfels provide a criterion to adaptively sample points and make `AiSDF` implicitly enforce surfel loss on surfel points. We seamlessly integrate this structural understanding inside the online SDF reconstruction framework. As a result, experiments demonstrate that `AiSDF` can reconstruct the details of the scene with overall structure while extracting the lightweight explicit planar map.

## References

[1] A. Dai, A. X. Chang, M. Savva, M. Halber, T. Funkhouser, and M. Nießner, "Scannet: Richly-annotated 3d reconstructions of indoor scenes," in *CVPR*, 2017.

[2] R. Cabral and Y. Furukawa, "Piecewise planar and compact floorplan reconstruction from images," in *CVPR*, 2014.

[3] T. Schöps, T. Sattler, and M. Pollefeys, "Surfelmeshing: Online surfel-based mesh reconstruction," *IEEE TPAMI*, 2019.

[4] P. Mittal, Y.-C. Cheng, M. Singh, and S. Tulsiani, "AutoSDF: Shape Priors for 3D Completion, Reconstruction and Generation," in *CVPR*, 2022.

[5] Y. Jiang, D. Ji, Z. Han, and M. Zwicker, "SDFDiff: Differentiable Rendering of Signed Distance Fields for 3D Shape Optimization," in *CVPR*, 2020.

[6] M. Zucker, N. Ratliff, A. D. Dragan, M. Pivtoraiko, M. Klingensmith, C. M. Dellin, J. A. Bagnell, and S. S. Srinivasa, "Chomp: Covariant hamiltonian optimization for motion planning," *IJRR*, 2013.

[7] E. Sucar, S. Liu, J. Ortiz, and A. J. Davison, "iMAP: Implicit mapping and positioning in real-time," in *ICCV*, 2021.

[8] Z. Zhu, S. Peng, V. Larsson, W. Xu, H. Bao, Z. Cui, M. R. Oswald, and M. Pollefeys, "Nice-SLAM: Neural implicit scalable encoding for SLAM," in *CVPR*, 2022.

[9] J. Ortiz, A. Clegg, J. Dong, E. Sucar, D. Novotny, M. Zollhoefer, and M. Mukadam, "iSDF: Real-Time Neural Signed Distance Fields for Robot Perception," in *RSS*, 2022.

[10] J. M. Coughlan and A. L. Yuille, "Manhattan world: Compass direction from a single image by bayesian inference," in *ICCV*, 1999.

[11] G. Schindler and F. Dellaert, "Atlanta World: An Expectation Maximization Framework for Simultaneous Low-Level Edge Grouping and Camera Calibration in Complex Man-Made Environments," in *CVPR*, 2004.

[12] S. Gupta, P. Arbelaez, and J. Malik, "Perceptual organization and recognition of indoor scenes from rgb-d images," in *CVPR*, 2013.

[13] L. Carlone, R. Tron, K. Daniilidis, and F. Dellaert, "Initialization techniques for 3D SLAM: a survey on rotation estimation and its use in pose graph optimization," in *ICRA*, 2015.

[14] P. Kim, B. Coltin, and H. J. Kim, "Low-drift visual odometry in structured environments by decoupling rotational and transnational motion," in *ICRA*, 2018.

[15] K. Joo, P. Kim, M. Hebert, I. S. Kweon, and H. J. Kim, "Linear RGB-D SLAM for structured environments," *IEEE TPAMI*, 2021.

[16] A. Szot, A. Clegg, E. Undersander, E. Wijmans, Y. Zhao, J. Turner, N. Maestre, M. Mukadam, D. Chaplot, O. Maksymets, A. Gokaslan, V. Vondrus, S. Dharur, F. Meier, W. Galuba, A. Chang, Z. Kira, V. Koltun, J. Malik, M. Savva, and D. Batra, "Habitat 2.0: Training home assistants to rearrange their habitat," in *NeurIPS*, 2021.

[17] H. Oleynikova, Z. Taylor, M. Fehr, R. Siegwart, and J. Nieto, "Voxblox: Incremental 3d euclidean signed distance fields for on-board mav planning," in *IROS*, 2017.

[18] J. Straub, O. Freifeld, G. Rosman, J. J. Leonard, and J. W. Fisher, "The Manhattan frame model—Manhattan world inference in the space of surface normals," *IEEE TPAMI*, 2017.

[19] K. Joo, T.-H. Oh, I. S. Kweon, and J.-C. Bazin, "Globally optimal inlier set maximization for Atlanta world understanding," *IEEE TPAMI*, 2019.

[20] N. Silberman, D. Hoiem, P. Kohli, and R. Fergus, "Indoor segmentation and support inference from rgbd images," in *ECCV*, 2012.

[21] W. Choi, Y.-W. Chao, C. Pantofaru, and S. Savarese, "Understanding indoor scenes using 3d geometric phrases," in *CVPR*, 2013.

[22] H. Wildenauer and A. Hanbury, "Robust camera self-calibration from monocular images of manhattan worlds," in *CVPR*, 2012.

[23] L. Mescheder, M. Oechsle, M. Niemeyer, S. Nowozin, and A. Geiger, "Occupancy networks: Learning 3d reconstruction in function space," in *CVPR*, 2019.

[24] 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," in *ECCV*, 2020.

[25] J. J. Park, P. Florence, J. Straub, R. Newcombe, and S. Lovegrove, "Deepsdf: Learning continuous signed distance functions for shape representation," in *CVPR*, 2019.

[26] R. Po, Z. Dong, A. W. Bergman, and G. Wetzstein, "Instant continual learning of neural radiance fields," in *ICCVW*, 2023.

[27] Z. Murez, T. v. As, J. Bartolozzi, A. Sinha, V. Badrinarayanan, and A. Rabinovich, "Atlas: End-to-end 3d scene reconstruction from posed images," in *ECCV*, 2020.

[28] J. Sun, Y. Xie, L. Chen, X. Zhou, and H. Bao, "Neuralrecon: Real-time coherent 3d reconstruction from monocular video," in *CVPR*, 2021.

[29] Z. Yan, Y. Tian, X. Shi, P. Guo, P. Wang, and H. Zha, "Continual neural mapping: Learning an implicit scene representation from sequential observations," in *CVPR*, 2021.

[30] A. Dai and M. Nießner, "Neural poisson: Indicator functions for neural fields," *arXiv*, 2022.

[31] L. Yariv, J. Gu, Y. Kasten, and Y. Lipman, "Volume rendering of neural implicit surfaces," in *NeurIPS*, 2021.

[32] 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," *NeurIPS*, 2021.

[33] D. Azinović, R. Martin-Brualla, D. B. Goldman, M. Nießner, and J. Thies, "Neural rgb-d surface reconstruction," in *CVPR*, 2022.

[34] J. Wang, P. Wang, X. Long, C. Theobalt, T. Komura, L. Liu, and W. Wang, "Neuris: Neural reconstruction of indoor scenes using normal priors," in *ECCV*, 2022.

[35] Z. Yu, S. Peng, M. Niemeyer, T. Sattler, and A. Geiger, "Monosdf: Exploring monocular geometric cues for neural implicit surface reconstruction," *NeurIPS*, 2022.

[36] H. Guo, S. Peng, H. Lin, Q. Wang, G. Zhang, H. Bao, and X. Zhou, "Neural 3D Scene Reconstruction with the Manhattan-world Assumption," in *CVPR*, 2022.

[37] K. Joo, T.-H. Oh, J. Kim, and I. S. Kweon, "Robust and globally optimal Manhattan frame estimation in near real time," *IEEE TPAMI*, 2018.

[38] A. Gropp, L. Yariv, N. Haim, M. Atzmon, and Y. Lipman, "Implicit geometric regularization for learning shapes," in *ICML*, 2020.

[39] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga, *et al.*, "Pytorch: An imperative style, high-performance deep learning library," *NeurIPS*, 2019.

[40] I. Loshchilov and F. Hutter, "Decoupled weight decay regularization," *ICLR*, 2018.

[41] S. *et al.*, "The Replica dataset: A digital replica of indoor spaces," *arXiv*, 2019.

[42] R. J. Griffin, G. Wiedebach, S. McCrory, S. Bertrand, I. Lee, and J. Pratt, "Footstep planning for autonomous walking over rough terrain," in *HUMANOIDS*, 2019.

[43] "Roomplan, Apple ARKit," https://machinelearning.apple.com/research/roomplan.