

# SG-NN: Sparse Generative Neural Networks for Self-Supervised Scene Completion of RGB-D Scans

Angela Dai

Christian Diller

Matthias Nießner

Technical University of Munich

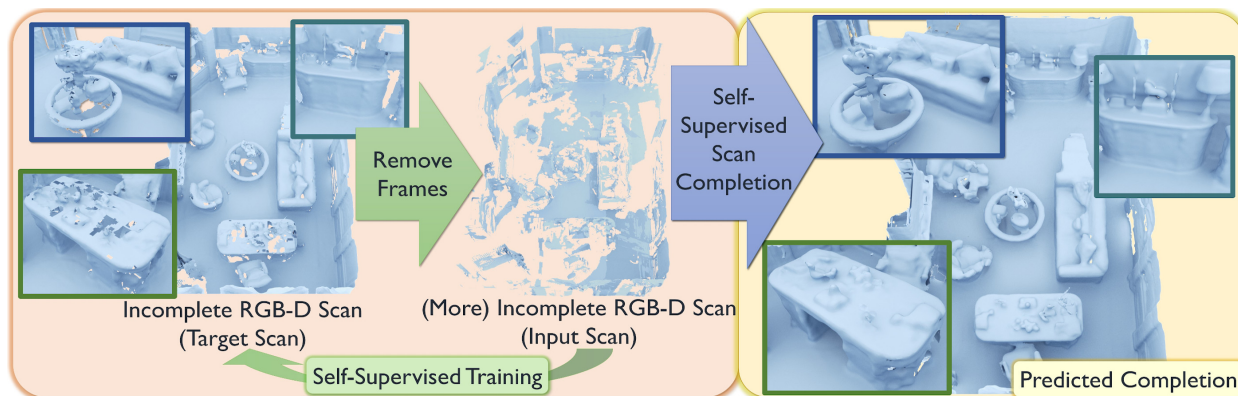


Figure 1: Our method takes as input a partial RGB-D scan and predicts a high-resolution 3D reconstruction while predicting unseen, missing geometry. Key to our approach is its self-supervised formulation, enabling training solely on real-world, incomplete scans. This not only obviates the need for synthetic ground truth, but is also capable of generating more complete scenes than any single target scene seen during training. To achieve high-quality surfaces, we further propose a new sparse generative neural network, capable of generating large-scale scenes at much higher resolution than existing techniques.

## Abstract

We present a novel approach that converts partial and noisy RGB-D scans into high-quality 3D scene reconstructions by inferring unobserved scene geometry. Our approach is fully self-supervised and can hence be trained solely on real-world, incomplete scans. To achieve self-supervision, we remove frames from a given (incomplete) 3D scan in order to make it even more incomplete; self-supervision is then formulated by correlating the two levels of partialness of the same scan while masking out regions that have never been observed. Through generalization across a large training set, we can then predict 3D scene completion without ever seeing any 3D scan of entirely complete geometry. Combined with a new 3D sparse generative neural network architecture, our method is able to predict highly-detailed surfaces in a coarse-to-fine hierarchical fashion, generating 3D scenes at 2cm resolution, more than twice the resolution of existing state-of-the-art methods as well as outperforming them by a significant margin in reconstruction quality.<sup>1</sup>

<sup>1</sup>Source code [available here](#).

## 1. Introduction

In recent years, we have seen incredible progress on RGB-D reconstruction of indoor environments using commodity RGB-D sensors such as the Microsoft Kinect, Google Tango, or Intel RealSense [22, 16, 23, 36, 5, 10]. However, despite remarkable achievements in RGB-D tracking and reconstruction quality, a fundamental challenge still remains – the incomplete nature of resulting 3D scans caused by inherent occlusions due to the physical limitations of the scanning process; i.e., even in a careful scanning session it is inevitable that some regions of a 3D scene remain unobserved. This unfortunately renders the resulting reconstructions unsuitable for many applications, not only those that require quality 3D content, such as video games or AR/VR, but also robotics where a completed 3D map significantly facilitates tasks such as grasping or querying 3D objects in a 3D environment.

In order to overcome the incomplete and partial nature of 3D reconstructions, various geometric inpainting techniques have been proposed, for instance, surface interpolation based on the Poisson equation [17, 18] or CAD shape-fitting techniques [1, 2, 8]. A very recent direction leverages

generative deep neural networks, often focusing volumetric representations for shapes [11] or entire scenes [31, 12]. These techniques show great promise since they can learn generalized patterns in a large variety of environments; however, existing data-driven scene completion methods rely on supervised training, requiring fully complete ground truth 3D models, thus depending on large-scale synthetic datasets such as ShapeNet [4] or SUNCG [31]. As a result, although we have seen impressive results from these approaches on synthetic test sets, domain transfer and application to real-world 3D scans remains a major limitation.

In order to address the shortcomings of supervised learning techniques for scan completion, we propose a new self-supervised completion formulation that can be trained only on (partial) real-world data. Our main idea is to learn to generate more complete 3D models from less complete data, while masking out any unknown regions; that is, from an existing RGB-D scan, we use the scan as the target and remove frames to obtain a more incomplete input. In the loss function, we can now correlate the difference in partialness between the two scans, and constrain the network to predict the delta while masking out unobserved areas. Although there is no single training sample which contains a fully-complete 3D reconstruction, we show that our network can nonetheless generalize to predict high levels of completeness through a combined aggregation of patterns across the entire training set. This way, our approach can be trained without requiring any fully-complete ground truth counterparts that would make generalization through a synthetic-to-real domain gap challenging.

Furthermore, we propose a new sparse generative neural network architecture that can predict high-resolution geometry in a fully-convolutional fashion. For training, we propose a progressively growing network architecture trained in coarse-to-fine fashion; i.e., we first predict the 3D scene at a low resolution, and then continue increasing the surface resolution through the training process. We show that our self-supervised, sparse generative approach can outperform state-of-the-art fully-supervised methods, despite their access to much larger quantities of synthetic 3D data.

We present the following main contributions:

- A self-supervised approach for scene completion, enabling training solely on incomplete, real-world scan data while predicting geometry more complete than any seen during training, by leveraging common patterns in the deltas of incompleteness.
- A generative formulation for sparse convolutions to produce a sparse truncated signed distance function representation at high resolution: we formulate this hierarchically to progressively generate a 3D scene in end-to-end fashion

## 2. Related Work

**RGB-D Reconstruction** Scanning and reconstructing 3D surfaces has a long history across several research communities. With the increase in availability of commodity range sensors, capturing and reconstructing 3D scenes has become a vital area of research. One seminal technique is the volumetric fusion approach of Curless and Levoy [7], operating on truncated signed distance fields to produce a surface reconstruction. It has been adopted by many state-of-the-art real-time reconstruction methods, from KinectFusion [22, 16] to VoxelHashing [23] and BundleFusion [10], as well as state-of-the-art offline reconstruction approaches [5].

These methods have produced impressive results in tracking and scalability of 3D reconstruction from commodity range sensors. However, a significant limitation that still remains is the partial nature of 3D scanning; i.e., a perfect scan is usually not possible due to occlusions and unobserved regions and thus, the resulting 3D representation cannot reach the quality of manually created 3D assets.

**Deep Learning on 3D Scans** With recent advances in deep learning and the improved availability of large-scale 3D scan datasets such as ScanNet [9] or Matterport [3], learned approaches on 3D data can be used for a variety of tasks like classification, segmentation, or completion.

Many current methods make use of convolutional operators that have been shown to work well on 2D data. When extended into 3D, they operate on regular grid representations such as distance fields [11] or occupancy grids [20]. Since dense volumetric grids can come with high computational and memory costs, several recent approaches have leveraged the sparsity of the 3D data for discriminative 3D tasks. PointNet [26, 27] introduced a deep network architecture for learning on point cloud data for semantic segmentation and classification tasks. Octree-based approaches have also been developed [29, 33, 34] that have been shown to be very memory efficient; however, generative tasks involving large, varying-sized environments seems challenging and octree generation has only been shown for single ShapeNet-style objects [28, 32]. Another option leveraging the sparsity of 3D geometric data is through sparse convolutions [14, 13, 6], which have seen success in discriminative tasks such as semantic segmentation, but not in the context of generative 3D modeling tasks, where the overall structure of the scene is unknown.

**Shape and Scene Completion** Completing 3D scans has been well-studied in geometry processing. Traditional methods, such as Poisson Surface Reconstruction [17, 18], locally optimize for a surface to fit to observed points, and work well for small missing regions. Recently, various

deep learning-based approaches have been developed with greater capacity for learning global structures of shapes, enabling compelling completion of larger missing regions in scans of objects [37, 11, 15, 35, 24]. Larger-scale completion of scans has been seen with SSCNet [31], operating on a depth image of a scene, and ScanComplete [12], which demonstrated scene completion on room- and building floor-scale scans. However, both these approaches operate on dense volumetric grids, significantly limiting their output resolutions. Moreover, these approaches are fully supervised with complete 3D scene data, requiring training on synthetic 3D scene data (where complete ground truth is known), in order to complete real-world scans.

An alternative approach for shape completion could through leveraging a single implicit latent space, as in DeepSDF [24] or Occupancy Networks [21]; however, it still remains a challenge as to how to scale a single latent space to represent large, varying-sized environments.

### 3. Method Overview

From an RGB-D scan of a 3D scene, our method learns to generate a high-quality reconstruction of the complete 3D scene, in a self-supervised fashion. The input RGB-D scan is represented as a truncated signed distance field (TSDF), as a sparse set of voxel locations within truncation and their corresponding distance values. The output complete 3D model of the scene is also generated as a sparse TSDF (similarly, a set of locations and per-voxel distances), from which a mesh can be extracted by Marching Cubes [19].

We design the 3D scene completion as a self-supervised process, enabling training purely on real-world scan data without requiring any fully-complete ground truth scenes. Since real-world scans are always incomplete due to occlusions and physical sensor limitations, this is essential for generating high-quality, complete models from real-world scan data. To achieve self-supervision, our main idea is to

formulate the training from incomplete scan data to less incomplete scan data; that is, from an existing RGB-D scan we can remove frames in order to create a more partial observation of the scene. This enables learning to complete in regions where scan geometry is known while ignoring regions of unobserved space. Crucially, our generative model can then learn to generate more complete models than seen in a specific sample of the target data.

To obtain an output high-resolution 3D model of a scene, we propose Sparse Generative Neural Networks (SG-NN), a generative model to produce a sparse surface representation of a scene. We build upon sparse convolutions [14, 13, 6], which have been shown to produce compelling semantic segmentation results on 3D scenes by operating only on surface geometry. In contrast to these discriminative tasks where the geometric structure is given as input, we develop our SG-NN to generate new, unseen 3D geometry suitable for generative 3D modeling tasks. This is designed in coarse-to-fine fashion, with a progressively growing network architecture which predicts each next higher resolution, finally predicting a high-resolution surface as a sparse TSDF. Since our sparse generative network operates in a fully-convolutional fashion, we can operate on 3D scans of varying spatial sizes.

### 4. Self-Supervised Completion

Our approach for self-supervision of scene completion of RGB-D scans is based on learning how to complete scan geometry in regions that have been seen, while ignoring unobserved regions. To this end, we can generate input and target TSDFs with similar scanning patterns as real-world scans; from an input scan composed of RGB-D frames  $\{f_0, \dots, f_n\}$ , we can generate the target TSDF  $\mathcal{S}_{\text{target}}$  through volumetric fusion [7] of  $\{f_0, \dots, f_n\}$ , and the input TSDF  $\mathcal{S}_{\text{input}}$  through volumetric fusion of a subset of the original frames  $\{f_k\} \subset \{f_0, \dots, f_n\}$ .

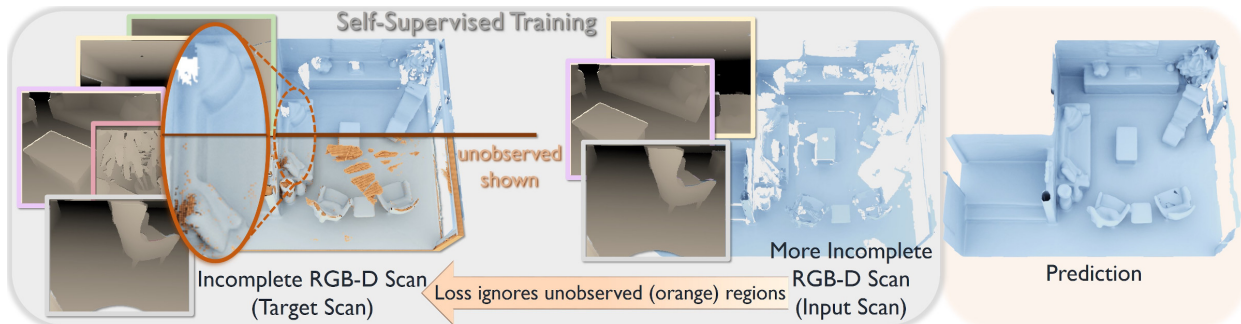


Figure 2: Our self-supervision approach for scan completion learns through deltas in partialness of RGB-D scans. From a given (incomplete) RGB-D scan, on the left, we produce a more incomplete version of the scan by removing some of its depth frames (middle). We can then train to complete the more incomplete scan (middle) using the original scan as a target (left), while masking out unobserved regions in the target scene (in orange). This enables our prediction to produce scenes that are more complete than the target scenes seen during training, as the training process effectively masks out incompleteness.

This produces input incomplete scans that maintain scanned data characteristics, as well as a correspondence between  $\mathcal{S}_{\text{input}}$  and  $\mathcal{S}_{\text{target}}$  going from a more incomplete scan to a less incomplete scan. Since  $\mathcal{S}_{\text{target}}$  remains nonetheless incomplete, we do not wish to use all of its data as the complete target for supervision, as this could result in contradictory signals in the training set (e.g., table legs have been seen in one scan but not in another, then it becomes unclear whether to generate table legs).

Thus, to effectively learn to generate a complete 3D model beyond even the completeness of the target training data, we formulate the completion loss only on observed regions in the target scan. That is, the loss is only considered in regions where  $\mathcal{S}_{\text{target}}(v) > -\tau$ , for a voxel  $v$  with  $\tau$  indicating the voxel size. Figure 2 shows an example  $\mathcal{S}_{\text{input}}$ ,  $\mathcal{S}_{\text{target}}$ , and prediction, with this self-supervision setup, we can learn to predict geometry that was unobserved in  $\mathcal{S}_{\text{target}}$ , e.g., occluded regions behind objects.

### 4.1. Data Generation

As input we consider an RGB-D scan comprising a set of depth images and their 6-DoF camera poses. For real-world scan data we use the Matterport3D [3] dataset, which contains a variety of RGB-D scans taken with a Matterport tripod setup. Note that for Matterport3D, we train and evaluate on the annotated room regions, whereas the raw RGB-D data is a sequence covering many different rooms, so we perform an approximate frame-to-room association by taking frames whose camera locations lie within the room.

From a given RGB-D scan, we construct the target scan  $\mathcal{S}_{\text{target}}$  using volumetric fusion [7] with 2cm voxels and truncation of 3 voxels. A subset of the frames is taken by randomly removing  $\approx 50\%$  of the frames (see Section 6 for more analysis of varying degrees of incompleteness in  $\mathcal{S}_{\text{input}}$ ,  $\mathcal{S}_{\text{target}}$ ). We can then again use volumetric fusion to

generate a more incomplete version of the scan  $\mathcal{S}_{\text{input}}$ .

At train time, we consider cropped views of these scans for efficiency, using random crops of size  $64 \times 64 \times 128$  voxels. The fully-convolutional nature of our approach enables testing on full scenes of varying sizes at inference time.

## 5. Generating a Sparse 3D Scene Representation

The geometry of a 3D scene occupies a very sparse set of the total 3D extent of the scene, so we aim to generate a 3D representation of a scene in a similarly sparse fashion. Thus we propose Sparse Generative Neural Networks (SG-NN) to hierarchically generate a sparse, truncated signed distance field representation of a 3D scene, from which we can extract the isosurface as the final output mesh.

An overview of our network architecture for the scene completion task is shown in Figure 3. The model is designed in encoder-decoder fashion, with an input partial scan first encoded to representative features at low spatial resolution, before generating the final TSDF output.

A partial scan, represented as a TSDF, is encoded with a series of 3D sparse convolutions [14, 13] which operate only on the locations where the TSDF is within truncation distance and using the distance values as input features. Each set of convolutions spatially compresses the scene by a factor of two. Our generative model takes the encoding of the scene and converts the features into a (low-resolution) dense 3D grid. The dense representation enables prediction of the full scene geometry at very coarse resolution; here, we use a series of dense 3D convolutions to produce a feature map  $F_0$  from which we also predict coarse occupancy  $O_0$  and TSDF  $S_0$  representations of the complete scene. We then construct a sparse representation of the predicted scene based on  $O_0$ : the features input to the next level are

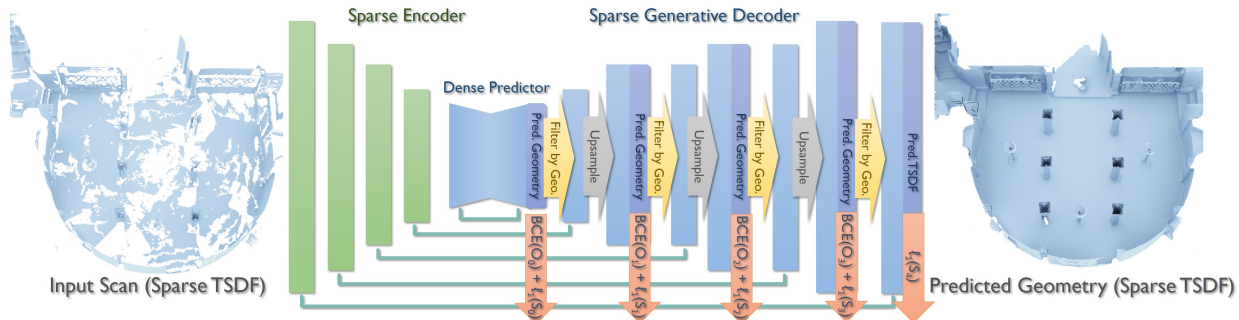


Figure 3: Our Sparse Generative Neural Network architecture for the task of scan completion. An input scan is encoded using a series of sparse convolutions, each set reducing the spatial dimensions by a factor of two. To generate high-resolution scene geometry, the coarse encoding is converted to a dense representation for a coarse prediction of the complete geometry. The predicted coarse geometry is converted to a sparse representation and input to our sparse, coarse-to-fine hierarchy, where each level of the hierarchy predicts the geometry of the next resolution (losses indicated in orange). The final output is a TSDF represented by sparse set of voxel locations and their corresponding distance values.



composed as  $\text{concat}(F_k, O_k, S_k) \forall \text{sigmoid}(O_k(v)) > 0.5$ . This can then be processed with sparse convolutions, then upsampled by a factor of two to predict the scene geometry at the next higher resolution. This enables generative, sparse predictions in a hierarchical fashion. To generate the final surface geometry, the last hierarchy level of our SG-NN outputs sparse  $O_n$ ,  $S_n$ , and  $F_n$ , which are then input to a final set of sparse convolutions to refine and predict the output signed distance field values.

**Sparse skip connections.** For scene completion, we also leverage skip connections between the encoder and decoder parts of the network architecture, connecting feature maps of same spatial resolution. This is in the same spirit as U-Net [30], but in our case the encoder and decoder feature maps are both sparse and typically do not contain the same set of sparse locations. Thus we concatenate features from the set of source locations which are shared with the destination locations, and use zero feature vectors for the destination locations which do not exist in the source.

**Progressive Generation.** In order to encourage more efficient and stable training, we train our generative model progressively, starting with the lowest resolution, and introducing each successive hierarchy level after  $N_{\text{level}}$  iterations. Each hierarchy level predicts the occupancy and TSDF of the next level, enabling successive refinement from coarse predictions, as shown in Figure 4.

**Loss.** We formulate the loss for the generated scene geometry on the final predicted TSDF locations and values, using an  $\ell_1$  loss with the target TSDF values at those locations. Following [11], we log-transform the TSDF values of the predictions and the targets before applying the  $\ell_1$  loss, in order to encourage more accurate prediction near the surface geometry. We additionally employ proxy losses at each hierarchy level for outputs  $O_k$  and  $S_k$ , using binary cross entropy with target occupancies and  $\ell_1$  with target TSDF values, respectively. This helps avoid a trivial solution of zero loss for the final surface with no predicted geometry. Note that for our self-supervised completion, we compute these losses only in regions of observed target values, as described in Section 4.

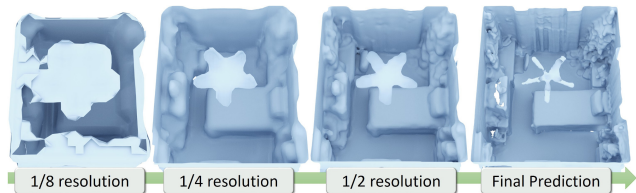


Figure 4: Progressive generation of a 3D scene using our SG-NN which formulates a generative model to predict a sparse TSDF as output.

## 5.1. Training

We train our SG-NN on a single NVIDIA GeForce RTX 2080, using the Adam optimizer with a learning rate of 0.001 and batch size of 8. We use  $N_{\text{level}} = 2000$  iterations for progressive introduction of each higher resolution output, and train our model for  $\approx 40$  hours until convergence.

## 6. Results and Evaluation

We evaluate our sparse generative neural network on scene completion for RGB-D scans on both real-world scans where no fully complete ground truth is available [3], as well as in a supervised setting on synthetic scans which have complete ground truth information [31]. We use the train/test splits provided by both datasets: 72/18 and 5519/155 trainval/test scenes comprising 1788/394 and 39600/1000 rooms, respectively. To measure completion quality, we follow [12] and use an  $\ell_1$  error metric between predicted and target TSDFs, where unobserved regions in the target are masked out. Note that unsigned distances are used in the error computation to avoid sign ambiguities. We measure the  $\ell_1$  distance in voxel units of the entire volume (*entire volume*), the unobserved region of the volume (*unobserved space*), near the target surface (*target*), and near the predicted surface (*predicted*), using a threshold of  $\leq 1$  to determine nearby regions, and a global truncation of 3. For all metrics, unobserved regions in the targets are ignored; note that on synthetic data where complete ground truth is available, we do not have any unobserved regions to ignore.

**Comparison to state of the art.** In Table 1, we compare to several state-of-the-art approaches for scan completion on real-world scans from the Matterport3D dataset [3]: the shape completion approach 3D-EPN [11], and the scene completion approach ScanComplete [12]. These methods both require fully-complete ground truth data for supervision, which is not available for the real-world scenes, so we train them on synthetic scans [31]. Since 3D-EPN and ScanComplete use dense 3D convolutions, limiting voxel resolution, we use 5cm resolution for training and evaluation of all methods. Our self-supervised approach enables training on incomplete real-world scan data, avoiding domain transfer while outperforming previous approaches that leverage large amounts of synthetic 3D data. Qualitative comparisons are shown in Figure 5.

To evaluate our SG-NN separate from its self-supervision, we also evaluate synthetic scan completion with full ground truth [31], in comparison to Poisson Surface Reconstruction [17, 18], SSCNet[31], 3D-EPN [11], and ScanComplete [12]. All data-driven approaches are fully supervised, using input scans from [12]. Similar to the real scan scenario, we train and evaluate at 5cm resolution due to resolution limitations of the prior learned approaches.

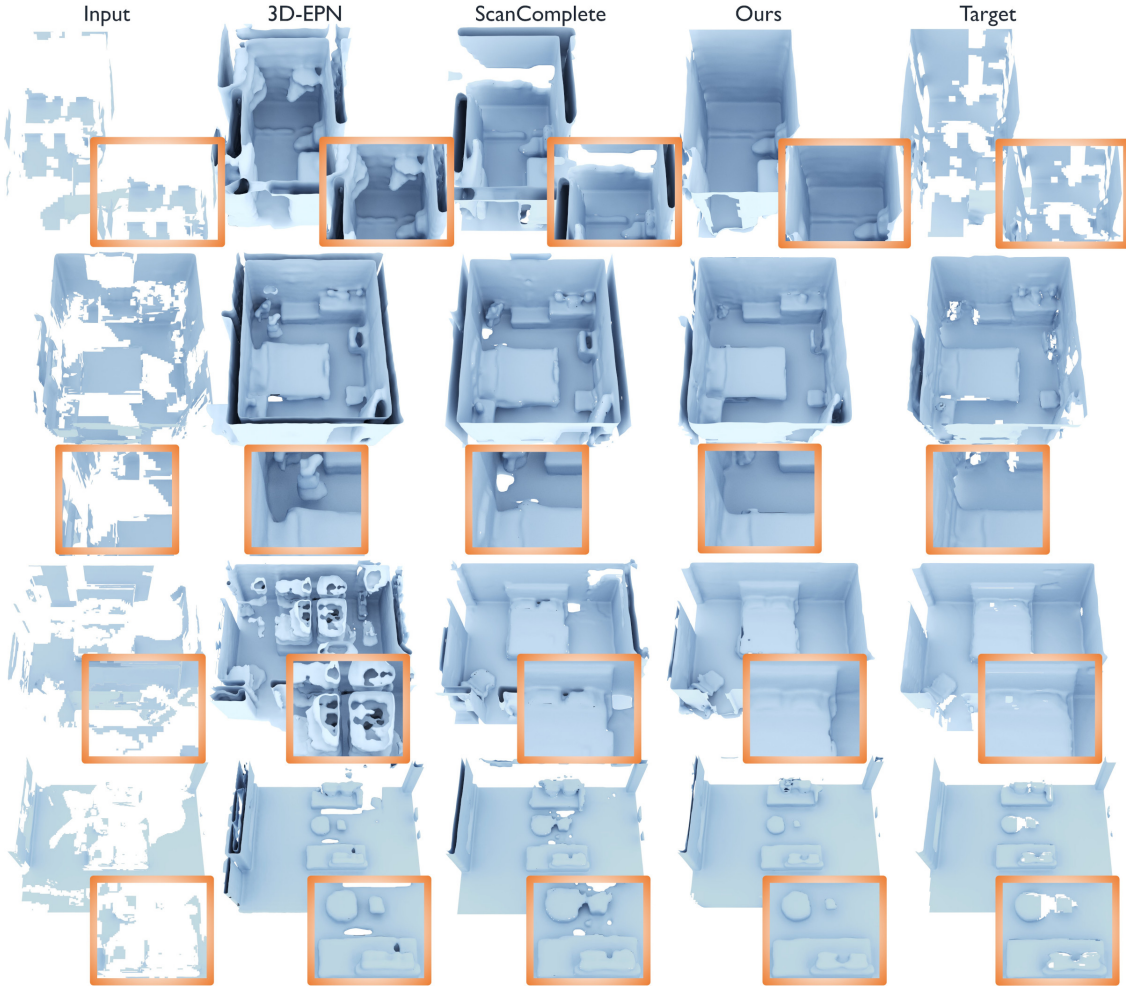


Figure 5: Comparison to state-of-the-art scan completion approaches on Matterport3D [3] data (5cm resolution), with input scans generated from a subset of frames. In contrast to the fully-supervised 3D-EPN [11] and ScanComplete [12], our self-supervised approach produces more accurate, complete scene geometry.

Method	$\ell_1$ error entire volume	$\ell_1$ error unobserved space	$\ell_1$ error target	$\ell_1$ error predicted
3D-EPN (unet) [11]	0.31	0.28	0.45	1.12
ScanComplete [12]	0.20	0.15	0.51	0.74
Ours	<b>0.17</b>	<b>0.14</b>	<b>0.35</b>	<b>0.67</b>

Table 1: Quantitative scan completion results on real-world scan data [3], with  $\ell_1$  distance measured in voxel units for 5cm voxels. Since target scans are incomplete, unobserved space in the target is masked out for all metrics. 3D-EPN [11] and ScanComplete [12] require full supervision, and so are trained on synthetic data [31]. Despite their access to large quantities of synthetic 3D data, our self-supervised approach outperforms these methods while training solely on real-world data.

In Table 2, we see that our sparse generative approach outperforms state of the art in a fully-supervised scenario.

**Can self-supervision predict more complete geometry than seen during training?** Our approach to self-supervision is designed to enable prediction of scene geometry beyond the completeness of the target scan data, by

leveraging knowledge of observed and unobserved space in RGB-D scans. To evaluate the completion quality of our method against the completeness of the target scene data, we perform a qualitative evaluation, as we lack fully complete ground truth to for quantitative evaluation. In Figure 7, we see that our completion quality can exceed the completeness of target scene data. We additionally evaluate our

Method	$\ell_1$ error entire volume	$\ell_1$ error unobserved space	$\ell_1$ error target	$\ell_1$ error predicted
Poisson Surface Reconstruction [17, 18]	0.53	0.51	1.70	1.18
SSCNet [31]	0.54	0.53	0.93	1.11
3D-EPN (unet) [11]	0.25	0.30	0.65	0.47
ScanComplete [12]	0.18	0.23	0.53	0.42
Ours	<b>0.15</b>	<b>0.16</b>	<b>0.50</b>	<b>0.28</b>

Table 2: Quantitative scan completion results on synthetic scan data [31], where complete ground truth is available to supervise all data-driven approaches.  $\ell_1$  distance is measured in voxel units for 5cm voxels.

approach with and without our self-supervision masking in Figure 7, where *w/o self-supervision masking* is trained using the same set of less-incomplete/more-incomplete scans but without the loss masking. This can perform effective completion in regions commonly observed in target scans, but often fails to complete regions that are commonly occluded. In contrast, our formulation for self-supervision using masking of unobserved regions enables predicting scene geometry even where the target scan remains incomplete.

**Comparison of our self-supervision approach to masking out by random crops.** In Table 3, we evaluate against another possible self-supervision approach: randomly cropping out target geometry to be used as incomplete inputs (*using crops for self-supervision*), similar to [25]. This scenario does not reflect the data characteristics of real-world scan partialness (e.g., from occlusions and lack of visibility), resulting in poor completion performance.

**What’s the impact of the input/output representation?** In Table 3, we evaluate the effect of a point cloud input (vs. TSDF input), as well as occupancy output (vs. TSDF output). We find that the TSDF representation has more potential descriptiveness in characterizing a surface (and its neighboring regions), resulting in improved performance in both input and output representation.

**What’s the impact of the degree of completeness of the target data during training?** In Figure 6, we evaluate the effect of the amount of completeness of the target data available for training. We create several incomplete versions of

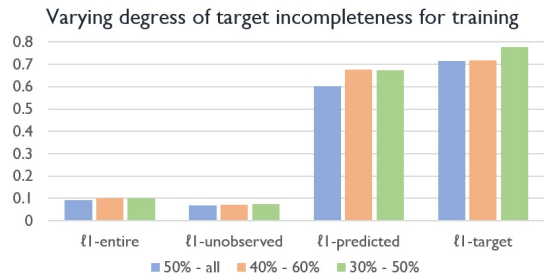


Figure 6: Evaluating varying target data completeness available for training. We generate various incomplete versions of the Matterport3D [3] scans using  $\approx$  30%, 40%, 50%, 60%, and 100% (all) of the frames associated with each room scene, and evaluate on the 50% incomplete scans. Our self-supervised approach remains robust to the level of completeness of the target training data.

the Matterport3D [3] scans using varying amounts of the frames available:  $\approx$  30%, 40%, 50%, 60%, and 100% (all) of the frames associated with each room scene. We train our approach using three different versions of input-target completeness: 50% – *all* (our default), 40% – 60%, and 30% – 50%. Even as the completeness of the target data decreases, our approach maintains robustness in predicting complete scene geometry.

**Limitations** Our SG-NN approach for self-supervised scan completion enables high-resolution geometric prediction of complete geometry from real-world scans. However, to generate the full appearance of a 3D scene, generation and inpainting of color is also required. Currently,

Method	$\ell_1$ error entire volume	$\ell_1$ error unobserved space	$\ell_1$ error target	$\ell_1$ error predicted
Using crops for self-supervision	0.13	0.09	1.25	0.68
Point cloud input	0.15	0.09	1.82	0.92
Occupancy output	0.13	0.10	0.89	0.86
2 hierarchy levels	0.10	0.08	0.74	0.68
Ours	<b>0.09</b>	<b>0.07</b>	<b>0.71</b>	<b>0.60</b>

Table 3: Ablation study of our self-supervision and generative model design choices on real-world scan data [3], with  $\ell_1$  distance measured in voxel units for 2cm voxels.



our method also does not consider or predict the semantic object decomposition of a scene; however, we believe this would be an interesting direction, specifically in the context for enabling interaction with a 3D environment (e.g., interior redesign or robotic understanding).

## 7. Conclusion

In this paper, we presented a self-supervised approach for completion of RGB-D scan geometry that enables training solely on incomplete, real-world scans while learning a generative geometric completion process capable of predicting 3D scene geometry more complete than any single target scene seen during training. Our sparse generative approach to generating a sparse TSDF representation of a sur-

face enables much higher output geometric resolution than previous on large-scale 3D scenes. Self-supervision allowing training only on real-world scan data for scan completion opens up new possibilities for various generative 3D modeling based only on real-world observations, perhaps mitigating the need for extensive synthetic data generation or domain transfer, and we believe this is a promising avenue for future research.

## Acknowledgments

This work was supported by the ZD.B, a Google Research Grant, a TUM-IAS Rudolf Mößbauer Fellowship, an NVidia Professorship Award, and the ERC Starting Grant Scan2CAD (804724).

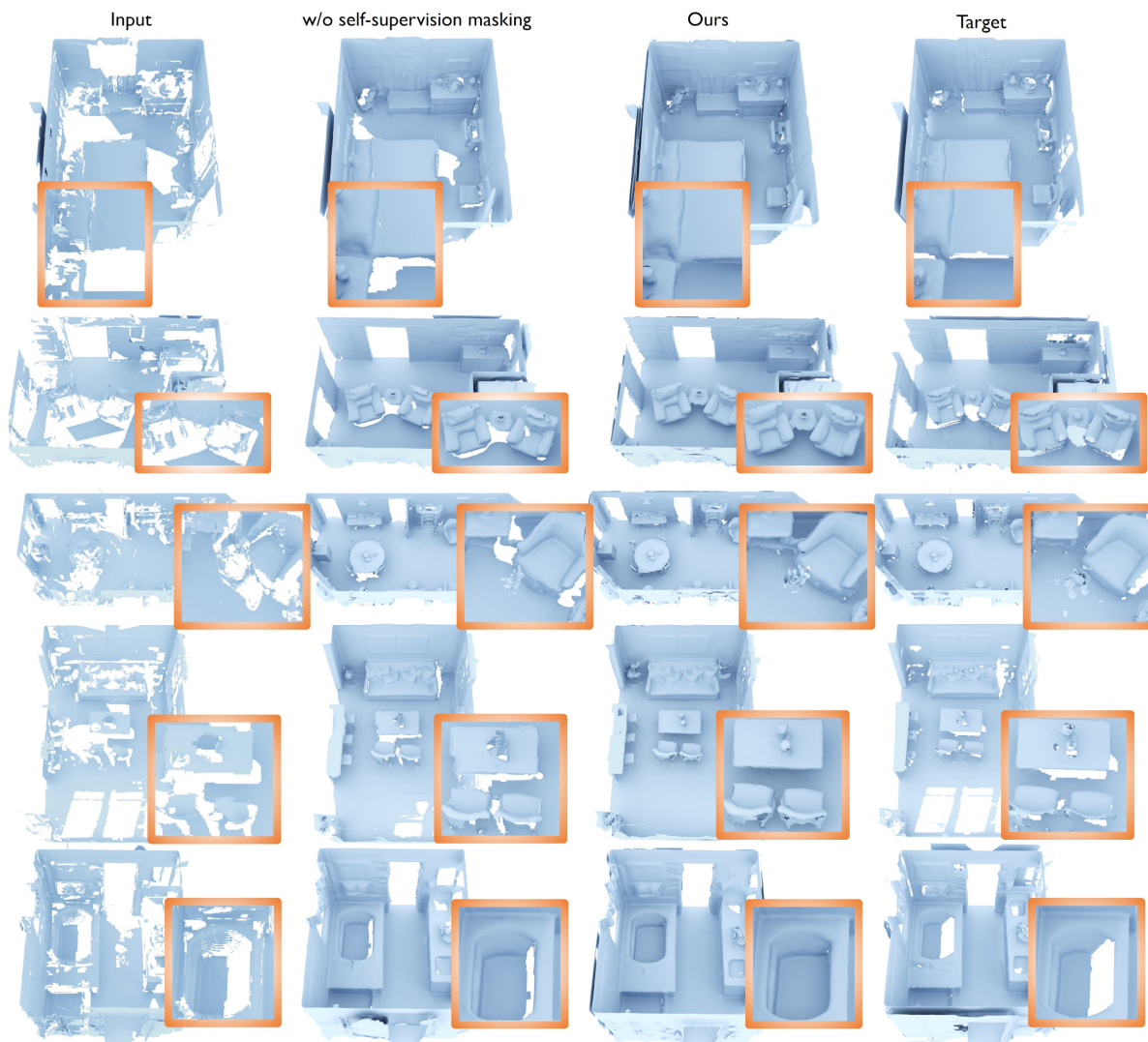


Figure 7: Scan completion results on Matterport3D [3] data (2cm resolution), with input scans generated from a subset of frames. Our self-supervision approach using loss masking enables more complete scene prediction than direct supervision using the target RGB-D scan, particularly in regions where occlusions commonly occur.



## References

- [1] Armen Avetisyan, Manuel Dahnert, Angela Dai, Manolis Savva, Angel X. Chang, and Matthias Nießner. Scan2cad: Learning CAD model alignment in RGB-D scans. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 2614–2623, 2019. **1**
- [2] Armen Avetisyan, Angela Dai, and Matthias Niessner. End-to-end cad model retrieval and 9dof alignment in 3d scans. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. **1**
- [3] Angel X. Chang, Angela Dai, Thomas A. Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from RGB-D data in indoor environments. In *2017 International Conference on 3D Vision, 3DV 2017, Qingdao, China, October 10-12, 2017*, pages 667–676, 2017. **2, 4, 5, 6, 7, 8, 11, 12**
- [4] 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. **2**
- [5] Sungjoon Choi, Qian-Yi Zhou, and Vladlen Koltun. Robust reconstruction of indoor scenes. In *2015 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 5556–5565. IEEE, 2015. **1, 2**
- [6] Christopher B. Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 3075–3084, 2019. **2, 3**
- [7] Brian Curless and Marc Levoy. A volumetric method for building complex models from range images. In *Proceedings of the 23rd annual conference on Computer graphics and interactive techniques*, pages 303–312. ACM, 1996. **2, 3, 4**
- [8] Manuel Dahnert, Angela Dai, Leonidas J. Guibas, and Matthias Niessner. Joint embedding of 3d scan and cad objects. In *The IEEE International Conference on Computer Vision (ICCV)*, October 2019. **1**
- [9] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017. **2**
- [10] Angela Dai, Matthias Nießner, Michael Zollhöfer, Shahram Izadi, and Christian Theobalt. Bundlefusion: Real-time globally consistent 3d reconstruction using on-the-fly surface reintegration. *ACM Trans. Graph.*, 36(3):24:1–24:18, 2017. **1, 2**
- [11] Angela Dai, Charles Ruizhongtai Qi, and Matthias Nießner. Shape completion using 3d-encoder-predictor cnns and shape synthesis. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6545–6554, 2017. **2, 3, 5, 6, 7**
- [12] Angela Dai, Daniel Ritchie, Martin Bokeloh, Scott Reed, Jürgen Sturm, and Matthias Nießner. Scancomplete: Large-scale scene completion and semantic segmentation for 3d scans. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 4578–4587, 2018. **2, 3, 5, 6, 7**
- [13] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *2018 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2018, Salt Lake City, UT, USA, June 18-22, 2018*, pages 9224–9232, 2018. **2, 3, 4**
- [14] Benjamin Graham and Laurens van der Maaten. Submanifold sparse convolutional networks. *arXiv preprint arXiv:1706.01307*, 2017. **2, 3, 4**
- [15] Xiaoguang Han, Zhen Li, Haibin Huang, Evangelos Kalogerakis, and Yizhou Yu. High-resolution shape completion using deep neural networks for global structure and local geometry inference. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 85–93, 2017. **3**
- [16] Shahram Izadi, David Kim, Otmar Hilliges, David Molyneaux, Richard A. Newcombe, Pushmeet Kohli, Jamie Shotton, Steve Hodges, Dustin Freeman, Andrew J. Davison, and Andrew W. Fitzgibbon. Kinectfusion: real-time 3d reconstruction and interaction using a moving depth camera. In *Proceedings of the 24th Annual ACM Symposium on User Interface Software and Technology, Santa Barbara, CA, USA, October 16-19, 2011*, pages 559–568, 2011. **1, 2**
- [17] 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*, pages 61–70, 2006. **1, 2, 5, 7**
- [18] Michael M. Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Trans. Graph.*, 32(3):29:1–29:13, 2013. **1, 2, 5, 7**
- [19] William E. Lorensen and Harvey E. Cline. Marching cubes: A high resolution 3d surface construction algorithm. In *Proceedings of the 14th Annual Conference on Computer Graphics and Interactive Techniques, SIGGRAPH 1987, Anaheim, California, USA, July 27-31, 1987*, pages 163–169, 1987. **3**
- [20] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems, IROS 2015, Hamburg, Germany, September 28 - October 2, 2015*, pages 922–928, 2015. **2**
- [21] Lars Mescheder, Michael Oechsle, Michael Niemeyer, Sebastian Nowozin, and Andreas Geiger. Occupancy networks: Learning 3d reconstruction in function space. In *Proceedings IEEE Conf. on Computer Vision and Pattern Recognition (CVPR)*, 2019. **3**
- [22] Richard A. Newcombe, Shahram Izadi, Otmar Hilliges, David Molyneaux, David Kim, Andrew J. Davison, Pushmeet Kohli, Jamie Shotton, Steve Hodges, and Andrew W. Fitzgibbon. Kinectfusion: Real-time dense surface mapping and tracking. In *10th IEEE International Symposium on Mixed and Augmented Reality, ISMAR 2011, Basel, Switzerland, October 26-29, 2011*, pages 127–136, 2011. **1, 2**

- [23] M. Nießner, M. Zollhöfer, S. Izadi, and M. Stamminger. Real-time 3d reconstruction at scale using voxel hashing. *ACM Transactions on Graphics (TOG)*, 2013. [1](#), [2](#)
- [24] Jeong Joon Park, Peter Florence, Julian Straub, Richard A. Newcombe, and Steven Lovegrove. Deepsdf: Learning continuous signed distance functions for shape representation. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2019, Long Beach, CA, USA, June 16-20, 2019*, pages 165–174, 2019. [3](#)
- [25] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A. Efros. Context encoders: Feature learning by inpainting. In *2016 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2016, Las Vegas, NV, USA, June 27-30, 2016*, pages 2536–2544. IEEE Computer Society, 2016. [7](#)
- [26] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 77–85, 2017. [2](#)
- [27] Charles Ruizhongtai 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 30: Annual Conference on Neural Information Processing Systems 2017, 4-9 December 2017, Long Beach, CA, USA*, pages 5099–5108, 2017. [2](#)
- [28] Gernot Riegler, Ali Osman Ulusoy, Horst Bischof, and Andreas Geiger. Octnetfusion: Learning depth fusion from data. In *2017 International Conference on 3D Vision, 3DV 2017, Qingdao, China, October 10-12, 2017*, pages 57–66, 2017. [2](#)
- [29] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 6620–6629, 2017. [2](#)
- [30] 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. [5](#)
- [31] Shuran Song, Fisher Yu, Andy Zeng, Angel X. Chang, Manolis Savva, and Thomas A. Funkhouser. Semantic scene completion from a single depth image. In *2017 IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2017, Honolulu, HI, USA, July 21-26, 2017*, pages 190–198, 2017. [2](#), [3](#), [5](#), [6](#), [7](#)
- [32] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2107–2115, 2017. [2](#)
- [33] Peng-Shuai Wang, Yang Liu, Yu-Xiao Guo, Chun-Yu Sun, and Xin Tong. O-CNN: octree-based convolutional neural networks for 3d shape analysis. *ACM Trans. Graph.*, 36(4):72:1–72:11, 2017. [2](#)
- [34] Peng-Shuai Wang, Chun-Yu Sun, Yang Liu, and Xin Tong. Adaptive O-CNN: a patch-based deep representation of 3d shapes. *ACM Trans. Graph.*, 37(6):217:1–217:11, 2018. [2](#)
- [35] Weiyue Wang, Qiangui Huang, Suyu You, Chao Yang, and Ulrich Neumann. Shape inpainting using 3d generative adversarial network and recurrent convolutional networks. In *IEEE International Conference on Computer Vision, ICCV 2017, Venice, Italy, October 22-29, 2017*, pages 2317–2325, 2017. [3](#)
- [36] Thomas Whelan, Stefan Leutenegger, Renato F. Salas-Moreno, Ben Glocker, and Andrew J. Davison. Elasticfusion: Dense SLAM without a pose graph. In *Robotics: Science and Systems XI, Sapienza University of Rome, Rome, Italy, July 13-17, 2015*, 2015. [1](#)
- [37] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1912–1920, 2015. [3](#)

## A. SG-NN Architecture Details

Figure 8 details our Sparse Generative Neural Network specification for scan completion. Convolution parameters are given as (nf\_in, nf\_out, kernel\_size, stride, padding), with stride and padding default to 1 and 0 respectively. Arrows indicate concatenation, and  $\oplus$  indicates addition. Each convolution (except the last) is followed by batch normalization and a ReLU.

## B. Varying Target Data Incompleteness

Here, we aim to evaluate how well our self-supervision approach performs as the completeness of the target data seen during training decreases. As long as there is enough variety in the completion patterns seen during training, our approach can learn to generate scene geometry with high levels of completeness. To evaluate this, we generate several versions of target scans from the Matterport3D [3]

room scenes with varying degrees of completeness; that is, we use  $\approx 50\%$ ,  $60\%$ , and  $100\%$  of the frames associated with each room scene to generate three different levels of completeness in the target scans, using  $\approx 30\%$ ,  $40\%$ , and  $50\%$  for the respective input scans. We provide a quantitative evaluation in the main paper, and a qualitative evaluation in Figure 9. Even as the level of completeness in the target data used decreases, our approach maintains robustness its completion, informed by the deltas in incompleteness as to the patterns of generating complete geometry.

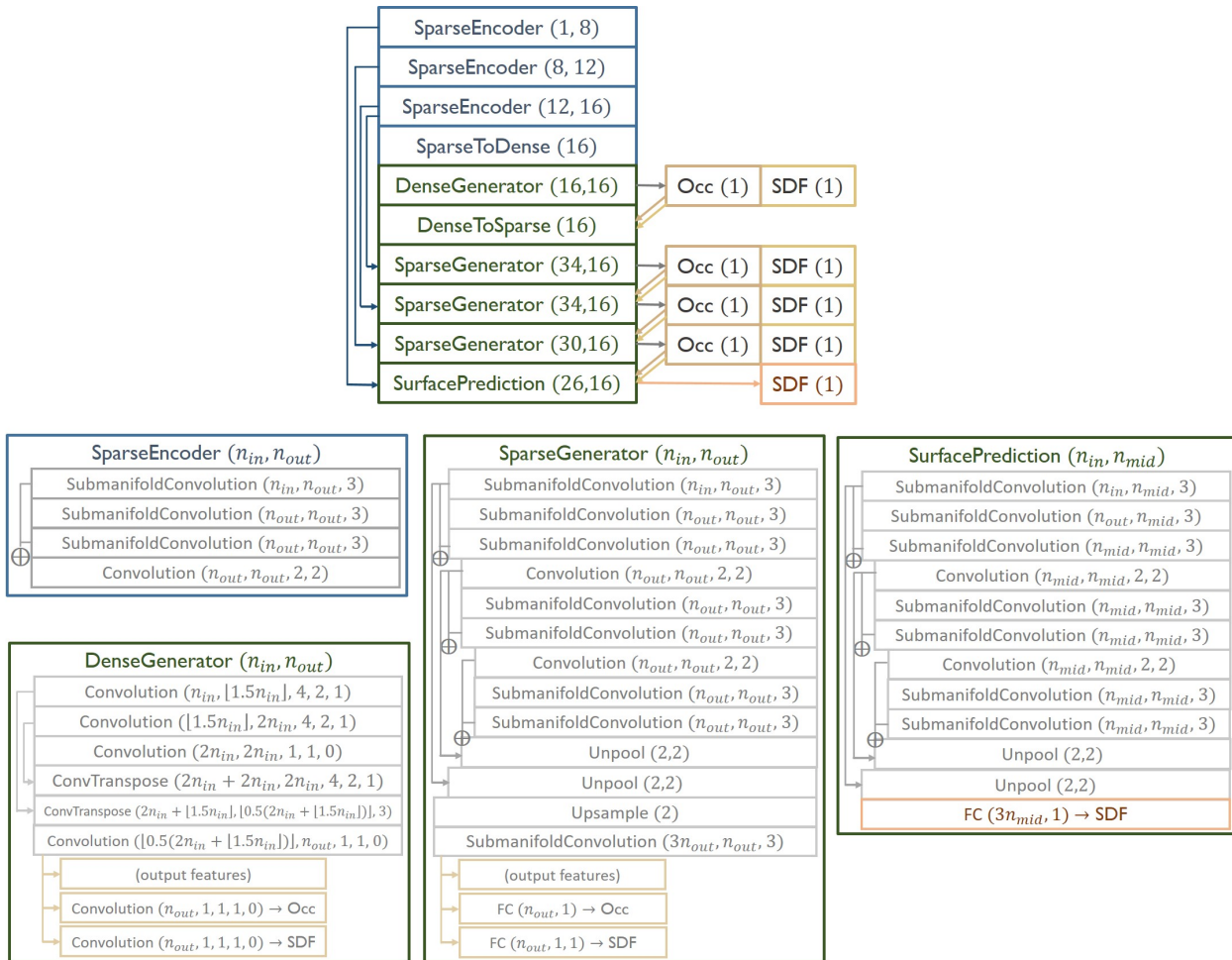


Figure 8: SG-NN architecture in detail. The final TSDF values are highlighted in orange, and intermediate outputs in yellow. Convolution parameters are given as (nf\_in, nf\_out, kernel\_size, stride, padding), with stride and padding default to 1 and 0. Arrows denote concatenation, and  $\oplus$  denotes addition.



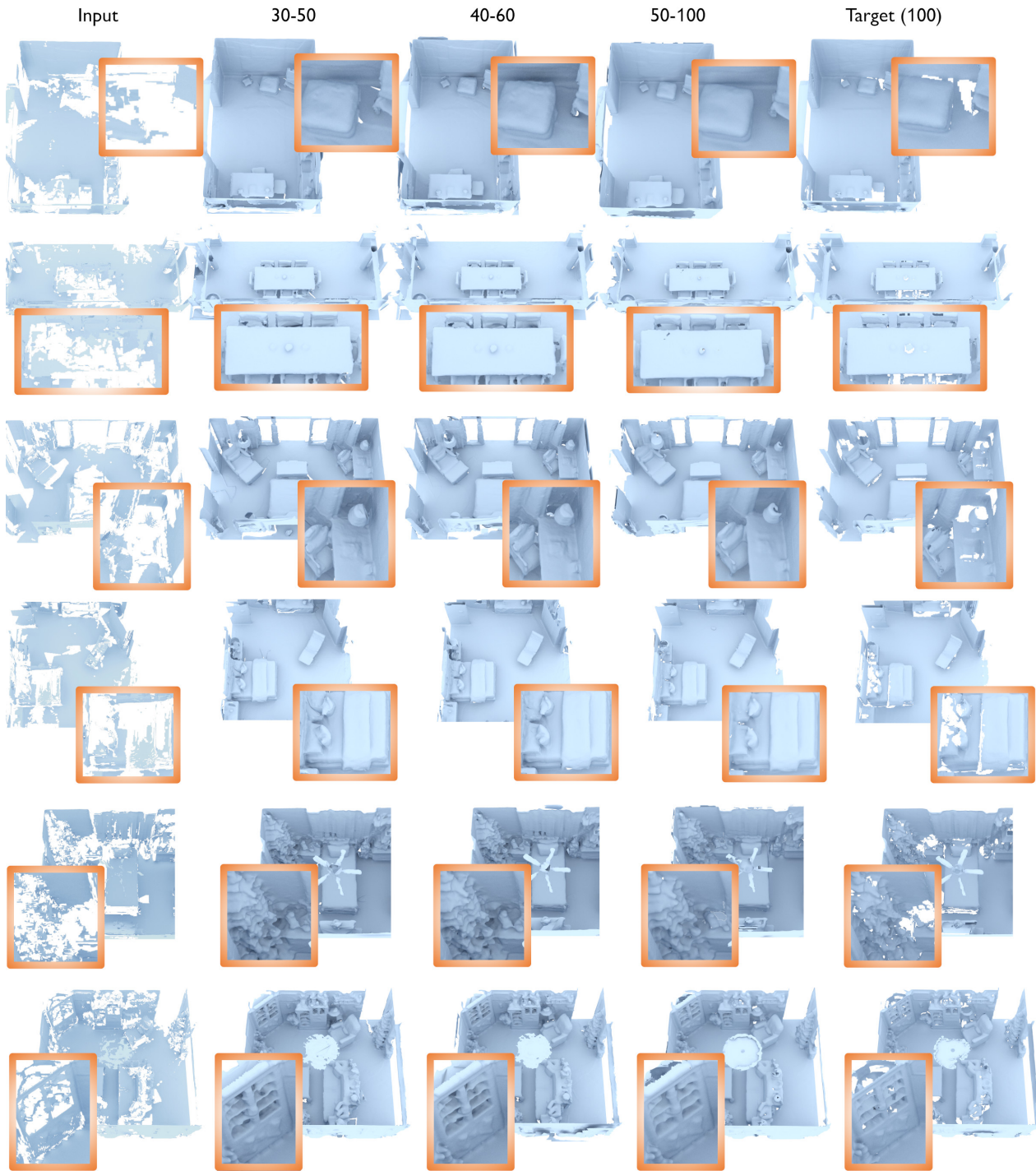


Figure 9: Qualitative evaluation of varying target data completeness available for training. We generate various incomplete versions of the Matterport3D [3] scans using  $\approx 30\%$ ,  $40\%$ ,  $50\%$ ,  $60\%$ , and  $100\%$  of the frames associated with each room scene, and evaluate on the  $50\%$  incomplete scans. Even as the level of completeness of the target data used during training decreases significantly, our self-supervised approach effectively learns the geometric completion process, maintaining robustness in generating complete geometry.