

ALTO: Alternating Latent Topologies for Implicit 3D Reconstruction

Zhen Wang^{1*} Shijie Zhou^{1*} Jeong Joon Park² Despoina Paschalidou²
 Suya You³ Gordon Wetzstein² Leonidas Guibas² Achuta Kadambi¹

¹University of California, Los Angeles ²Stanford University ³DEVCOM Army Research Laboratory

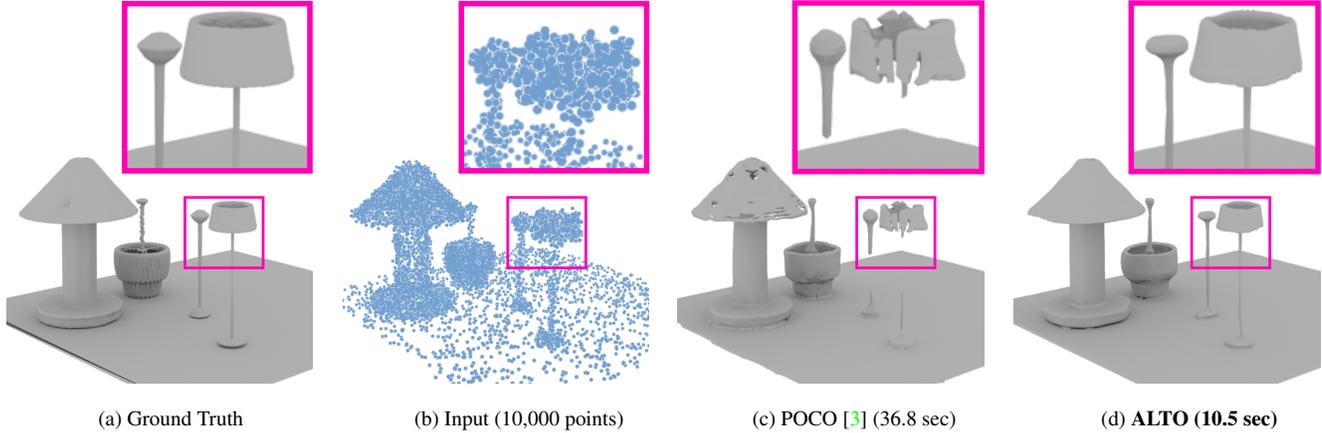


Figure 1. **Rethinking latent topologies for fast and detailed implicit 3D reconstructions.** Recent work (POCO CVPR’22 [3]) has used latent encodings for each point to preserve 3D detail. We introduce ALTO, which can alternate between latent topologies like grid latents and point latents to speed up inference and recover more detail, like the 3D reconstruction of a thin lamp-post. Scene from [49].

Abstract

This work introduces alternating latent topologies (ALTO) for high-fidelity reconstruction of implicit 3D surfaces from noisy point clouds. Previous work identifies that the spatial arrangement of latent encodings is important to recover detail. One school of thought is to encode a latent vector for each point (point latents). Another school of thought is to project point latents into a grid (grid latents) which could be a voxel grid or triplane grid. Each school of thought has tradeoffs. Grid latents are coarse and lose high-frequency detail. In contrast, point latents preserve detail. However, point latents are more difficult to decode into a surface, and quality and runtime suffer. In this paper, we propose ALTO to sequentially alternate between geometric representations, before converging to an easy-to-decode latent. We find that this preserves spatial expressiveness and makes decoding lightweight. We validate ALTO on implicit 3D recovery and observe not only a performance improvement over the state-of-the-art, but a runtime improvement of 3-10 \times . Project website at

<https://visual.ee.ucla.edu/alto.html/>.

1. Introduction

Reconstructing surfaces from noisy point clouds is an active problem in 3D computer vision. Today, conditional neural fields offer a promising way to learn surfaces from noisy point clouds. Alternatives like voxel regression or mesh estimation are limited by cubic complexity and the requirement of a mesh template, respectively. Recent work has successfully used conditional neural fields to reconstruct 3D surfaces as an occupancy function. A conditional neural field takes as input a query coordinate and conditions this on a latent representation, e.g., feature grids. The spatial expressiveness of the latent representation impacts the overall surface reconstruction quality.

To achieve spatial expression, a neural field is conditioned on a latent space of features (**latents**) from the conditional input. In 3D surface reconstruction the input point cloud is transformed into latents arranged in some topological structure. **Point latents** occur when each point in the input point cloud is assigned a latent vector [3]. **Triplane latents** are formed when point latents are projected into a

*Equal contribution.

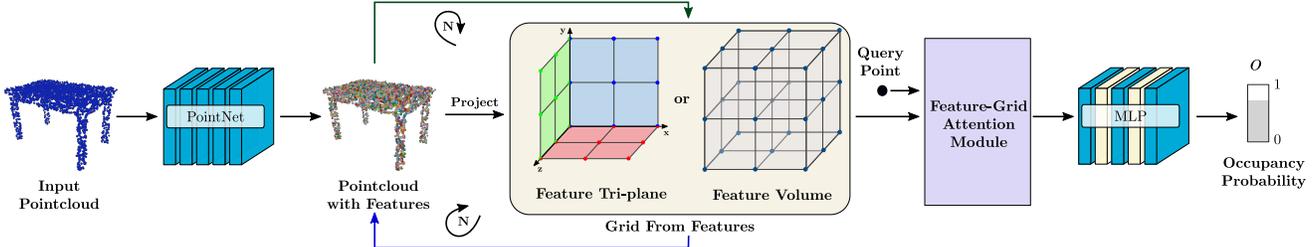


Figure 2. **An overview of our method.** Given input surface points, we obtain an implicit occupancy field with iterative alternation between features in the forms of points and 2D or 3D grids (Sec. 3.2). Then we decode the occupancy values for query points with a learned attention-based interpolation from neighboring grids (Sec. 3.3).

3-axis grid [37, 49]. The triplane latent is not as spatially expressive as freeform points, but the lower spatial complexity makes it easier to decode. **Voxel latents** are another type of grid latent where latents are arranged in a feature volume [49, 62].

To reconstruct detailed surfaces, recent state-of-the-art methods try to preserve point latents as long as possible. Because point latents are spatially expressive, methods based on point latents are considered state-of-the-art for detailed surface reconstruction [3, 16]. However, using point latents in this way has some tradeoffs. It is difficult to correlate a query with the unstructured topology of a point-based latent space, placing a burden on the decoder. Results from POCO [3] are shown in Fig. 1 where runtime and high-quality detail like thin lampposts remain out of reach.

In this paper, we seek to blend the upside of different latent topologies, while minimizing downside. We present an alternating latent topology (ALTO) method. In contrast to previous work, our method does not stay with either point [3] or grid latents [49], but instead alternates back and forth between point and grid latents before converging to a final grid for ease-of-decoding.

Our method is general. We can plug-in the ALTO component to existing grid-based conditional models [8, 49] to boost detail recovery. While we have shown that our method can generate occupancy fields, we expect gain of high-fidelity details for other neural fields, such as semantic or affordance fields [30, 66], where similar conditional techniques can be adopted.

We summarize our **contributions** as follows:

- We introduce an iterative technique to blend the strengths of different latent topologies for high-fidelity conditional neural fields generation.
- We propose an attention-based decoder that replaces naive linear interpolation of feature-grids or computationally expensive point-wise attention while keeping compute burden in check.
- We demonstrate performance and runtime improvements over the highest-quality previous method [3], as

well as performance improvements over all other baselines.

2. Related Work

In this section, we discuss the most relevant literature on learning-based 3D reconstruction methods. Based on their output, existing learning-based approaches can be categorized as implicit or explicit-based representations. In this work, we primarily focus on implicit-based representations as they are closely related to our method.

2.1. Explicit Representations

A common shape representation is 2.5D depth maps, which can be inferred using 2D CNNs [15, 26, 27, 32, 47, 74, 75, 78]. However, 2.5D depth maps cannot capture the full 3D geometry. In contrast, voxels [4, 9, 12, 19, 40, 52, 60, 69–71] naturally capture 3D object geometry, by discretizing the shape into a regular grid. As voxel-based methods exhibit cubic space complexity that results in high memory and computation requirements, several works tried to circumvent this with efficient space partitioning techniques [25, 40, 53, 54, 63]. Although these methods allow for increasing the voxel resolution and hence capturing more complex geometries, their application is still limited. Recently, a promising new direction explored learning grid deformations to better capture geometric details [20]. An alternative representation relies on pointclouds. Point-based approaches [1, 18, 29, 51, 65, 73] discretize the 3D space using points and are more light-weight and memory efficient. However, as they lack surface connectivity, they require additional post-processing steps (e.g. using Poisson Surface Reconstruction [34]) for generating the final mesh. Instead, mesh-based methods [6, 13, 22, 24, 31, 36, 45, 67, 73] naturally yield smooth surfaces but they typically require a template mesh [67], which makes scaling them to arbitrarily complex topologies difficult. Other works, proposed to also represent the geometry as an atlas of mappings [14, 24, 39], which can result in non-watertight meshes. To address limitations with learning explicit representations, implicit models [7, 41, 46, 58] emerged as an alternative more com-

pact representation that yield 3D geometries at infinitely high resolutions using iso-surfacing operations (i.e. marching cubes). In this work, we capture 3D geometries implicitly, using an occupancy field [41], as it faithfully can capture complex topologies.

2.2. Neural Implicit Representations

Unlike explicit representations that discretize the output space using voxels, points or mesh vertices, implicit representations represent the 3D shape and appearance implicitly, in the latent vectors of a neural network that learns a mapping between a query point and a context vector to either a signed distance value [2, 23, 42, 46, 61] or a binary occupancy value [7, 41, 59]. However, while these methods typically rely on a single global latent code, they cannot capture local details and struggle scaling to more complicated geometries. To address this, several works [56, 57, 72] explored pixel-aligned implicit representations, that rely on both global and local image features computed along a viewing direction. While, these approaches are able to capture fine-grained geometric details, they rely on features that are computed from images, hence are limited to image-based inputs with known camera poses.

Our work falls in the category of methods that perform 3D reconstructions from points. Among the first to explore this direction were [8, 28, 49]. To increase the expressivity of the underlying representation and to be able to capture complex geometries, instead of conditioning on a global latent code, these works condition on local per-point features. For example, Jiang et al. [28] leverage shape priors by conditioning on a patch-based representation of the point cloud. Other works [37, 49, 62], utilize grid-based convolutional features extracted from feature planes [37], feature volumes [8, 62] or both [49]. An alternative line of work, [68] introduce a test-time optimization mechanism to refine the per-point features predicted on a feature volume. Concurrently, POCO [3], propose to estimate per-point features, which are then refined based on the per-point features of their neighboring points using an attention-based mechanism. A similar idea is also explored in Points2Surf [17] that introduces a patch-based mechanism to decide the sign of the implicit function. Our work is closely related to POCO [3] that can faithfully capture higher-frequency details due its point-wise latent coding. However, the lack of a grid-like structure places extra complexity on the attention-based aggregation module, which results in a higher computation cost. AIR-Net [21] applied local attention to reduce the computation but limits its operating range to objects. In this work, we propose conditioning on a hybrid representation of points and grid latents. In particular, instead of fusing points and grids, we demonstrate that it is the point and grid *alternation* between points and grids that enables recovery of more detail than POCO [3], while reducing com-

pute time by an order of magnitude.

2.3. Obtaining Implicit Fields from Images

As the previous methods require 3D supervision, several recent works propose combining implicit representations with surface [44, 77] or volumetric [43] rendering techniques in order to learn 3D object geometry and texture from images. Among the most extensively used implicit representations are Neural Radiance Fields (NeRF) [43] that combine an implicit representation with volumetric rendering to perform a novel view synthesis. In particular, they employ a neural network that maps a 3D point along a viewing direction to a color and a density value. NeRFs are trained using only posed images. On the other hand [44, 76] combined occupancy fields [41] and signed distance functions (SDFs) with surface rendering in order to recover the geometry and appearance of 3D objects.

Though we only demonstrate a case study in occupancy field in this paper, our proposed method is general and can be readily plugged into all kinds of other neural fields for better 3D point feature encoding. We leave that as future exploration.

3. Method

There are three insights that motivate our approach: (1) conditioning on the right topology of the latent space is important; (2) previous neural fields for surface reconstruction condition on point or grid latents; (3) both point and grid latents have complementary strengths and weaknesses. Point latents are more spatially expressive but grid latents are easier to decode into a surface.

It might seem like a simple concatenation of point and grid latents would be sufficient. The problem is that point latents remain difficult to decode (even if concatenated with a grid latent). Therefore, our insight is to alternate between grid and point latents, and converge to a grid latent. For feature triplane latents, the alternation also permits communication between the individual planes, which in previous, grid-based works would have been fed into independent hourglass U-Nets [49].

In this section, we introduce a version of ALTO as a point-grid alternating U-Net. An overview of the method is shown in Fig. 2. We demonstrate how the convolutional grid form is learned in Sec. 3.1, our point-grid alternating network in Sec. 3.2, our attention-based decoder in Sec. 3.3 and training and inference in more detail in Sec. 3.4.

3.1. Convolutional Feature Grids

The input to our method is a noisy un-oriented point cloud $\mathcal{P} = \{\mathbf{p}_i \in \mathbb{R}^3\}_{i=1}^S$, where S is the number of input points. We first use a shallow Point-Net [50] to obtain the initial point features as in ConvONet [49]. These point

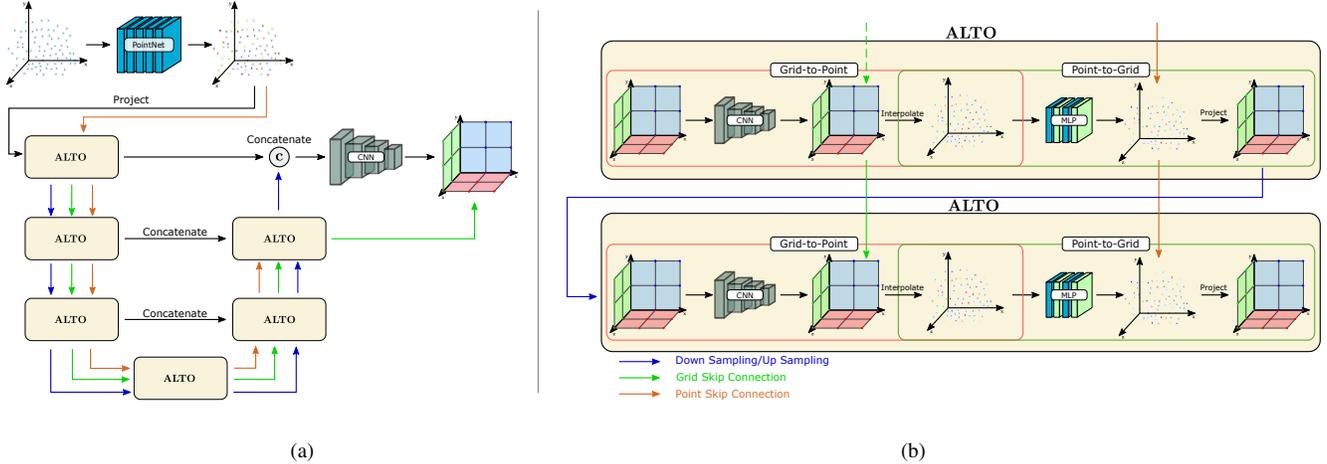


Figure 3. **An illustration of our ALTO encoder.** (a) As an example, we show the ALTO block instantiated by alternating between two latent topologies: point and triplanes via an “in-network” fashion, i.e. within each level of an hourglass framework U-Net. ‘Concatenate’ refers to concatenation of the ALTO block output triplane in the downsampling stage and the ALTO block input triplane in the corresponding upsampling stage. (b) We expand on ALTO block to illustrate the sequential grid-to-point and point-to-grid conversion. There are skip connections for both point and grid features between two consecutive levels in the ALTO U-Net.

features are then projected into three 2D grids (feature triplanes) $\{\mathbf{\Pi}_{xy}, \mathbf{\Pi}_{xz}, \mathbf{\Pi}_{yz}\} \in \mathbb{R}^{H \times W \times d}$ or 3D grids (feature volumes) $\mathbf{V} \in \mathbb{R}^{H \times W \times D \times d}$ before feeding into a 2D or 3D convolutional hourglass (U-Net) networks [10, 55]. d is the number of feature channels. For feature volumes, we set $H = W = D = 64$ due to memory overhead of 3D-CNN and for feature triplanes, H and W can be set as high as 128 depending on the task.

3.2. ALTO Latent to Blend Grid and Point Latents

Without loss of generality, we demonstrate ALTO in the context of blending grid and point latents. Note that naive concatenation of latents would not work, as the point latents are difficult to decode. The goal is to use ALTO to blend point latent characteristics into a grid latent via alternation. The alternating block is illustrated in Fig. 3 and incorporated into a U-Net architecture. At each alternation, we first do grid-to-point conversion where convolutional grid features are transformed into point features, followed by point-to-grid conversion where extracted point features are transformed back into grid features for next alternation.

Grid-to-Point Conversion: At each alternation, to aggregate local neighborhood information, we use convolutional operations for the grid features. We then project each point p orthographically onto the canonical planes and query the feature values through bilinear interpolation for 2D grid and trilinear interpolation for 3D grid. For triplane latents, we sum together the interpolated features from each individual plane.

Point-to-Grid Conversion: At each alternation, given the interpolated point features, we then process point features with an MLP in order to model individual point fea-

ture with finer granularity. For feature triplanes grid form, an MLP also gives an additional benefit of having individual plane features communicate with each other. The MLP is implemented with two linear layers and one ReLU non-linearity. Projected point features falling within the same pixel or voxel cell will be aggregated using average pooling. If using triplane latents with each plane discretized at $H \times W$, this results in planar features with dimensionality $H \times W \times d$ or if using voxel latents we obtain dimensionality $H \times W \times D \times d$, where d is the number of features.

We also adopt skip connections for both point features and grid features between two consecutive ALTO blocks, as illustrated in Fig. 3b. The alternation needs to be implemented carefully to minimize runtime. Naively, we can alternate between triplanes or voxel latents using a U-Net and point latents using MLP, but that would require multiple network passes. Instead, we incorporate the point-grid alternating *inside* each block of a U-Net, i.e. replacing original convolution-only block with ALTO block. We call this single U-Net, the ALTO U-Net. This also enables point and grid features blended at multiple scales and the number of alternation blocks depends on the number of levels of U-Net.

3.3. Decoding ALTO latents using Attention

As discussed in the previous section, ALTO provides a way to get a single latent that blends characteristics of different topologies. The advantage is that the final latent that ALTO converges to (hereafter, ALTO latent) can take on the topology that is easier to decode.

For example, in the case of using ALTO to blend point and grid characteristics, we would like ALTO to converge to

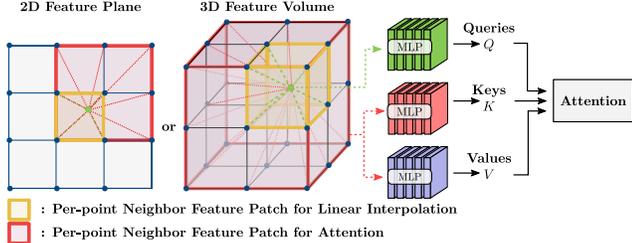


Figure 4. **Attention-based decoder on neighboring grids (2D or 3D)**. To obtain features of each query point for decoding occupancy value, we use learned interpolation from neighboring grids that improves occupancy prediction, while being more efficient than expansive point-wise attention mechanism (e.g. POCO [3]).

a final output in the simpler grid topology. Then, given the ALTO latent in grid form and any query point $\mathbf{q} \in \mathbb{R}^3$ in 3D space, our goal is to decode the learned feature and estimate the occupancy probability of each query point. ALTO benefits from attention on the decoder side. The ALTO latent is in grid form, but has spatial expressivity coming from the blended-in point latents. Standard grid latent decoding, e.g., bi-/tri-linear interpolation used in previous work [37, 49] would not preserve this spatial expressivity.

To decode an ALTO latent, we propose an efficient attention-based mechanism to replace the previous approach of linear interpolation on feature grids. While attention is not new, we leverage *grid latent attention* to avoid heavy runtime issues of *point latent attention* [3] that applies attention over a point-wise 3D neighborhood. As illustrated in Fig. 4, we consider the nearest grids (indices denoted as $\mathcal{N}(\mathbf{q})$), where $|\mathcal{N}| = 9$ for triplane representation and $|\mathcal{N}| = 27$ for volume representation, we call these areas as per-point neighbor feature patches $\mathbf{C}_{\{i \in \mathcal{N}\}}$. We define the query Q , key K , and value V for our attention as follows:

$$\begin{aligned} Q &= \text{MLP}(\psi(\mathbf{q})), \\ K &= \text{MLP}(\mathbf{C}_{\{i \in \mathcal{N}(\mathbf{q})\}}), \\ V &= \text{MLP}(\mathbf{C}_{\{i \in \mathcal{N}(\mathbf{q})\}}), \end{aligned} \quad (1)$$

where $\psi(\mathbf{q})$ is the linear interpolated feature value of the query points. Additionally, we compute the displacement vector $\mathbf{d} \in \mathbb{R}^2$ or \mathbb{R}^3 which represents the spatial relationship between the projected query point coordinate and the nearest feature grid points. We use the subtraction relation [80] in our attention scoring function:

$$A = \text{softmax}(\text{MLP}((Q - K) + \gamma(\mathbf{d}))), \quad (2)$$

where $\gamma(\mathbf{d})$ works as a learnable positional encoding. In our implementation, γ is an MLP with two fully-connected layers and activated by ReLU. We compute the attention-based interpolated per-point feature F as:

$$F = A \odot (V + \gamma(\mathbf{d})). \quad (3)$$

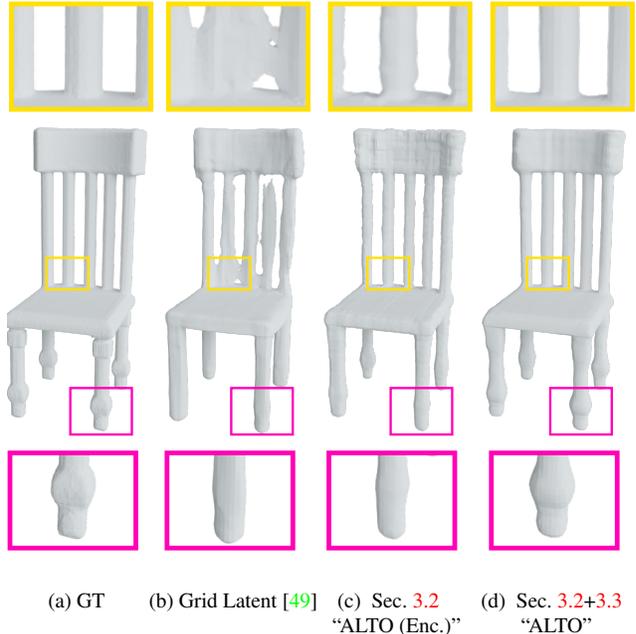


Figure 5. **Ablation analysis on ShapeNet**. Note the top inset showing the poles in the chair back (yellow). ALTO (Enc.) is ALTO (Encoder Only) and uses the latent space encoding proposed in Sec. 3.2 with a standard decoder. The full ALTO method includes also the attention-based decoder in Sec. 3.3.

Note that the same positional encoding from above is added to V and \odot denotes the element-wise product operation. For the case of triplane representation, we use a single-head attention to extract the feature F from each individual plane. The per-triplane features are then concatenated and used for the occupancy prediction. For the case of the volume representation, we use multi-head attention for h independently learned subspaces of Q , K , and V , where h is a hyperparameter varying based on the experiment. Additional details are provided in the supplementary.

Finally, we predict the occupancy of \mathbf{q} using a small fully-connected occupancy network:

$$f_{\theta}(F) \rightarrow [0, 1]. \quad (4)$$

The network f_{θ} consists of several ResNet blocks as in [49]. The major difference to the original occupancy decoder in [49] is that we do not bring in the absolute 3D coordinate of \mathbf{q} as input since it theoretically breaks the translational equivalence property.

This concludes our description of our latent space encoding and attention decoding. In Fig. 5 observe that the architectures we have proposed progressively improve detail from a standard grid formulation.

Method	input points 3K				input points 1K				input points 300			
	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow
ONet [41]	0.761	0.87	0.891	0.785	0.772	0.81	0.894	0.801	0.778	0.80	0.895	0.806
ConvONet [49]	0.884	0.44	0.938	0.942	0.859	0.50	0.929	0.918	0.821	0.59	0.907	0.883
POCO [3]	0.926	0.30	0.950	0.984	0.884	0.40	0.928	0.950	0.808	0.61	0.892	0.869
ALTO (Encoder Only)	0.931	0.30	0.950	0.981	0.889	0.39	0.932	0.951	0.842	0.52	0.908	0.903
ALTO	0.930	0.30	0.952	0.980	0.905	0.35	0.940	0.964	0.863	0.47	0.922	0.924

Table 1. **Performance on ShapeNet with various point density levels.** Input noisy point cloud with 3K, 1K and 300 input points from left to right. ALTO is our proposed method and ALTO (Encoder only) is an ablation that uses only our encoder with a non-attention based decoder.

3.4. Training and Inference

At training time, we uniformly sample query points \mathcal{Q} and minimize the binary cross-entropy loss between the predicted occupancy value and ground-truth occupancy values written as:

$$\mathcal{L}(\hat{o}_{\mathbf{q}}, o_{\mathbf{q}}) = - \sum_{\mathbf{q} \in \mathcal{Q}} [o_{\mathbf{q}} \log(\hat{o}_{\mathbf{q}}) + (1 - o_{\mathbf{q}}) \log(1 - \hat{o}_{\mathbf{q}})] \quad (5)$$

Our model is implemented in PyTorch [48] and uses the Adam optimizer [35] with a learning rate of 10^{-4} . During inference, we use a form of Marching Cubes [38] to obtain the mesh.

4. Experimental Evaluation

4.1. Datasets, Metrics, and Baselines

Object Level Datasets: For evaluation on object-level reconstruction, we use ShapeNet [5]. In particular, ShapeNet [5] contains watertight meshes of object shapes in 13 classes. For fair comparison, we use the same train/val splits and 8500 objects for testing as described in [3, 49]. Points are obtained by randomly sampling from each mesh and adding Gaussian noise with zero mean and standard deviation of 0.05.

Scene-Level Datasets: For scene level evaluation, we use Synthetic Rooms dataset [49] and ScanNet-v2 [11]. In total, we use 5000 synthetic room scenes with walls, floors and ShapeNet objects randomly placed together. We use an identical train/val/test split as prepared in prior work [3, 49], and Gaussian noise with zero mean and 0.05 standard deviation. ScanNet-v2 contains 1513 scans from real-world scenes that cover a wide range of room types. Since the provided meshes in ScanNet-v2 are not watertight, models are trained on the Synthetic Rooms dataset and tested on ScanNet-v2, which also enables some assessment of Sim2Real performance of various methods.

Evaluation Metrics: The quantitative evaluation metrics used in data tables are standard metrics that enable us to form comparisons to prior work. These include: volumetric IoU, Chamfer- L_1 distance $\times 10^2$, and normal consistency (NC). A detailed definition of each metric can be found in

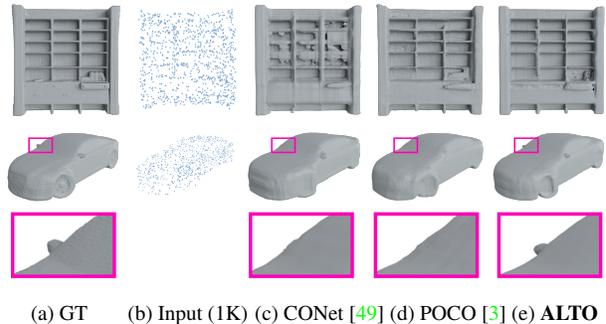


Figure 6. **Object-level comparisons on ShapeNet.** On the car, ALTO recovers the detail of having both side mirrors.

the Occupancy Networks paper [41]. We also include an F-Score metric [64] with threshold value 1%.

Baselines for Comparison: As noted by Boulch et al. [3], baseline methods often perform better in the settings of the original papers. In the same spirit, we thus strictly adapt the protocol of the state-of-the-art (SOTA) paper POCO [3] for evaluation protocol. In addition to POCO, other baselines we include are SPSR [33], ONet [41] and ConvONet [49]. Note that we omit NFK [68] from our evaluations, as they have not made their code publicly available.

Our Method: “Our method” is ALTO. ALTO combines Sec. 3.2 + Sec. 3.3, and is shown in Fig. 5d. Figures/tables use ALTO to denote the proposed method. If we are considering an ablation analysis we will use ALTO with parentheses, e.g., “ALTO (Encoder Only)”, and the table caption will specify the ablation. To demonstrate ALTO with different latent topologies, for object-level reconstruction, we use alternation between point and feature triplanes and the resolution of initial individual plane $H = W = 64$. For scene-level reconstruction, we use alternation between point and feature volumes and the resolution $H = W = D = 64$.

4.2. Results of Object-level Reconstruction

Qualitative Object-level Comparisons: Qualitative results of object-level reconstruction are provided in Fig. 6. We observe that [49] obtains a blurry reconstruction. Note

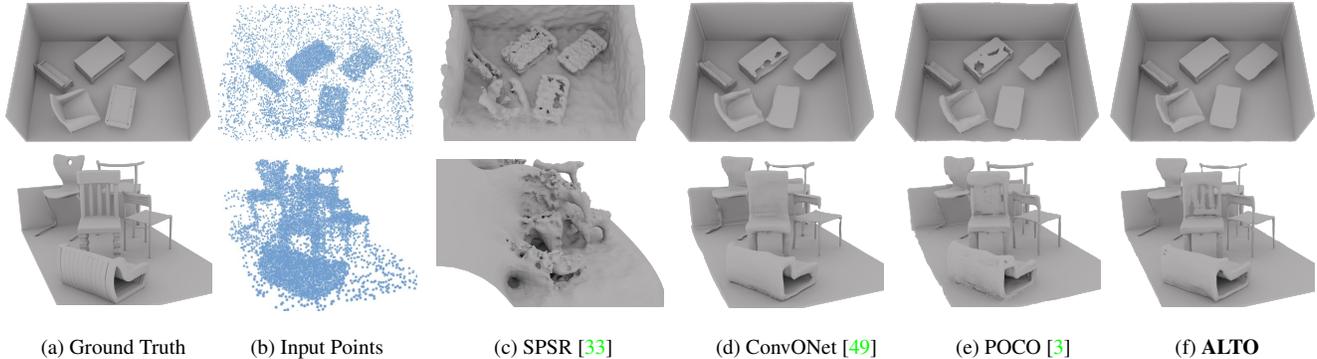


Figure 7. **Qualitative comparison on scene-level reconstruction Synthetic Room Dataset.** Learning-based methods are trained and tested on 10K noisy points. ALTO can reconstruct the (top scene) double-deck table and (bottom scene) details in the chair.

Method	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow
ONet [41]	0.475	2.03	0.783	0.541
SPSR [33]	-	2.23	0.866	0.810
SPSR trimmed [33]	-	0.69	0.890	0.892
ConvONet [49]	0.849	0.42	0.915	0.964
DP-ConvONet [37]	0.800	0.42	0.912	0.960
POCO [3]	0.884	0.36	0.919	0.980
ALTO	0.914	0.35	0.921	0.981

Table 2. **Synthetic Room Dataset.** Input points 10K with noise added. Boldface font represents the preferred results.

Method	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow
ConvONet [49]	0.818	0.46	0.906	0.943
POCO [3]	0.801	0.57	0.904	0.812
ALTO	0.882	0.39	0.911	0.969

Table 3. **Performance on Synthetic Room Dataset (sparser input point cloud with 3K input points).** Boldface font represents the preferred results.

that the width of the bookcase dividers are thicker and there are spurious blobs on the shelves. The SOTA baseline [3] is able to recover some detail, such as the wheel geometry in the car, but loses both side mirrors in the reconstruction. ALTO seems to have a higher fidelity reconstruction due to combining ideas from point and grid latents.

Quantitative Object-level Comparisons: ALTO’s performance metrics at various point density levels are listed in Tab. 1, for 3K, 1K and 300 input points. When point clouds are sparser, ALTO performs better than POCO on all four metrics. At high point density, ALTO outperforms POCO on three of four metrics. An ablation of just using the ALTO latent encoding and a traditional interpolating decoder [49] is also conducted in the table.

4.3. Scene-level Reconstruction

Qualitative Scene-level Comparisons: ALTO achieves detailed qualitative results compared to baselines. Fig. 7 depicts scene-level reconstruction on the Synthetic Room dataset introduced in [49]. In the first row of Fig. 7 the baselines of ConvONet and POCO both have holes in the coffee table. In the second row of Fig. 7 the high-frequency detail in the wooden slats of the chair is fully blurred out by ConvONet. The advantages of ALTO are even more apparent for fine detail, such as the thin lamp-posts shown in Fig. 1. ALTO reduces the quantization effect due to the grid discretization in the grid form by using iterative alternation between grid and point latents, encoding more fine-grained local features for conditional occupancy field generation.

Quantitative Scene-level Comparisons: ALTO scores higher on quantitative values for scene-level metrics, shown in Tab. 2 and Tab. 3. In the sparse setting, for the baselines methods, we find that ConvONet [49] is quantitatively superior to the SOTA of POCO [3] because oversmoothing tends to improve quantitative results on noisy point clouds. Nonetheless, ALTO performs better than both baselines because ALTO limits spurious noise without resorting to as much oversmoothing.

4.4. Real-world Scene Generalization

A final experiment in the main paper is to assess the performance of our model in real-world scans from ScanNet-v2 [11]. All models are trained on Synthetic Rooms and tested on ScanNet-v2 to demonstrate generalization capability of our method along with baselines. We demonstrate the qualitative results of the setting where models trained on both the same input points of synthetic dataset as ScanNet test set ($N_{\text{Train}}=N_{\text{Test}}=3\text{K}$) in Fig. 8. Our method is qualitatively superior to SPSR [33], ConvONet [49], and POCO [3]. As in the Synthetic Rooms dataset, we observe that ConvONet oversmooths surfaces (sometimes causing entire objects to disappear, like the conference table in the

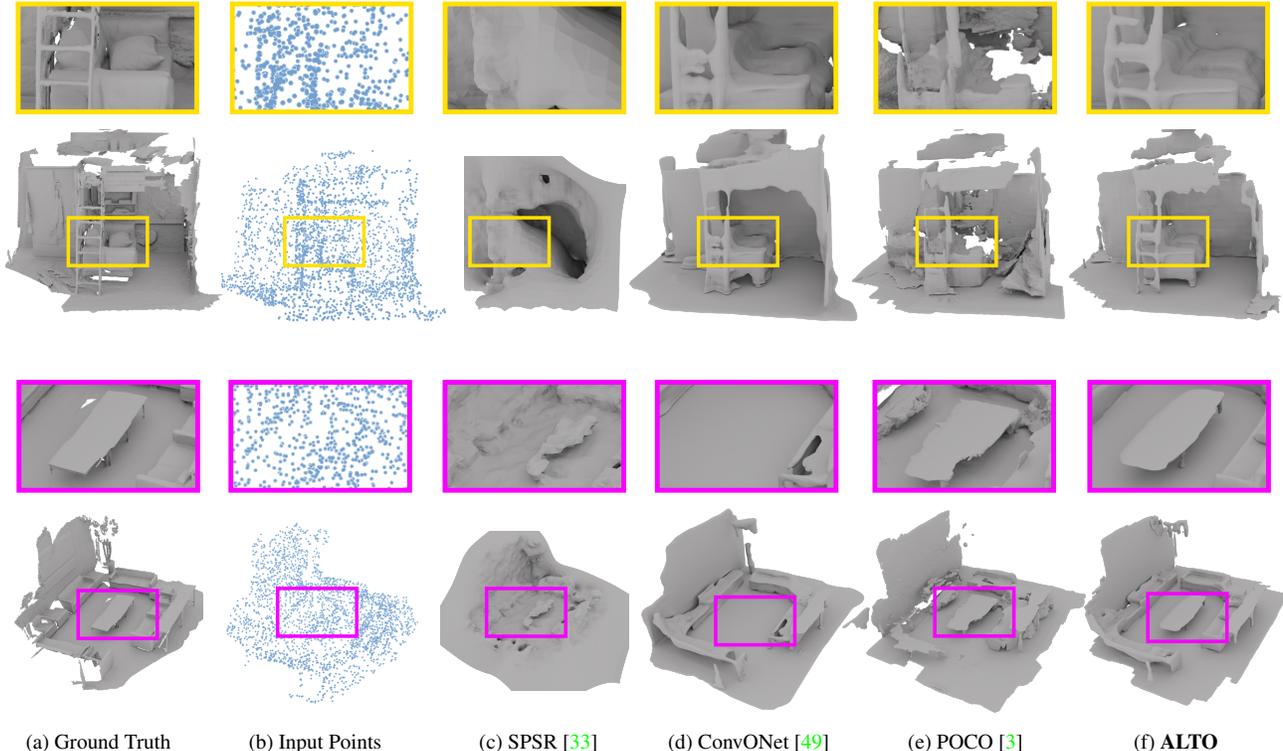


Figure 8. **Cross-dataset evaluation of ALTO and baselines by training on Synthetic Rooms [49] and testing on real-world ScanNet-v2 [11].** Note the large conference-room table is missing in ConvONet [49] (purple inset). The ladder (yellow inset) is a high-frequency surface and we believe our method is qualitatively closest. Please zoom in if browsing with PDF.

Method	$N_{\text{Train}}=10\text{K}, N_{\text{Test}}=3\text{K}$		$N_{\text{Train}}=N_{\text{Test}}=3\text{K}$	
	Chamfer- L_1 ↓	F-score ↑	Chamfer- L_1 ↓	F-score ↑
ConvONet [49]	1.01	0.719	1.16	0.669
POCO [3]	0.93	0.737	1.15	0.667
ALTO	0.87	0.746	0.92	0.726

Table 4. **ScanNet-v2.** We test the generalization capability of all the methods on real-world scans ScanNet using models trained on both the same input points of synthetic dataset as test set ($N_{\text{Train}}=N_{\text{Test}}=3\text{K}$) and different point density level ($N_{\text{Train}}=10\text{K}, N_{\text{Test}}=3\text{K}$). Boldface font represents the preferred results.

purple inset of Fig. 8). In contrast, POCO retains some detail but is noisier. The quantitative results in Tab. 4 are consistent with qualitative results. The cross point density test results ($N_{\text{Train}}=10\text{k}, N_{\text{Test}}=3\text{k}$) also demonstrates the superiority of our method on generalization when there are abundant input points in synthetic dataset used for training and low point density in real-world inference.

5. Discussion and Conclusion

In summary, this paper has adopted a different philosophy from the SOTA in surface detail recovery. We rely nei-

Method	# Parameters	Inference time (s)
ConvONet [49]	4,166,657	1.6
POCO [3]	12,790,454	36.1
ALTO	4,787,905	3.6

Table 5. **Runtime comparison.** We report the number of parameters and inference time corresponding to Fig. 8. ALTO is much faster than POCO and recovers more detail [3]. ALTO is also only slightly slower than fast methods that are not as spatially expressive [49].

ther on point latents [3] or grid latents [49] alone, but alternate between topologies. The output of ALTO is a spatially expressive latent that is also topologically easy-to-decode into a 3D surface. This breaks a Pareto tradeoff that previous works have posited between spatial expressiveness and decoding complexity. For this reason, it is not surprising that our method reconstructs more detailed 3D surfaces with faster runtimes than state-of-the-art (Fig. 1 and Tab. 5).

The idea of alternating latent topologies could have implications beyond surface reconstruction. Concurrent research has introduced unusual latent topologies, known as **irregular latents** that show compelling performance ben-

efits for neural fields [79]. One can imagine alternating not only between point, triplane, and voxel latents, but also throwing irregular latents in the mix. We are also curious to see if alternating topologies can improve performance on a wide range of tasks in neural fields that require spatially expressive latents, such as semantic or affordance fields.

6. Acknowledgement

This project was supported by the US DoD LUCI (Laboratory University Collaboration Initiative) fellowship and partially supported by ARL grants W911NF-20-2-0158 and W911NF-21-2-0104 under the cooperative A2I2 program. D.P. is supported by the Swiss National Science Foundation under grant number P500PT_206946. G.W. is supported by Samsung, Stanford HAI, and a PECASE from the ARO. L.G. is also supported by an ONR Vannevar Bush Faculty Fellowship. A.K. is also supported by a DARPA Young Faculty Award, NSF CAREER Award, and Army Young Investigator Award.

References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3d point clouds. In *Proc. of the International Conf. on Machine learning (ICML)*, 2018. 2
- [2] Matan Atzmon and Yaron Lipman. SAL: sign agnostic learning of shapes from raw data. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, pages 2562–2571, 2020. 3
- [3] Alexandre Boulch and Renaud Marlet. Poco: Point convolution for surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6302–6314, 2022. 1, 2, 3, 5, 6, 7, 8
- [4] André Brock, Theodore Lim, James M. Ritchie, and Nick Weston. Generative and discriminative voxel modeling with convolutional neural networks. *arXiv.org*, 1608.04236, 2016. 2
- [5] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015. 6
- [6] Wenzheng Chen, Huan Ling, Jun Gao, Edward Smith, Jaako Lehtinen, Alec Jacobson, and Sanja Fidler. Learning to predict 3d objects with an interpolation-based differentiable renderer. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 2
- [7] Zhiqin Chen and Hao Zhang. Learning implicit fields for generative shape modeling. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3
- [8] Julian Chibane, Thiemo Alldieck, and Gerard Pons-Moll. Implicit functions in feature space for 3d shape reconstruction and completion. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2, 3
- [9] Christopher Bongsso Choy, Danfei Xu, JunYoung Gwak, Kevin Chen, and Silvio Savarese. 3d-r2n2: A unified approach for single and multi-view 3d object reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2016. 2
- [10] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, pages 424–432. Springer, 2016. 4
- [11] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 6, 7, 8
- [12] Angela Dai, Christian Diller, and Matthias Nießner. Sg-nn: Sparse generative neural networks for self-supervised scene completion of rgb-d scans. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 849–858, 2020. 2
- [13] Boyang Deng, Kyle Genova, Soroosh Yazdani, Sofien Bouaziz, Geoffrey Hinton, and Andrea Tagliasacchi. Cvxnets: Learnable convex decomposition. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [14] Zhantao Deng, Jan Bednařik, Mathieu Salzmann, and Pascal Fua. Better patch stitching for parametric surface reconstruction. 2020. 2
- [15] Simon Donne and Andreas Geiger. Learning non-volumetric depth fusion using successive reprojections. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [16] Philipp Erler, Paul Guerrero, Stefan Ohrhallinger, Niloy J Mitra, and Michael Wimmer. Points2surf learning implicit surfaces from point clouds. In *European Conference on Computer Vision*, pages 108–124. Springer, 2020. 2
- [17] Philipp Erler, Paul Guerrero, Stefan Ohrhallinger, Niloy J. Mitra, and Michael Wimmer. Points2surf learning implicit surfaces from point clouds. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 3
- [18] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [19] Matheus Gadelha, Subhransu Maji, and Rui Wang. 3d shape induction from 2d views of multiple objects. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2017. 2
- [20] Jun Gao, Wenzheng Chen, Tommy Xiang, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Learning deformable tetrahedral meshes for 3d reconstruction. In *Advances in Neural Information Processing Systems (NIPS)*, 2020. 2
- [21] Simon Giebenhain and Bastian Goldlücke. Air-nets: An attention-based framework for locally conditioned implicit representations. In *2021 International Conference on 3D Vision (3DV)*, pages 1054–1064. IEEE, 2021. 3

- [22] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh R-CNN. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [23] Amos Gropp, Lior Yariv, Niv Haim, Matan Atzmon, and Yaron Lipman. Implicit geometric regularization for learning shapes. In *Proc. of the International Conf. on Machine learning (ICML)*, 2020. 3
- [24] Thibault Groueix, Matthew Fisher, Vladimir G. Kim, Bryan C. Russell, and Mathieu Aubry. AtlasNet: A papier-mâché approach to learning 3d surface generation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [25] Christian Häne, Shubham Tulsiani, and Jitendra Malik. Hierarchical surface prediction for 3d object reconstruction. *arXiv.org*, 1704.00710, 2017. 2
- [26] Wilfried Hartmann, Silvano Galliani, Michal Havlena, Luc Van Gool, and Konrad Schindler. Learned multi-patch similarity. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2
- [27] Xinyu Huang, Xinjing Cheng, Qichuan Geng, Binbin Cao, Dingfu Zhou, Peng Wang, Yuanqing Lin, and Ruigang Yang. The apolloescape dataset for autonomous driving. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [28] Chiyu Jiang, Avneesh Sud, Ameesh Makadia, Jingwei Huang, Matthias Nießner, and Thomas Funkhouser. Local implicit grid representations for 3d scenes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [29] Li Jiang, Shaoshuai Shi, Xiaojuan Qi, and Jiaya Jia. GAL: geometric adversarial loss for single-view 3d-object reconstruction. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [30] Zhenyu Jiang, Yifeng Zhu, Maxwell Svetlik, Kuan Fang, and Yuke Zhu. Synergies between affordance and geometry: 6-dof grasp detection via implicit representations. *arXiv preprint arXiv:2104.01542*, 2021. 2
- [31] Angjoo Kanazawa, Shubham Tulsiani, Alexei A. Efros, and Jitendra Malik. Learning category-specific mesh reconstruction from image collections. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [32] Abhishek Kar, Christian Häne, and Jitendra Malik. Learning a multi-view stereo machine. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2
- [33] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 6, 7, 8
- [34] Michael M. Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the Fourth Eurographics Symposium on Geometry Processing, Cagliari, Sardinia, Italy, June 26-28, 2006*, volume 256, pages 61–70, 2006. 2
- [35] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014. 6
- [36] Yiyi Liao, Simon Donne, and Andreas Geiger. Deep marching cubes: Learning explicit surface representations. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [37] Stefan Lionar, Daniil Emtsev, Dusan Svilarkovic, and Songyou Peng. Dynamic plane convolutional occupancy networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1829–1838, 2021. 2, 3, 5, 7
- [38] 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. 6
- [39] Qianli Ma, Shunsuke Saito, Jinlong Yang, Siyu Tang, and Michael J. Black. SCALE: Modeling clothed humans with a surface codec of articulated local elements. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 2
- [40] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *Proc. IEEE International Conf. on Intelligent Robots and Systems (IROS)*, 2015. 2
- [41] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. 2, 3, 6, 7
- [42] Mateusz Michalkiewicz, Jhony K Pontes, Dominic Jack, Mahsa Baktashmotlagh, and Anders Eriksson. Implicit surface representations as layers in neural networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 3
- [43] Ben Mildenhall, Pratul P Srinivasan, Matthew Tancik, Jonathan T Barron, Ravi Ramamoorthi, and Ren Ng. NeRF: Representing scenes as neural radiance fields for view synthesis. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2020. 3
- [44] Michael Niemeyer, Lars Mescheder, Michael Oechsle, and Andreas Geiger. Differentiable volumetric rendering: Learning implicit 3d representations without 3d supervision. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [45] Junyi Pan, Xiaoguang Han, Weikai Chen, Jiapeng Tang, and Kui Jia. Deep mesh reconstruction from single RGB images via topology modification networks. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [46] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. DeepSDF: Learning continuous signed distance functions for shape representation. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2, 3
- [47] Despoina Paschalidou, Ali Osman Ulusoy, Carolin Schmitt, Luc van Gool, and Andreas Geiger. Raynet: Learning volumetric 3d reconstruction with ray potentials. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [48] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, et al. Pytorch: An imperative style, high-performance deep learning library. *Ad-*

- vances in neural information processing systems*, 32, 2019. 6
- [49] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision*, pages 523–540. Springer, 2020. 1, 2, 3, 5, 6, 7, 8
- [50] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. 3
- [51] Charles R Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in Neural Information Processing Systems (NIPS)*, 2017. 2
- [52] Danilo Jimenez Rezende, S. M. Ali Eslami, Shakir Mohamed, Peter Battaglia, Max Jaderberg, and Nicolas Heess. Unsupervised learning of 3d structure from images. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 2
- [53] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. OctNetFusion: Learning depth fusion from data. In *Proc. of the International Conf. on 3D Vision (3DV)*, 2017. 2
- [54] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2017. 2
- [55] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. 4
- [56] Shunsuke Saito, Zeng Huang, Ryota Natsume, Shigeo Morishima, Angjoo Kanazawa, and Hao Li. Pifu: Pixel-aligned implicit function for high-resolution clothed human digitization. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 3
- [57] Shunsuke Saito, Tomas Simon, Jason M. Saragih, and Hanbyul Joo. Pifuhd: Multi-level pixel-aligned implicit function for high-resolution 3d human digitization. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 3
- [58] Vincent Sitzmann, Julien N.P. Martel, Alexander W. Bergman, David B. Lindell, and Gordon Wetzstein. Implicit neural representations with periodic activation functions. In *Proc. NeurIPS*, 2020.
- [59] Vincent Sitzmann, Michael Zollhöfer, and Gordon Wetzstein. Scene representation networks: Continuous 3d-structure-aware neural scene representations. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 3
- [60] David Stutz and Andreas Geiger. Learning 3d shape completion from laser scan data with weak supervision. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2018. 2
- [61] Towaki Takikawa, Joey Litalien, Kangxue Yin, Karsten Kreis, Charles T. Loop, Derek Nowrouzezahrai, Alec Jacobson, Morgan McGuire, and Sanja Fidler. Neural geometric level of detail: Real-time rendering with implicit 3d shapes. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2021. 3
- [62] Jiapeng Tang, Jiabao Lei, Dan Xu, Feiying Ma, Kui Jia, and Lei Zhang. Sa-convnet: Sign-agnostic optimization of convolutional occupancy networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 6504–6513, 2021. 2, 3
- [63] M. Tatarchenko, A. Dosovitskiy, and T. Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2017. 2
- [64] Maxim Tatarchenko, Stephan R Richter, René Ranftl, Zhuwen Li, Vladlen Koltun, and Thomas Brox. What do single-view 3d reconstruction networks learn? In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 3405–3414, 2019. 6
- [65] Hugues Thomas, Charles R. Qi, Jean-Emmanuel Deschaud, Beatriz Marcotegui, François Goulette, and Leonidas J. Guibas. Kpconv: Flexible and deformable convolution for point clouds. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [66] Suhani Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi SM Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *arXiv preprint arXiv:2111.13260*, 2021. 2
- [67] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [68] Francis Williams, Zan Gojcic, Sameh Khamis, Denis Zorin, Joan Bruna, Sanja Fidler, and Or Litany. Neural fields as learnable kernels for 3d reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 18500–18510, 2022. 3, 6
- [69] Jiajun Wu, Ilker Yildirim, Joseph J. Lim, Bill Freeman, and Joshua B. Tenenbaum. Galileo: Perceiving physical object properties by integrating a physics engine with deep learning. In *Advances in Neural Information Processing Systems (NIPS)*, 2015. 2
- [70] Jiajun Wu, Chengkai Zhang, Tianfan Xue, Bill Freeman, and Josh Tenenbaum. Learning a probabilistic latent space of object shapes via 3d generative-adversarial modeling. In *Advances in Neural Information Processing Systems (NIPS)*, 2016. 2
- [71] Haozhe Xie, Hongxun Yao, Xiaoshuai Sun, Shangchen Zhou, and Shengping Zhang. Pix2vox: Context-aware 3d reconstruction from single and multi-view images. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [72] Qiangeng Xu, Weiyue Wang, Duygu Ceylan, Radomír Mech, and Ulrich Neumann. DISN: deep implicit surface network for high-quality single-view 3d reconstruction. In *Advances in Neural Information Processing Systems (NIPS)*, 2019. 3
- [73] Guandao Yang, Xun Huang, Zekun Hao, Ming-Yu Liu, Serge J. Belongie, and Bharath Hariharan. Pointflow: 3d

- point cloud generation with continuous normalizing flows. In *Proc. of the IEEE International Conf. on Computer Vision (ICCV)*, 2019. 2
- [74] Yao Yao, Zixin Luo, Shiwei Li, Tian Fang, and Long Quan. Mvsnet: Depth inference for unstructured multi-view stereo. *Proc. of the European Conf. on Computer Vision (ECCV)*, 2018. 2
- [75] Yao Yao, Zixin Luo, Shiwei Li, Tianwei Shen, Tian Fang, and Long Quan. Recurrent mvsnet for high-resolution multi-view stereo depth inference. *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. 2
- [76] Lior Yariv, Matan Atzmon, and Yaron Lipman. Universal differentiable renderer for implicit neural representations. *arXiv.org*, 2003.09852, 2020. 3
- [77] Lior Yariv, Yoni Kasten, Dror Moran, Meirav Galun, Matan Atzmon, Ronen Basri, and Yaron Lipman. Multiview neural surface reconstruction by disentangling geometry and appearance. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2020. 3
- [78] Zehao Yu and Shenghua Gao. Fast-mvsnet: Sparse-to-dense multi-view stereo with learned propagation and gauss-newton refinement. In *Proc. IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2020. 2
- [79] Biao Zhang, Matthias Nießner, and Peter Wonka. 3dilg: Irregular latent grids for 3d generative modeling. *arXiv preprint arXiv:2205.13914*, 2022. 9
- [80] Hengshuang Zhao, Li Jiang, Jiaya Jia, Philip HS Torr, and Vladlen Koltun. Point transformer. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 16259–16268, 2021. 5

ALTO: Alternating Latent Topologies for Implicit 3D Reconstruction

Supplemental Material

Zhen Wang^{1*} Shijie Zhou^{1*} Jeong Joon Park² Despoina Paschalidou²
 Suya You³ Gordon Wetzstein² Leonidas Guibas² Achuta Kadambi¹

¹University of California, Los Angeles ²Stanford University ³DEVCOM Army Research Laboratory

Supplementary Content

This supplement is organized as follows:

- Section **A** contains network architecture details;
- Section **B** contains more details on the training and inference settings;
- Section **C** contains more ablation studies of our method;
- Section **D** contains both quantitative and qualitative results on ShapeNet dataset;
- Section **E** contains more qualitative results on Synthetic Room dataset;
- Section **F** contains additional qualitative results on ScanNet dataset;
- Section **G** contains the code link of the comparison baselines; and
- Section **H** contains discussion on the limitation of our method and future work.

A. Network Architecture

PointNet: Given the input un-oriented point cloud $\mathcal{P} = \{\mathbf{p}_i \in \mathbb{R}^3\}_{i=1}^S$, where S is the number of input points, we map the input coordinates to point features using a fully-connected layer and a ResNet-FC [3] block. Instead of using global features as in [10], we use locally-pooled features to fuse local features. Specifically, we aggregate features within the same plane or voxel cell from a 2D triplanar or 3D volumetric grids using max-pooling. We concatenate the locally pooled features with the feature before pooling and then input to the next ResNet block. To obtain the final point features, there are totally 5 ResNet blocks used.

Our ALTO U-Net: Our alternation U-Net architecture is similar to traditional U-Net [2, 11], except that we replace the convolution-only block with our ALTO block where point and grid (either 2D or 3D) features are converted back and forth as depicted in Fig. 3 of main paper. The input and output feature dimensions is set to be 32. There is no ALTO block in the final block of the U-Net.

Our Attention-based Decoder: For the triplane representation, we implement 3 single-head attention for 3 feature planes respectively, where the hidden dimension is equal to the feature dimension 32. For the volume representation, we implement a multi-head attention with h heads. To maximize the flexibility of our method for different datasets and experiments, we set the number of heads h as a hyperparameter and the hidden dimension as $h \times \text{feature dimension}(32)$. The following occupancy network consisting of 5 stacked ResNet-FC blocks with skip connections is used to predict the occupancy probability of query point features. For all experiments, we use a hidden dimension equal to the attention output feature dimension and 5 ResNet blocks for the occupancy network.

*Equal contribution.

Total # of alternation blocks	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow
0	0.831	0.55	0.912	0.892
3	0.847	0.50	0.914	0.910
6	0.863	0.47	0.922	0.924

Table A. Ablation study of total number of ALTO alternation blocks on ShapeNet dataset with 300 input points.

Method	IoU \uparrow	Chamfer- L_1 \downarrow	NC \uparrow	F-score \uparrow
ConvONet (3×128^2) [9]	0.805	0.44	0.903	0.948
ConvONet (64^3) [9]	0.849	0.42	0.915	0.964
ALTO (3×128^2 , Encoder Only)	0.834	0.43	0.906	0.960
ALTO (3×128^2)	0.895	0.37	0.910	0.974
ALTO (64^3 , Encoder Only)	0.903	0.36	0.920	0.978
ALTO (64^3)	0.914	0.35	0.921	0.981

Table B. Ablation study of our attention-based decoder for different latent topologies used (i.e. point-triplane and point-voxel alternations) on Synthetic Room dataset. Input points 10K with noise added. Boldface font represents the preferred results.

B. Training and Inference Details

Object-Level Reconstruction: For object-level reconstruction in ShapeNet, we use alternation between latent topologies: point and triplane, because triplane representation is found to tend to give better results for object-level reconstruction in ConvONet [9]. The dimension of each 2D feature plane is set as 64^2 . The depth of our ALTO U-Net is 4, and we do not downsample or upsample in the top two levels of the U-Net, so the lowest resolution of the U-Net is 16^2 .

Scene-Level Reconstruction: For scene-level reconstruction, we use alternation between two topologies: point and feature volume. The dimension of the feature volume is set as 64^3 . The depth of our ALTO U-Net is 4, and similarly we do not downsample or upsample in the top two levels, so the lowest resolution of the U-Net is 16^3 . At decoder stage, we set the hyperparameter $h = 4$ for experiments on Synthetic Room dataset and $h = 1$ for experiments on ScanNet dataset which we find the best performance in practice.

Mesh Generation: We use a form of Marching Cubes (MC) [7] to evaluate occupancy values from implicit representations on a 3D grid. As a result of Marching Cube, the vertices are usually placed in the middle of segments, which causes discretization effects [1]. To deal with this issue, we apply the refinement method from POCO [1], which takes both the generated vertices and their floor to predict their occupancy values again. After that, we compare two values, mask out non-perfect vertices, take the average between the generated vertices and their floor, and repeat 10 times to improve the granularity. For object-level reconstruction, we use resolution 128 and for scene-level reconstruction, we use resolution 256 for marching cubes.

Hardware: We describe the detailed setups that have been used for inference evaluation:

- CUDA version: 11.1
- PyTorch version: 1.9.0
- GPU: single NVIDIA GeForce RTX 3090
- CPU: AMD RYZEN PRO 3955WX 16-Cores CPU

C. Ablation Studies

In Tab. A, we report the performance of method with different number of alternations between point and grid forms within each block in the ALTO U-Net. 0 represents no point-grid alternations (i.e. staying with only grid form), 3 represents that there is only point-grid alternation in the top two levels of our ALTO U-Net, and 6 represents that there is point-grid

Method	IoU \uparrow				Chamfer- L_1 \downarrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.734	0.849	0.902	0.908	0.64	0.34	0.23	0.22
Bench	0.682	0.830	0.865	0.890	0.67	0.35	0.28	0.26
Cabinet	0.855	0.940	0.960	0.965	0.82	0.46	0.37	0.34
Car	0.830	0.886	0.921	0.924	1.04	0.75	0.41	0.43
Chair	0.720	0.871	0.919	0.925	0.95	0.46	0.33	0.32
Display	0.799	0.927	0.956	0.962	0.82	0.36	0.28	0.27
Lamp	0.546	0.785	0.877	0.868	1.59	0.59	0.33	0.34
Loudspeaker	0.826	0.918	0.957	0.953	1.18	0.64	0.41	0.41
Rifle	0.668	0.846	0.897	0.898	0.66	0.28	0.19	0.19
Sofa	0.865	0.936	0.963	0.966	0.73	0.42	0.30	0.29
Table	0.739	0.888	0.924	0.937	0.76	0.38	0.31	0.29
Telephone	0.896	0.955	0.968	0.977	0.46	0.27	0.22	0.21
Vessel	0.729	0.865	0.927	0.924	0.94	0.43	0.25	0.26
mean	0.761	0.884	0.926	0.931	0.87	0.44	0.30	0.30

Method	NC \uparrow				F-score \uparrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.886	0.931	0.944	0.949	0.829	0.965	0.994	0.992
Bench	0.871	0.921	0.928	0.941	0.827	0.964	0.988	0.991
Cabinet	0.913	0.956	0.961	0.967	0.833	0.956	0.979	0.982
Car	0.874	0.893	0.894	0.917	0.747	0.849	0.946	0.940
Chair	0.886	0.943	0.956	0.959	0.730	0.939	0.985	0.985
Display	0.926	0.968	0.975	0.976	0.795	0.971	0.994	0.993
Lamp	0.809	0.900	0.929	0.924	0.581	0.892	0.975	0.962
Loudspeaker	0.903	0.939	0.952	0.951	0.727	0.892	0.964	0.955
Rifle	0.849	0.929	0.949	0.949	0.818	0.980	0.998	0.996
Sofa	0.928	0.958	0.967	0.971	0.832	0.953	0.989	0.987
Table	0.917	0.959	0.966	0.968	0.824	0.967	0.991	0.990
Telephone	0.970	0.983	0.985	0.987	0.930	0.989	0.998	0.998
Vessel	0.857	0.919	0.940	0.940	0.734	0.931	0.989	0.982
mean	0.891	0.938	0.950	0.954	0.785	0.942	0.984	0.981

Table C. Performance on ShapeNet with input noisy point cloud 3K. Boldface font represents the preferred results.

alternations in each level of our ALTO U-Net. As we can see the results, we can observe the trend that increasing the number of ALTO blocks improves the results for all the metrics.

We also report the results of the ablation study of our attention-based decoder on synthetic room dataset in Tab. B. As demonstrated in the table, with our attention-based decoder, it improves results for both triplanar (3×128^3) and volumetric representations (64^3).

D. Additional Results on ShapeNet

D.a. Quantitative results

We show per-category quantitative results in ShapeNet with various point density levels: 3K input points (Tab. C), 1K input points (Tab. D) and 300 input points (Tab. E). It is notable that when point clouds get sparser, ALTO performs better than POCO on all four metrics for all categories.

D.b. Qualitative results

Besides 1K input points for ShapeNet as we show in Fig. 6 of the main paper, we show additional qualitative results in ShapeNet with 3K input points in Fig. A and 300 input points in Fig. B.

E. Additional Results on Synthetic Room Dataset

We show additional qualitative results in Synthetic Room dataset with 10K inputs points in Fig. C and 3K inputs points in Fig. D.

Method	IoU \uparrow				Chamfer- L_1 \downarrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.748	0.825	0.850	0.872	0.59	0.39	0.32	0.29
Bench	0.702	0.798	0.804	0.856	0.62	0.40	0.38	0.30
Cabinet	0.862	0.926	0.936	0.953	0.76	0.50	0.46	0.37
Car	0.837	0.867	0.878	0.901	0.99	0.83	0.60	0.50
Chair	0.736	0.837	0.867	0.894	0.89	0.55	0.44	0.39
Display	0.812	0.911	0.930	0.946	0.78	0.41	0.34	0.31
Lamp	0.567	0.741	0.807	0.820	1.44	0.68	0.50	0.50
Loudspeaker	0.831	0.899	0.923	0.933	1.14	0.72	0.54	0.48
Rifle	0.680	0.801	0.850	0.862	0.63	0.36	0.27	0.25
Sofa	0.873	0.921	0.937	0.952	0.69	0.47	0.38	0.33
Table	0.757	0.858	0.880	0.913	0.70	0.44	0.38	0.33
Telephone	0.897	0.946	0.953	0.968	0.46	0.29	0.26	0.23
Vessel	0.736	0.840	0.880	0.893	0.91	0.51	0.37	0.33
mean	0.772	0.859	0.884	0.905	0.82	0.50	0.40	0.35

Method	NC \uparrow				F-score \uparrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.894	0.922	0.920	0.933	0.850	0.946	0.970	0.976
Bench	0.882	0.911	0.902	0.925	0.849	0.943	0.956	0.979
Cabinet	0.925	0.949	0.945	0.957	0.852	0.939	0.951	0.972
Car	0.904	0.885	0.867	0.889	0.763	0.819	0.868	0.912
Chair	0.893	0.931	0.930	0.946	0.753	0.902	0.943	0.965
Display	0.930	0.961	0.962	0.970	0.805	0.956	0.976	0.984
Lamp	0.820	0.885	0.895	0.905	0.606	0.845	0.924	0.926
Loudspeaker	0.914	0.929	0.928	0.936	0.740	0.863	0.908	0.926
Rifle	0.859	0.916	0.928	0.936	0.828	0.957	0.984	0.987
Sofa	0.937	0.950	0.950	0.960	0.846	0.932	0.961	0.974
Table	0.918	0.950	0.949	0.961	0.842	0.947	0.964	0.979
Telephone	0.972	0.980	0.979	0.984	0.940	0.983	0.990	0.994
Vessel	0.866	0.906	0.913	0.923	0.740	0.899	0.952	0.961
mean	0.901	0.929	0.928	0.940	0.801	0.918	0.950	0.964

Table D. Performance on ShapeNet with input noisy point cloud 1K. Boldface font represents the preferred results.

F. Additional Results on ScanNet

We demonstrate the Sim2Real qualitative results with the model trained on Synthetic Room dataset and tested on ScanNet in Fig. 8 of the main paper. We show in Fig. E of the supplement material the Sim2Real results with different point density levels (i.e. $N_{\text{Train}}=10\text{k}$, $N_{\text{Test}}=3\text{k}$) to further demonstrate the generalization capability of our method ALTO.

G. Comparison Code Links

We list all the links of the code of the comparisons baselines in Tab. F. Our code is attached as part of the supplement materials and will be uploaded at <https://github.com/cvpr2023-submission/ALTO> upon acceptance.

H. Limitation and Future Work

For our current method, we are not learning a probabilistic generative model that can learn the distribution of the input data, which limits the diversity of the shapes our model can generate. Moreover, we are uniformly sampling points as in previous work such as [9]. More efficient sampling strategy that samples more points on densely populated regions and less on sparsely populated regions can be adopted to capture more details on the fine-grained areas.

As our method is general in encoding 3D point features, it can be generalized to not just occupancy fields, but also radiance fields trained from images. Similarly, it can be applied to a broader range of neural fields such as semantic field [12] and affordance field [4].

Method	IoU \uparrow				Chamfer- L_1 \downarrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.760	0.782	0.744	0.825	0.57	0.48	0.57	0.39
Bench	0.716	0.743	0.707	0.801	0.60	0.50	0.56	0.39
Cabinet	0.867	0.900	0.889	0.927	0.73	0.52	0.58	0.46
Car	0.834	0.843	0.817	0.867	0.99	0.76	0.83	0.67
Chair	0.736	0.787	0.776	0.840	0.89	0.67	0.71	0.52
Display	0.817	0.885	0.878	0.917	0.76	0.47	0.49	0.38
Lamp	0.567	0.663	0.681	0.747	1.38	1.02	0.93	0.76
Loudspeaker	0.827	0.870	0.867	0.901	1.16	0.78	0.79	0.64
Rifle	0.691	0.757	0.742	0.801	0.61	0.43	0.45	0.35
Sofa	0.872	0.898	0.893	0.926	0.69	0.52	0.53	0.42
Table	0.758	0.813	0.794	0.868	0.72	0.52	0.57	0.42
Telephone	0.916	0.939	0.927	0.952	0.41	0.31	0.33	0.27
Vessel	0.748	0.797	0.795	0.846	0.85	0.63	0.60	0.47
mean	0.778	0.821	0.808	0.863	0.80	0.59	0.61	0.47

Method	NC \uparrow				F-score \uparrow			
	ONet [8]	ConvONet [9]	POCO [1]	ALTO	ONet [8]	ConvONet [9]	POCO [1]	ALTO
Airplane	0.897	0.901	0.867	0.914	0.864	0.902	0.867	0.938
Bench	0.878	0.886	0.864	0.906	0.860	0.912	0.882	0.947
Cabinet	0.916	0.931	0.917	0.943	0.856	0.916	0.896	0.943
Car	0.875	0.864	0.835	0.873	0.757	0.810	0.766	0.850
Chair	0.889	0.905	0.885	0.923	0.754	0.850	0.833	0.910
Display	0.926	0.947	0.938	0.956	0.813	0.926	0.916	0.957
Lamp	0.813	0.853	0.834	0.875	0.618	0.771	0.781	0.857
Loudspeaker	0.897	0.911	0.897	0.916	0.737	0.832	0.819	0.871
Rifle	0.863	0.890	0.883	0.909	0.838	0.919	0.918	0.952
Sofa	0.928	0.935	0.924	0.946	0.846	0.906	0.899	0.941
Table	0.917	0.933	0.917	0.945	0.839	0.913	0.894	0.947
Telephone	0.970	0.975	0.970	0.978	0.942	0.975	0.971	0.984
Vessel	0.860	0.879	0.867	0.898	0.758	0.850	0.851	0.909
mean	0.895	0.908	0.892	0.922	0.806	0.883	0.869	0.924

Table E. Performance on ShapeNet with input noisy point cloud 300. Boldface font represents the preferred results.

Methods	Links
SPSR [5]	https://github.com/mmolero/pypoisson
ONet [8]	https://github.com/autonomousvision/occupancy_networks
ConvONet [9]	https://github.com/autonomousvision/convolutional_occupancy_networks
DP-ConvONet [6]	https://github.com/dsvilarkovic/dynamic_plane_convolutional_onet
POCO [1]	https://github.com/valeoai/POCO

Table F. The link for the baseline methods we compare.

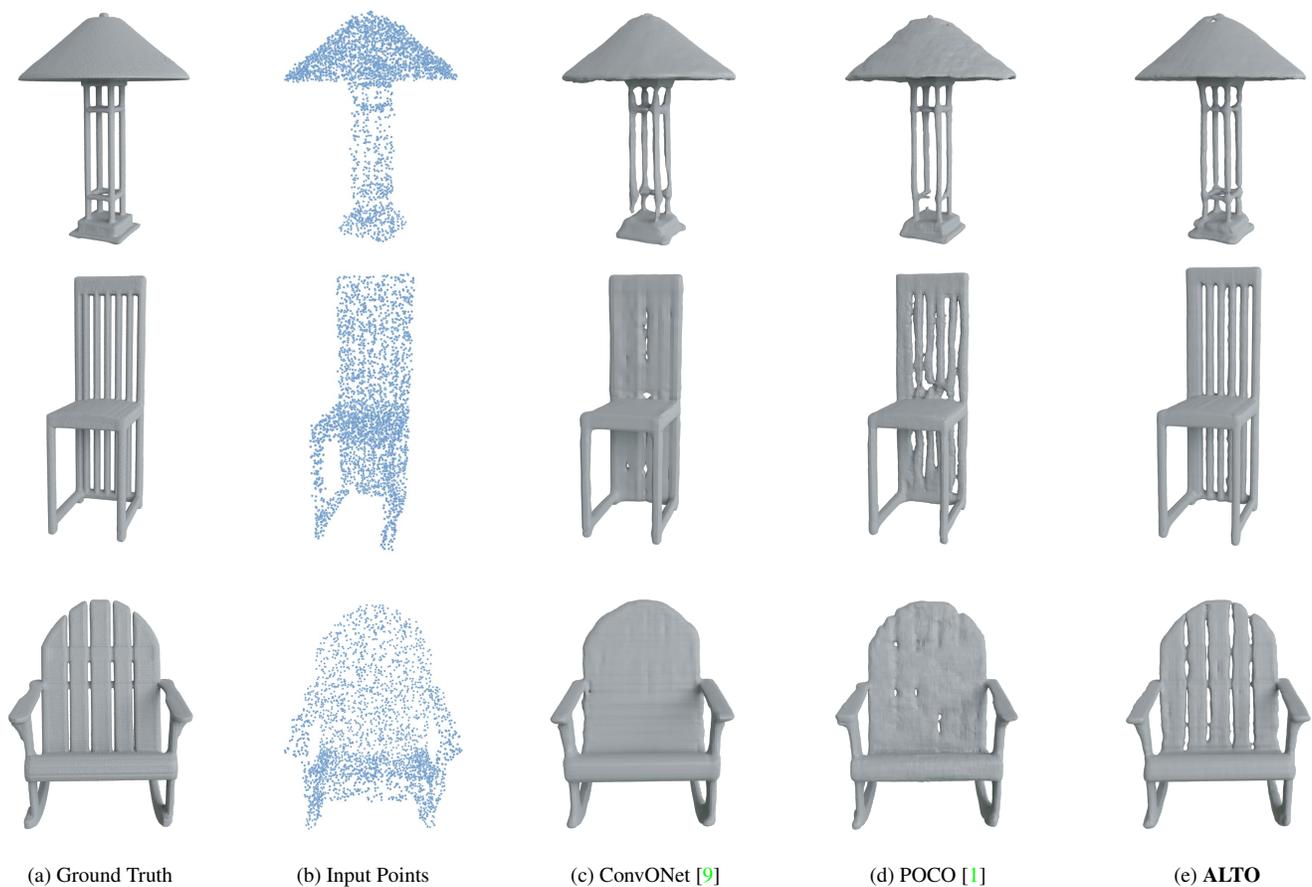


Figure A. **Qualitative comparison on object-level reconstruction ShapeNet dataset.** Trained and tested on 3k noisy points.

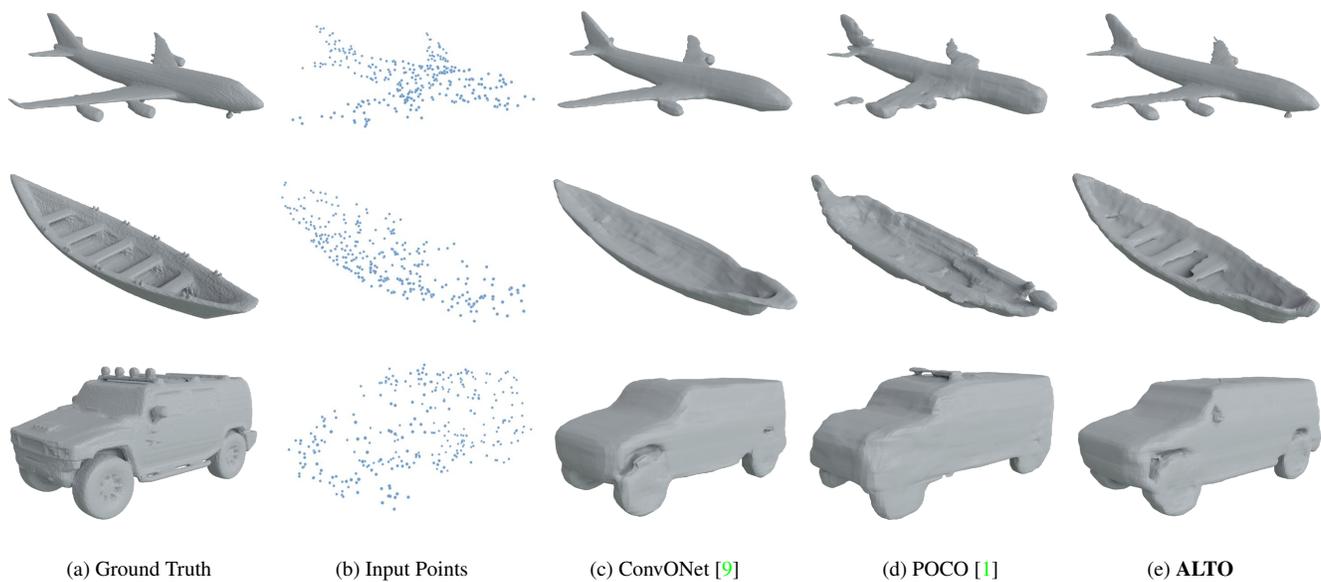


Figure B. **Qualitative comparison on object-level reconstruction ShapeNet dataset.** Trained and tested on 300 noisy points.

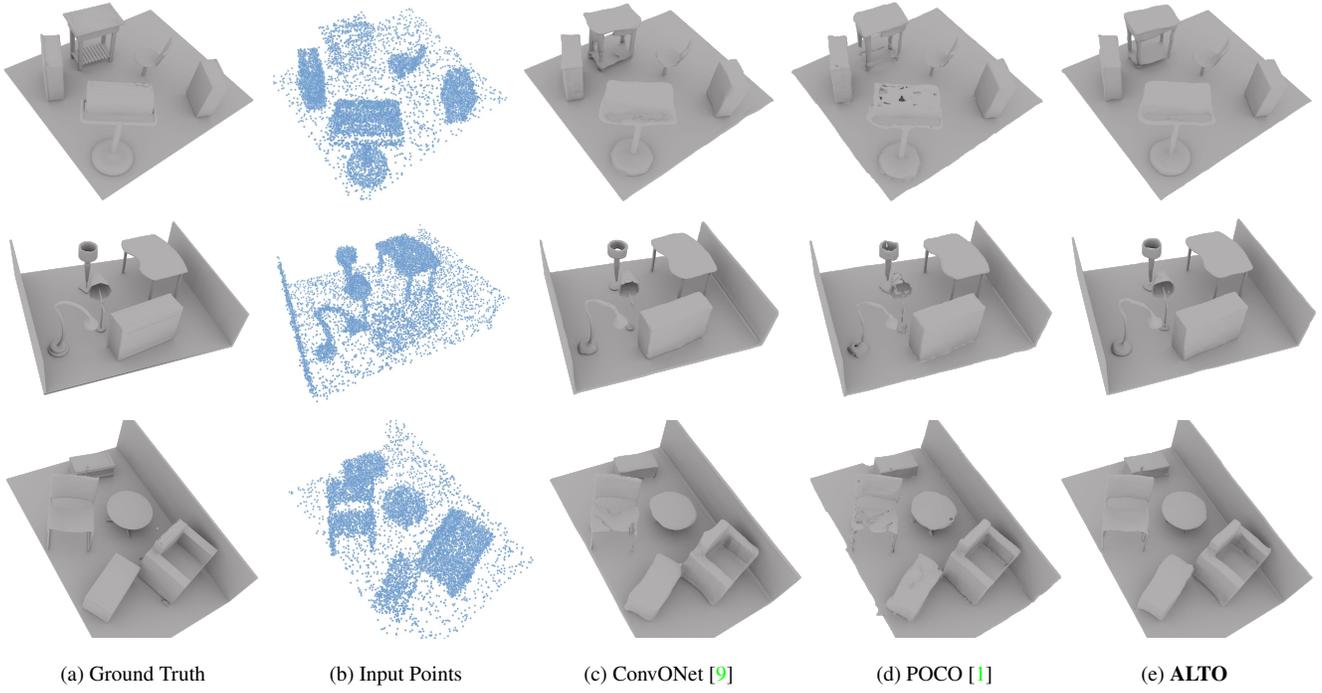


Figure C. **Qualitative comparison on scene-level reconstruction Synthetic Room dataset.** Trained and tested on 10k noisy points.

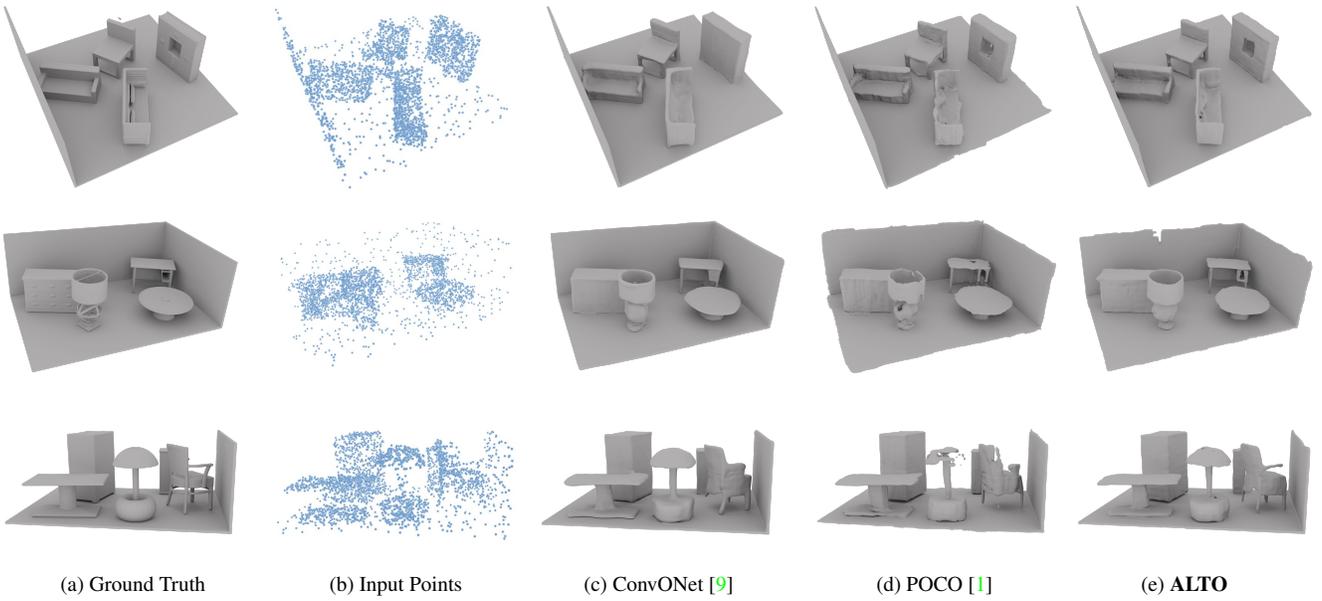


Figure D. **Qualitative comparison on scene-level reconstruction Synthetic Room dataset.** Trained and tested on 3K noisy points.

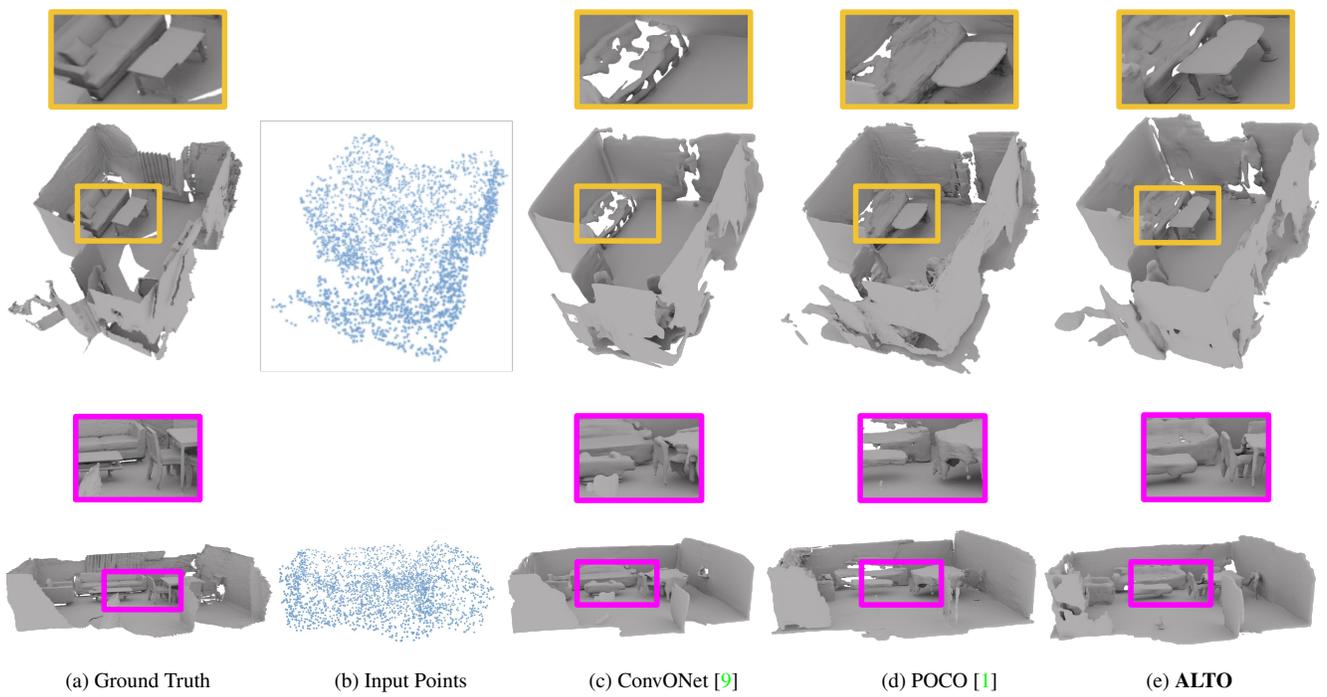


Figure E. Qualitative comparison on scene-level reconstruction ScanNet.

References

- [1] Alexandre Boulch and Renaud Marlet. Poco: Point convolution for surface reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6302–6314, 2022. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [2] Özgün Çiçek, Ahmed Abdulkadir, Soeren S Lienkamp, Thomas Brox, and Olaf Ronneberger. 3d u-net: learning dense volumetric segmentation from sparse annotation. In *International conference on medical image computing and computer-assisted intervention*, pages 424–432. Springer, 2016. [1](#)
- [3] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. [1](#)
- [4] Zhenyu Jiang, Yifeng Zhu, Maxwell Svetlik, Kuan Fang, and Yuke Zhu. Synergies between affordance and geometry: 6-dof grasp detection via implicit representations. *arXiv preprint arXiv:2104.01542*, 2021. [4](#)
- [5] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. [5](#)
- [6] Stefan Lionar, Daniil Emtsev, Dusan Svilarukovic, and Songyou Peng. Dynamic plane convolutional occupancy networks. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pages 1829–1838, 2021. [5](#)
- [7] 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. [2](#)
- [8] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4460–4470, 2019. [3](#), [4](#), [5](#)
- [9] Songyou Peng, Michael Niemeyer, Lars Mescheder, Marc Pollefeys, and Andreas Geiger. Convolutional occupancy networks. In *European Conference on Computer Vision*, pages 523–540. Springer, 2020. [2](#), [3](#), [4](#), [5](#), [6](#), [7](#), [8](#)
- [10] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 652–660, 2017. [1](#)
- [11] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *International Conference on Medical image computing and computer-assisted intervention*, pages 234–241. Springer, 2015. [1](#)
- [12] Suhani Vora, Noha Radwan, Klaus Greff, Henning Meyer, Kyle Genova, Mehdi SM Sajjadi, Etienne Pot, Andrea Tagliasacchi, and Daniel Duckworth. Nesf: Neural semantic fields for generalizable semantic segmentation of 3d scenes. *arXiv preprint arXiv:2111.13260*, 2021. [4](#)