

UnScene3D: Unsupervised 3D Instance Segmentation for Indoor Scenes

David Rozenberszki¹ Or Litany^{2,3} Angela Dai¹

¹Technical University of Munich ²Technion ³NVIDIA

<https://rozdavid.github.io/unscene3d>



Figure 1. We propose UnScene3D, a fully-unsupervised 3D instance segmentation method, effectively separating semantic instances without requiring any manual annotations. We utilize geometric primitives to ensure crisp masks, and due to our self-training loop, we can also obtain a dense set of predictions, even in cluttered indoor scenarios.

Abstract

3D instance segmentation is fundamental to geometric understanding of the world around us. Existing methods for instance segmentation of 3D scenes rely on supervision from expensive, manual 3D annotations. We propose UnScene3D, the first fully unsupervised 3D learning approach for class-agnostic 3D instance segmentation of indoor scans. UnScene3D first generates pseudo masks by leveraging self-supervised color and geometry features to find potential object regions. We operate on a basis of geometric oversegmentation, enabling efficient representation and learning on high-resolution 3D data. The coarse proposals are then refined through self-training our model on its predictions. Our approach improves over clustering-based alternatives to unsupervised 3D instance segmentation methods by more than 300% Average Precision score, demonstrating effective instance segmentation even in challenging, cluttered 3D scenes.

1. Introduction

The increasing availability of commodity RGB-D sensors, now widely available on iPhones as well as with the Microsoft Kinect or Intel RealSense, has enabled consumer-

level capture of 3D geometry of real-world environments. To enable applications in robotics, autonomous navigation, and mixed reality in such scenes, semantic 3D scene understanding is necessary. In particular, 3D instance segmentation is critical to 3D perception, providing dense instance mask predictions, thus enabling physical and geometric reasoning about objects in an environment. While various 3D deep learning approaches have been developed for 3D instance segmentation [5, 14, 17, 18, 21, 22, 30, 32, 42–45, 50–52, 54, 57, 58], they require full supervision from expensive, manual, dense annotations on 3D scenes.

We introduce UnScene3D, a novel approach designed for class-agnostic 3D instance segmentation. Our aim is to identify objects in real-world 3D scans by predicting their dense instance masks, without any constraints to a predefined set of class categories. Moreover, we avoid expensive data annotation requirements by operating in an unsupervised fashion, instead leveraging self-supervised 2D and 3D features for segmentation.

UnScene3D comprises two essential elements. First, we observe that for RGB-D scan data, self-supervised representation learning methods [19, 60] can provide an innate signal indicating object-ness through feature similarity. We thus generate pseudo masks over 3D segment primitives, based on multimodal analysis of self-supervised color and geometry features from the RGB-D data. By considering

mesh segments rather than voxels or points, our approach efficiently scales with high-resolution 3D data in large scene environments while inherently promoting contiguous segmentation masks. As we require strong features for these initial coarse estimates, we fuse information from both geometric and 2D color features in a complementary fashion. Second, following the pseudo mask generation, we train our model through iterative self-training on both the initial pseudo masks and the current confident model predictions. Through multiple rounds of self-training with noise robust losses achieve improved object recognition and segmentation. At inference time, we do not require any 2D color signal and can produce class-agnostic 3D instance segmentation for a new geometric observation of a 3D environment. Experiments on challenging, cluttered indoor environments from the ScanNet [10], S3DIS [1] and ARKit [2] datasets show that UnScene3D improves significantly over unsupervised, clustering-based state of the art. In summary, our contributions are:

- We propose an unsupervised 3D instance segmentation approach for indoor RGB-D scans, without requiring any human annotation.
- We generate sparse 3D pseudo masks for unsupervised training based on a multi-modal fusion of color and geometric signal from RGB-D scan data. We achieve robustness and efficiency through a geometry-aware scene coarsening.
- Our generated pseudo masks are iteratively refined by self-training for 3D instances to improve 3D instance segmentation performance.

2. Related Work

Self-supervised 3D pretraining While significant progress has been made in fully supervised 3D instance segmentation [8, 14, 16–18, 20, 42, 44, 50, 51] the amount of densely annotated 3D data is scarce. Inspired by success in the 2D domain, various 3D pretraining methods have been developed to boost semantic and instance segmentation performance when fine-tuning with annotated semantic labels. Such methods leverage instance discrimination based on different camera views [19, 60], local augmentations [62], or multiple LIDAR sweeps [39]. While these methods can provide powerful 3D feature extraction, they do not construct any notion of object instances.

Weakly-supervised 3D segmentation Classical methods have leveraged object template information to match or retrieve templates to local geometry in a scene [4, 25, 28, 31, 36, 37], thereby identifying potential object locations. Other methods formulated 3D dense instance segmentation with only 3D box annotation [6, 41] or single-point supervision and active-learning [34, 53]. More recent methods have focused on exploiting knowledge from powerful pre-trained

vision-language models to inform text-guided queries in 3D scenes [12, 24, 33, 40, 46]; however, such methods still rely on large-scale annotated data in the 2D domain.

Clustering-based segmentation There has been very little work done in fully unsupervised 3D instance segmentation, but classical clustering methods have been used to group regions with similar geometric properties together. A particularly notable approach is the density-based clustering of DBSCAN [13] and its hierarchical counterpart HDBSCAN [35]. These methods can be used to group point clusters in a 3D scene based on point normals and colors. The ScanNet dataset [10] showed that the Felzenswalb algorithm [15] originally developed for image over-segmentation, can generate useful geometric segment clusters. We also exploit such geometric primitives to guide dimensionality reduction and feature aggregation.

Finally, recent methods have been developed to detect instances with self-supervised pretrained features in driving scenarios. These methods often leverage the unique properties of such data including dynamics and instance sparsity. Song et. al. [48] identify object instances through motion, showing promise for self-driving scenarios, but limited to moving objects. Nunes et. al. [38] additionally propose a clustering and graph cut based refinement on pre-trained 3D features, focusing on sparse outdoor scenarios to identify spatially separate objects. Our solution aims to segment instances in complex, cluttered indoor environments.

Unsupervised 2D instance segmentation Classical graph-cut algorithms [7, 11, 47, 59] can be used to detect objects in scenes, employing low-level feature clustering to identify self-similar regions. Recent advances in self-supervised feature learning have been employed in 2D unsupervised instance segmentation methods, which use two-stage training pipelines to achieve remarkable segmentation results [55, 56]. These methods first generate a set of coarse pseudo masks building on the insights of graph-cut algorithms and then refine them with a series of self-training iterations. In particular, FreeSolo [55] uses multi-branch feature extraction to obtain self-similar regions as mask proposals, producing a dense set of initial pseudo-annotated instances. CutLER [56] uses the normalized cut (NCut) algorithm [47] with deep self-supervised features from DINO [3] to identify multiple prominent regions as pseudo masks. Inspired by such approaches we also leverage pseudo mask generation and self-training, but to handle high-dimensional, noisy real-world 3D scan data, we employ a multi-modal feature reasoning and geometric graph coarsening for robust unsupervised 3D instance segmentation.

3. Method

Problem definition We propose an unsupervised learning-based method for 3D instance segmentation. We operate on a set of training 3D scenes $\{X_i\}_{i=1}^{n_i}$, represented as mesh graphs $G = (V, E)$, of vertices V and triangular face edges E , where each scene X_i contains an unknown set of n_i objects in the i^{th} scene. We aim to train a model that can predict for a previously unseen input scene X , a set of 3D masks representing the different object instances in that scene.

Method overview In order to achieve unsupervised 3D instance segmentation we first break down the scenes into N geometric primitives S_N , which we use to initialize an adjacency matrix W to extract an initial set of pseudo masks M^0 , representing instance hypotheses based on combining 2D and 3D inputs $\mathcal{F}_{2D} / \mathcal{F}_{3D} \in R^{N \times D_{2D}/3D}$, where D_{2D}, D_{3D} are the dimensions of the 2D/3D self-supervised features. We regularize the per-segment similarities over geometric primitives for mitigating noise and enabling efficient 3D reasoning. We then employ a series of self-training cycles, updating pseudo mask supervision with new predicted masks, in order to produce final 3D instances. An overview of our approach is shown in Figure 2.

3.1. Initial pseudo mask generation

In order to initiate self-training, we first generate an initial set of pseudo masks, leveraging complementary information from 2D and 3D signal in $\{X_i\}$.

3.1.1 Feature aggregation

To encourage effective initial pseudo mask generation, we employ joint reasoning across both self-supervised color and geometry features, as they can provide complementary information regarding objects. As RGB-D scans often contain color image information and reconstructed 3D scan geometry, we can associate both 2D and 3D features in 3D by back-projecting the 2D extracted features using the corresponding depth and camera pose information for each image. Both 2D and 3D features are extracted through state-of-the-art self-supervised feature learning methods [3, 19]. As real-world camera estimation often contains small misalignment errors and noise or oversmoothing in reconstructed scan geometry, these self-supervised features can often also contain high-frequency noise, which we address in Sec. 3.1.2 when reasoning over these features. Note that while we employ both 2D and 3D signal when available for training, we do not require any aligned color image inputs for inference, enabling more general applicability.

3.1.2 3D Graph Cut

To generate pseudo masks from the 2D and 3D self-supervised features, we employ graph cut to estimate class-agnostic instances from the background. More precisely, we leverage the principle of Normalized Cut [47] (NCut), which employs eigenvalue decomposition from an adjacency matrix $W \in R^{N \times N}$ over a graph to identify self-similar regions potentially representing semantic instances, where a set of potential instances can be extracted iteratively. Given a graph representing the 3D scene, we build an adjacency matrix W and self-supervised features with a corresponding degree matrix $D \in R^{N \times N}$, where $D(i, i) = \sum_j W(i, j)$ and $(D - W)v = \lambda Dv$. In this system, finding the second smallest eigenvalue λ and its corresponding eigenvector v is a close approximation for the minimized cost. From v , we obtain foreground separation by taking all node activations where the eigenvector components were larger than their mean. To identify multiple foreground objects, this process is repeated iteratively.

Unfortunately, applying this approach directly to the 3D scenes $\{X_i\}$ in common 3D representations such as voxels or points is not only computationally infeasible, but unreliable due to the noise in camera pose estimation and geometric reconstruction of 3D scan data. Thus, we propose to regularize the graph cut across geometric primitives.

3.1.3 Geometric Primitives

To employ efficient reasoning across high-dimensional 3D data and enable robust 3D regularization of noisy features, we propose to operate on geometric primitives acquired through a graph coarsening process. For a 3D scene X_i we construct the graph $G = (V, E)$ where V and E being the mesh vertices and face edges. Then, nodes with similar normals and colors are aggregated and clustered based on the mesh topology following [15] and resulting in a set $S_N = \{C_1 \dots C_N\}$ and $\bigcup(S_N) = V$ where C_n represent a single primitive. This reduces the graph size by multiple orders of magnitude, and enables effective regularization of noise in the used self-supervised 2D and 3D features.

3.1.4 NCut on Geometric Primitives

After addressing the challenge of dimensionality reduction and effectively mitigating speckle noise in our features using geometric primitives, we can leverage the capabilities of the Normalized Cut algorithm to achieve a clean partitioning of scene graphs. For this, we iteratively apply NCut to our aggregated features for the extraction of initial pseudo masks denoted as M . Starting with an empty set $M^0 = \{\}$, we iteratively compute the adjacency matrix over S_N and retrieve the masks $m \subset S_N$. We start

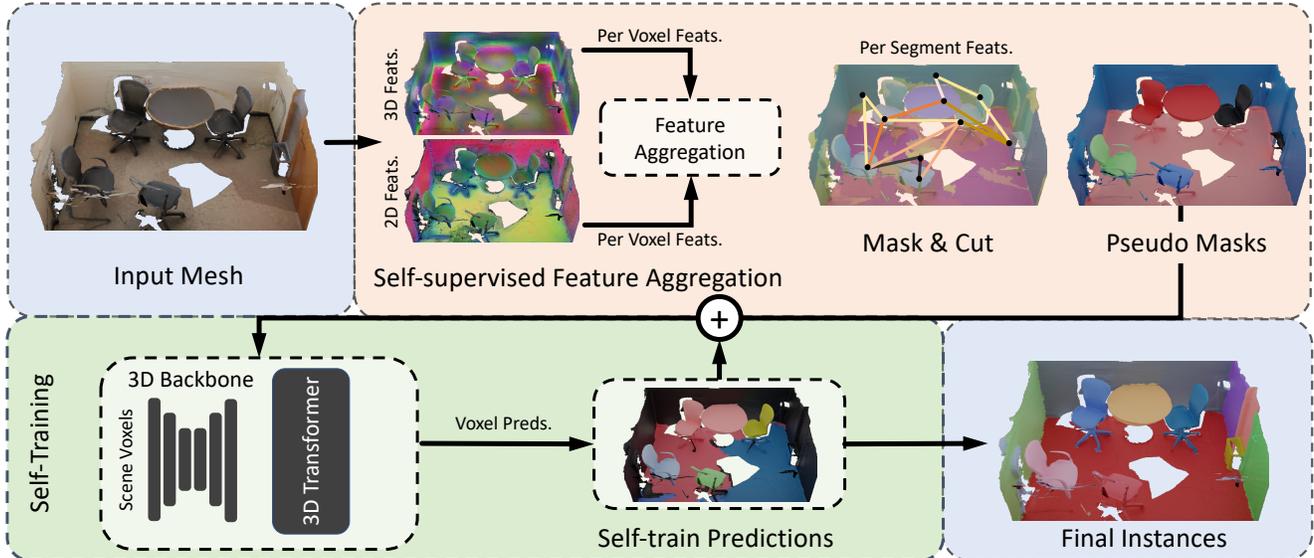


Figure 2. UnScene3D first generates a set of pseudo masks (top) to initiate self-training (bottom) for unsupervised 3D instance segmentation. We leverage features from 3D self-supervised pre-training in combination with 2D self-supervised features on an input mesh. These multi-modal features are then aggregated on geometric primitives, integrating low- and high-level signals for pseudo mask segmentation. These initial pseudo masks are then used as supervision for a 3D transformer-based model to produce updated instance masks that are integrated into the supervision of multiple self-training cycles. Finally, we obtain clean and dense instance segmentation without using any manual annotations.

from N geometric segments with their corresponding D -dimensional features $\mathcal{F} \in \mathcal{R}^{N \times D}$, and construct the similarity matrix $A = \text{sim}(\mathcal{F})$, where sim denotes cosine similarity. Additionally, for the multi-modal setup we calculate similarity matrices A_{2D} and A_{3D} independently and take their weighted average to obtain the final scores. Empirically, we found this to be more robust than direct feature fusion of the different modalities, due to their different statistical characteristics. We obtain W_j introduced in Section 3.1.2 by thresholding A at τ_{cut} , where j denotes the j^{th} NCut iteration. Using W_j , we solve for the second eigenvector v_j and threshold it to retrieve the partition m_j . We keep all separated foregrounds in M^0 , where for each upcoming iteration, we mask out the row and column vectors from W_j , where $m_i \in M^0$ was already accepted as a foreground instance and i being the previous segment ids. This allows greedy separation of instances in order of confidence in every cut iteration. Examples of our generated pseudo masks are visualized in Figures 5 and 6.

As the adjacency graph is unaware of the mesh connectivity, NCut often results in masks that span spatially separated scene regions. In 3D, we can leverage knowledge of physical distance and connectivity of G to constrain masks to be contiguous in the coarsened scene connectivity graph. We thus filter masks m_j that have separated components, keeping only the parts \tilde{m}_j that contain the item with the maximum absolute value in v_j . Separation based on connectivity is performed before saving \tilde{m}_j into M^0 , thus allowing for repeated detection of the dropped part over the

next NCut iterations. Finally, we iterate until the maximum number of instances $M^0 = \{m_i\}_{i=1}^{N_m}$ are obtained, or there are no segments left in the scene. Moreover, we favor generating a reliable set of masks at the cost of restricting to a sparse initial set (i.e., missing potential instances rather than generating noisy masks for them) through a stricter τ_{cut} or lower number of instances.

3.2. Self-Training

Our initial pseudo masks can provide a set of proposed instances M^0 ; however, these pseudo masks are quite sparse in the scenes and sometimes over- or under-split nearby instances. We thus refine the pseudo mask data through an iterative self-training strategy, producing final instance segmentation predictions M^l with more dense and complete instance proposals.

We leverage a state-of-the-art 3D transformer-based backbone [45] for our self-training from pseudo mask data as mask-head supervision, while the class-head is collapsed to *foreground* and *background* classes. Through multiple training cycles we save the proposals of the t^{th} iteration into M^t , from the self-trained model, and save these masks as an extension to the original pseudo dataset obtaining $M^t \supseteq M^0$. From the second training iteration, we can extract the most confident K predictions and sample these new instance proposals as an addition to the pseudo annotations. Further, we only accept new instances if the added information value is larger than the minimum threshold, measured by simple segment IoU scores. This way, we can effectively

densify the originally sparse annotations, but without limiting the quality of the originally clean pseudo masks.

3.3. Implementation Details

Backbones. We use a Res16UNet34C sparse-voxel UNet implemented in the MinkowskiEngine [8] for 3D pre-trained feature extraction as well as for the 3D transformer during self-training. For the pretrained features we use our own trained weights of [19] for compatibility reasons.

Self-training. We employ the 3D transformer architecture of [45], initialized from scratch. The first self-training cycle is trained for 600 epochs with a batch size of 8 until convergence, which takes ≈ 3 days on a single NVIDIA RTX A6000 GPU. Further self-training cycles are all initialized from the previous state and finetuned for an additional 50 epochs in ≈ 4 hours and for a total of 4 training cycles to produce the final set of instance predictions S . For the Hungarian assignment, we take the original weighted combination of dice and binary cross-entropy losses and only apply the DropLoss condition in the backpropagation phase.

4. Experiments

We demonstrate the effectiveness of UnScene3D for unsupervised class-agnostic 3D instance segmentation on challenging real-world 3D scan datasets containing a large diversity of objects and significant clutter. We train our method and all learned baselines on ScanNet [10], using the official train split. Note that no semantic annotation data is used for training, only the RGB-D reconstructions. Additionally, we show that our approach trained on ScanNet data can effectively transfer to class-agnostic 3D instance segmentation on ARKitScenes [2] data.

Datasets. We train and evaluate UnScene3D on RGB-D scan data from ScanNet [10], using the official train split. We use the raw RGB images, and registered camera poses for training our approach, while the semantic annotations are used only for evaluation. We use the official ScanNet train split for both the pre-trained 3D features from [19] and our self-training iterations. We additionally evaluate our method on ARKitScenes [2], on an 884/120 train/test split of indoor LIDAR scans. For ARKitScenes, we use 3D pre-trained features from ScanNet, followed by pseudo mask generation and self-training on the ARKitScenes train scenes. We convert the LIDAR scan data to meshes with Poisson Surface Reconstruction [26, 27] prior to our graph coarsening. Note that all baselines using learned features are trained on the same ScanNet data as ours.

Evaluation metrics. We evaluate class-agnostic 3D instance segmentation performance with the widely-used Average Precision score on the full-resolution mesh vertices.

<i>ScanNet</i>	AP@25	AP@50	AP
HDBSCAN [35]	32.1	5.5	1.6
Nunes et al. [38]	30.5	7.3	2.3
Felzenszwalb [15]	38.9	12.7	5.0
CutLER Projection [56]	7.0	0.2	0.3
Ours	58.5	32.2	15.9

Table 1. Unsupervised class-agnostic 3D instance segmentation on ScanNet [10]. Our approach improves significantly over baselines (3x improvement in AP) due to our pseudo mask generation and self-training strategy.

Following the strategy of the supervised benchmark [10] we report scores at IoU scores of 25% and 50% (AP@25, AP@50) and averaged over all overlaps between [50% and 95%] at 5% steps (AP). Note that since predictions are class agnostic, all methods evaluate only instance mask AP values without considering any semantic class labels. For ScanNet, we evaluate against ground truth instance masks from the established 20-class benchmark. Since ARKitScenes does not contain any ground truth instance mask annotations, we evaluate all methods qualitatively.

Comparison to the state of the art. We evaluate our approach in comparison to state-of-the-art traditional clustering methods HDBSCAN [35] and Felzenszwalb’s algorithm [15], in addition to the unsupervised approach of Nunes et. al. [38] leveraging learned feature clustering and refinement. All baselines are provided with input mesh vertices, colors, and normals, while our approach and Nunes et. al. also operate on sparse voxel scene representations. Table 1 and Figure 3 show comparisons on ScanNet data; our UnScene3D approach improves significantly over state of the art by effectively leveraging signal from self-supervised 3D features to guide our model through self-training. Note that since Nunes et. al. has been designed for outdoor applications, even while leveraging ScanNet-trained features, it uses ground removal and relies on physical object separation, making segmentation difficult in cluttered scenes.

Additionally, we demonstrate the importance of reasoning in 3D, and compare with a state-of-the-art unsupervised 2D instance segmentation approach CutLER [56] run on the RGB frames of the scans, and projected to 3D using the corresponding camera poses. Here, the difficulty lies in resolving view inconsistencies, occlusions, and lack of knowledge of geometric structure resulting in poor 3D segmentation performance despite plausible 2D proposals.

Evaluation on other datasets We quantitatively evaluate UnScene3D on the Area_5 of the S3DIS dataset [1] using only 3D features pretrained on [10]. Comparison with 3D-only state-of-the-art can be seen in Table 2.

We additionally compare with state of the art on ARKitScenes [2] data in Figure 7. Here we show only qualitative

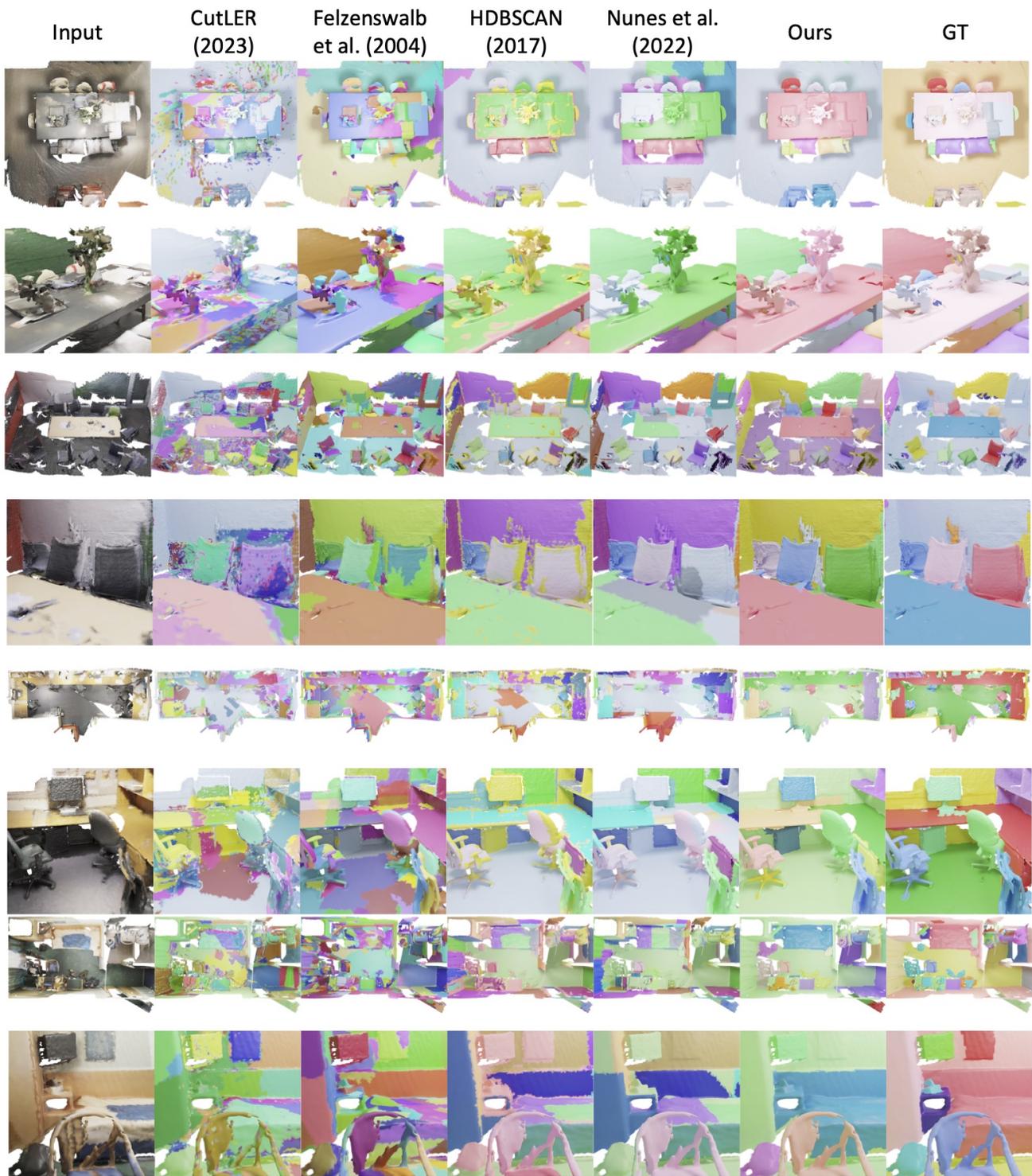


Figure 3. Qualitative comparison on ScanNet [10] scenes with projected predictions from the 2D method CutLER [56], traditional clustering-based methods Felzenszwalb [15] and HDBSCAN [35], and the GraphCut-based cluster refinement method [38]. Our approach leverages strong pseudo mask prediction and a self-training strategy to produce cleaner, more accurate instance segmentation.

<i>S3DIS</i>	AP@25	AP@50	AP
HDBSCAN [35]	27.9	11.2	5.0
Felzenswalb [15]	23.5	10.7	5.0
Nunes et al. [38]	20.1	10.5	5.5
Ours	52.6	40.3	21.4

Table 2. Evaluation on S3DIS dataset (Area_5). UnScene3D is able to adapt to other datasets as well and shows a significant improvement over previous SOTA methods.

results due to the absence of ground truth instance mask annotations. UnScene3D effectively produces cleaner, more accurate segmentations in these complex environments.

UnScene3D as data-efficient pretraining UnScene3D is able to learn powerful object properties and dense segmentation even in a fully unsupervised fashion. We demonstrate the potential of our strong learned features for downstream 3D instance segmentation with limited annotated data. We follow the setup introduced by CSC [19] with limited reconstructions available for downstream fine-tuning. We show our method as a strong pretraining strategy in Figure 4, notably outperforming both training from scratch as well as the state-of-the-art 3D pretraining of CSC. For more details we refer to our supplementary material.

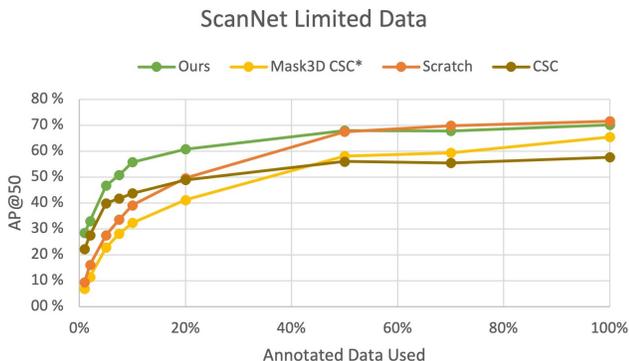


Figure 4. Our unsupervised self-training produces strong 3D features that can serve as a powerful pretraining strategy for 3D instance segmentation in limited data scenarios. UnScene3D significantly outperforms state-of-the-art self-supervised 3D pretraining [19] on ScanNet instance segmentation.

What is the effect of multi-modal signal for pseudo mask generation? We evaluate the effect self-supervised color and geometry signals for generating pseudo annotations in Table 3. We consider using only self-supervised geometric features (3D), only self-supervised color features (2D) that are projected to the 3D scans, and both together (both). We find that the color and geometry provide complementary signals. We also note that color features are only used for the initial pseudo mask generation, during self-training iterations and test time only 3D features were used.

	Modality	AP@25	AP@50	AP	AP Final
FreeMask	3D	14.4	3.6	1.3	2.0
Ours	3D	45.4	16.7	9.2	13.3
FreeMask	2D	31.1	15.1	6.8	13.8
Ours	2D	51.3	21.8	9.4	15.7
FreeMask	both	23.7	10.1	5.7	12.1
Ours	both	52.9	23.2	10.4	15.9

Table 3. We compare pseudo mask generation from 3D-only features (3D), color-only features (2D), and both color and geometry (both) signal, as well as with pseudo annotation generation algorithm FreeMask [55]. In this table we report method performances after a single iteration of self-training initialized from the different pseudo annotation methods and the final AP scores after 4 self-training iterations.



Figure 5. Initial pseudo masks generated by UnScene3D in comparison with a 3D-lifted FreeMask [55]. FreeMask tends to produce a larger set of noisier pseudo masks, while we rely on a cleaner but sparser set for our self-training.

What is the effect of pseudo annotations? We also evaluate the effect of our pseudo mask generation in Table 3 and Figure 5, in comparison to the 3D adaptation of the FreeMask [55] approach operating on our geometric segments. FreeMask tends to estimate a larger but noisier set of initial pseudo masks, while our approach is focusing on a sparser set of more reliable pseudo masks and produces significantly better performance. The strong difference in performance can be explained by the nature of the samples. While a sparser set of examples can be extended with multiple iterations of self-training, noisy samples will propagate through the full pipeline, and thus directly degrade the final performance. Further details of our adaptations of the FreeMask 3D method can be found in our supplemental.

What is the impact of self-training? We observe that while self-training iterations are always improving the qualitative performance, their effective added information value is saturating after a limited number of cycles. We report on Table 4 through the first 4 steps, and observe a significant relative improvement in both modalities.



Figure 6. UnScene3D employs self-training to refine the initial sparse set of proposals. We can see consistent improvement over both the number of predicted instances and the quality of the instance masks. Here we show results using the pseudo annotations obtained from both modalities.

	3D Only			3D & 2D		
	AP@25	AP@50	AP	AP@25	AP@50	AP
S^0 pseudo masks	13.8	4.7	2	19.9	10.0	5.9
1 st Self-train	45.4	16.7	9.2	52.9	23.2	10.4
2 nd Self-train	50.0	24.1	12.0	56.5	29.8	15.0
3 rd Self-train	52.2	25.8	12.8	58.8	31.9	15.9
4 st Self-train	52.7	26.2	13.3	58.5	32.2	15.9

Table 4. Multiple iterations of self-training significantly improve performance, saturating around 4 iterations.

Limitations While UnScene3D offers a promising step towards unsupervised 3D instance segmentation, various limitations remain. We rely on a mesh representation for graph coarsening, but believe this could be extended to alternative representations through neighborhood reasoning. Additionally, our graph coarsening step may cause very small objects (e.g., pens, cell phones) to be missed in the pseudo annotation generation. Finally, employing a fixed set of pseudo masks from the initial stage that are used

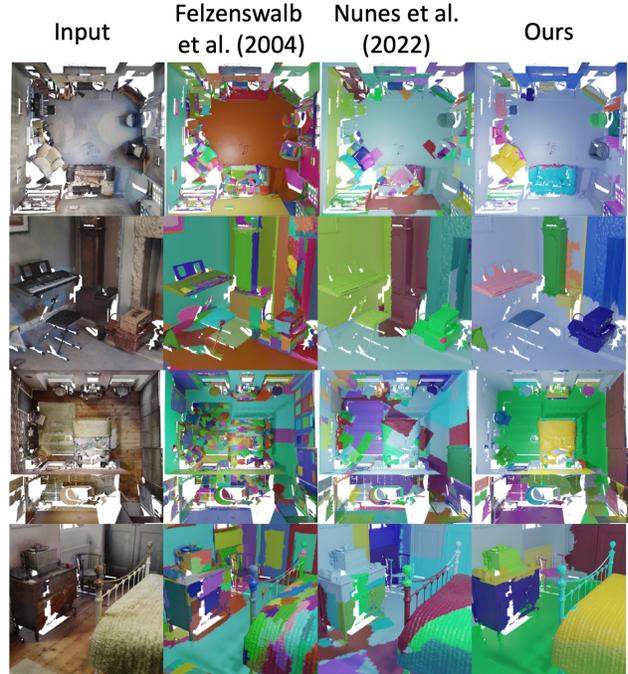


Figure 7. As UnScene3D does not require any human annotation, so we can also train and test our method on the ARKitScenes [2] dataset. We leverage 3D features followed by a series of self-training iterations for cleaner, more accurate instance segmentation. Qualitative results show consistently better results than our baselines.

through self-training could reinforce noisy predictions.

5. Conclusion

We introduced UnScene3D, a novel approach towards achieving fully-unsupervised 3D instance segmentation in cluttered indoor scenes. Our approach effectively combined low-level geometric properties to regularize multi-modal self-supervised deep features for initial pseudo mask extraction, and our self-training notably improved performance by refining these proposals to a more complete, dense set of instances. As 3D instance segmentation is a crucial aspect of 3D scene understanding, UnScene3D’s ability to achieve this without requiring any manual annotations opens up new possibilities for 3D semantic understanding.

6. Acknowledgements

This project is funded by the Bavarian State Ministry of Science and the Arts and coordinated by the Bavarian Research Institute for Digital Transformation (bidt), the ERC Starting Grant SpatialSem (101076253), and supported in part by a Google research gift. Or Litany is a Taub fellow and is supported by the Azrieli Foundation Early Career Faculty Fellowship.

References

- [1] Iro Armeni, Ozan Sener, Amir R. Zamir, Helen Jiang, Ioannis Brilakis, Martin Fischer, and Silvio Savarese. 3d semantic parsing of large-scale indoor spaces. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1534–1543, 2016. 2, 5
- [2] Gilad Baruch, Zhuoyuan Chen, Afshin Dehghan, Tal Dimry, Yuri Feigin, Peter Fu, Thomas Gebauer, Brandon Joffe, Daniel Kurz, Arik Schwartz, and Elad Shulman. ARK-scenes - a diverse real-world dataset for 3d indoor scene understanding using mobile RGB-d data. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*, 2021. 2, 5, 8, 12, 13
- [3] Mathilde Caron, Hugo Touvron, Ishan Misra, Hervé Jégou, Julien Mairal, Piotr Bojanowski, and Armand Joulin. Emerging properties in self-supervised vision transformers. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2021. 2, 3, 16
- [4] Kang Chen, Yu-Kun Lai, Yu-Xin Wu, Ralph Martin, and Shi-Min Hu. Automatic semantic modeling of indoor scenes from low-quality rgb-d data using contextual information. *ACM Transactions on Graphics*, 33(6), 2014. 2
- [5] Shaoyu Chen, Jiemin Fang, Qian Zhang, Wenyu Liu, and Xinggang Wang. Hierarchical aggregation for 3d instance segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 15467–15476, 2021. 1
- [6] Julian Chibane, Francis Engelmann, Tuan Anh Tran, and Gerard Pons-Moll. Box2mask: Weakly supervised 3d semantic instance segmentation using bounding boxes. In *European Conference on Computer Vision (ECCV)*. Springer, 2022. 2
- [7] Sunil Chopra and M. R. Rao. The partition problem. *Mathematical Programming*, 59:87–115, 1993. 2
- [8] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4d spatio-temporal convnets: Minkowski convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3075–3084, 2019. 2, 5
- [9] Angela Dai and Matthias Nießner. 3dmv: Joint 3d-multi-view prediction for 3d semantic scene segmentation. In *European Conference on Computer Vision*, 2018. 12
- [10] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proc. Computer Vision and Pattern Recognition (CVPR), IEEE*, 2017. 2, 5, 6, 12, 14
- [11] Michel Deza and Monique Laurent. Geometry of cuts and metrics. In *Algorithms and Combinatorics*, 2009. 2
- [12] Runyu Ding, Jihan Yang, Chuhui Xue, Wenqing Zhang, Song Bai, and Xiaojuan Qi. Pla: Language-driven open-vocabulary 3d scene understanding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023. 2
- [13] Martin Ester, Hans-Peter Kriegel, Jörg Sander, and Xiaowei Xu. A density-based algorithm for discovering clusters in large spatial databases with noise. In *Proceedings of the Second International Conference on Knowledge Discovery and Data Mining*, page 226–231. AAAI Press, 1996. 2
- [14] Siqi Fan, Qiulei Dong, Fenghua Zhu, Yisheng Lv, Peijun Ye, and Fei-Yue Wang. Scf-net: Learning spatial contextual features for large-scale point cloud segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14504–14513, 2021. 1, 2
- [15] Pedro F Felzenszwalb and Daniel P Huttenlocher. Efficient graph-based image segmentation. *International journal of computer vision*, 59:167–181, 2004. 2, 3, 5, 6, 7
- [16] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. *CVPR*, 2018. 2
- [17] Lei Han, Tian Zheng, Lan Xu, and Lu Fang. Ocuseg: Occupancy-aware 3d instance segmentation. In *2020 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2937–2946, 2020. 1
- [18] Ji Hou, Angela Dai, and Matthias Nießner. 3d-sis: 3d semantic instance segmentation of rgb-d scans. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 4421–4430, 2019. 1, 2, 12
- [19] Ji Hou, Benjamin Graham, Matthias Nießner, and Saining Xie. Exploring data-efficient 3d scene understanding with contrastive scene contexts. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15587–15597, 2021. 1, 2, 3, 5, 7, 12, 15, 16
- [20] Ji Hou, Xiaoliang Dai, Zijian He, Angela Dai, and Matthias Nießner. Mask3d: Pre-training 2d vision transformers by learning masked 3d priors. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 13510–13519, 2023. 2, 12
- [21] Qingyong Hu, Bo Yang, Linhai Xie, Stefano Rosa, Yulan Guo, Zhihua Wang, Niki Trigoni, and Andrew Markham. Randa-net: Efficient semantic segmentation of large-scale point clouds. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11108–11117, 2020. 1
- [22] Le Hui, Linghua Tang, Yaqi Shen, Jin Xie, and Jian Yang. Learning superpoint graph cut for 3d instance segmentation. In *NeurIPS*, 2022. 1
- [23] Maximilian Jaritz, Jiayuan Gu, and Hao Su. Multi-view pointnet for 3d scene understanding. In *ICCV Workshop 2019*, 2019. 12
- [24] Krishna Murthy Jatavallabhula, Alihusein Kuwajerwala, Qiao Gu, Mohd Omama, Tao Chen, Shuang Li, Ganesh Iyer, Soroush Saryazdi, Nikhil Keetha, Ayush Tewari, Joshua B. Tenenbaum, Celso Miguel de Melo, Madhava Krishna, Liam Paull, Florian Shkurti, and Antonio Torralba. Conceptfusion: Open-set multimodal 3d mapping. *arXiv*, 2023. 2
- [25] Andrej Karpathy, Stephen Miller, and Li Fei-Fei. Object discovery in 3d scenes via shape analysis. In *2013 IEEE international conference on robotics and automation*, pages 2088–2095. IEEE, 2013. 2
- [26] Michael Kazhdan and Hugues Hoppe. Screened poisson surface reconstruction. *ACM Transactions on Graphics (ToG)*, 32(3):1–13, 2013. 5

- [27] Michael Kazhdan, Matthew Bolitho, and Hugues Hoppe. Poisson surface reconstruction. In *Proceedings of the fourth Eurographics symposium on Geometry processing*, page 0, 2006. 5
- [28] Young Min Kim, Niloy J Mitra, Dong-Ming Yan, and Leonidas Guibas. Acquiring 3d indoor environments with variability and repetition. *ACM Transactions on Graphics (TOG)*, 31(6):1–11, 2012. 2
- [29] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *arXiv preprint arXiv:2304.02643*, 2023. 15
- [30] Maksim Kolodiazhnyi, Danila Rukhovich, Anna Vorontsova, and Anton Konushin. Top-down beats bottom-up in 3d instance segmentation, 2023. 1
- [31] Yangyan Li, Angela Dai, Leonidas Guibas, and Matthias Nießner. Database-assisted object retrieval for real-time 3d reconstruction. In *Computer graphics forum*, pages 435–446. Wiley Online Library, 2015. 2
- [32] Zhihao Liang, Zhihao Li, Songcen Xu, Mingkui Tan, and Kui Jia. Instance segmentation in 3d scenes using semantic superpoint tree networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 2783–2792, 2021. 1
- [33] Minghua Liu, Yin hao Zhu, H. Cai, Shizhong Han, Z. Ling, Fatih Murat Porikli, and Hao Su. Partslip: Low-shot part segmentation for 3d point clouds via pretrained image-language models. *2023 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21736–21746, 2022. 2
- [34] Zhengzhe Liu, Xiaojuan Qi, and Chi-Wing Fu. One thing one click: A self-training approach for weakly supervised 3d semantic segmentation. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 1726–1736, 2021. 2
- [35] Leland McInnes and John Healy. Accelerated hierarchical density based clustering. In *2017 IEEE International Conference on Data Mining Workshops (ICDMW)*, pages 33–42. IEEE, 2017. 2, 5, 6, 7
- [36] Yoshikatsu Nakajima, Byeongkeun Kang, Hideo Saito, and Kris Kitani. Incremental class discovery for semantic segmentation with rgbd sensing. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019. 2
- [37] Liangliang Nan, Ke Xie, and Andrei Sharf. A search-classify approach for cluttered indoor scene understanding. *ACM Trans. Graph.*, 31(6), 2012. 2
- [38] Lucas Nunes, Xieyuanli Chen, Rodrigo Marcuzzi, Aljosa Osep, Laura Leal-Taixé, Cyrill Stachniss, and Jens Behley. Unsupervised class-agnostic instance segmentation of 3d lidar data for autonomous vehicles. *IEEE Robotics and Automation Letters*, 7(4):8713–8720, 2022. 2, 5, 6, 7
- [39] Lucas Nunes, Rodrigo Marcuzzi, Xieyuanli Chen, Jens Behley, and Cyrill Stachniss. Segcontrast: 3d point cloud feature representation learning through self-supervised segment discrimination. *IEEE Robotics and Automation Letters*, 7(2):2116–2123, 2022. 2, 12
- [40] Songyou Peng, Kyle Genova, Chiyu Jiang, Andrea Tagliasacchi, Marc Pollefeys, Thomas Funkhouser, et al. Openscene: 3d scene understanding with open vocabularies. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 815–824, 2023. 2
- [41] Yinyin Peng, Hui Feng, Tao Chen, and Bo Hu. Point cloud instance segmentation with inaccurate bounding-box annotations. *Sensors (Basel, Switzerland)*, 23, 2023. 2
- [42] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. *Advances in neural information processing systems*, 30, 2017. 1, 2, 16
- [43] Dario Reithage, Johanna Wald, Jurgen Sturm, Nassir Navab, and Federico Tombari. Fully-convolutional point networks for large-scale point clouds. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 596–611, 2018.
- [44] David Rozenberszki, Or Litany, and Angela Dai. Language-grounded indoor 3d semantic segmentation in the wild. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2022. 2
- [45] Jonas Schult, Francis Engelmann, Alexander Hermans, Or Litany, Siyu Tang, and Bastian Leibe. Mask3D for 3D Semantic Instance Segmentation. In *International Conference on Robotics and Automation (ICRA)*, 2023. 1, 4, 5, 15
- [46] Nur Muhammad Mahi Shafullah, Chris Paxton, Lerrel Pinto, Soumith Chintala, and Arthur Szlam. Clip-fields: Weakly supervised semantic fields for robotic memory. *arXiv preprint arXiv: Arxiv-2210.05663*, 2022. 2
- [47] Jianbo Shi and Jitendra Malik. Normalized cuts and image segmentation. *IEEE Transactions on pattern analysis and machine intelligence*, 22(8):888–905, 2000. 2, 3
- [48] Ziyang Song and Bo Yang. OGC: Unsupervised 3D Object Segmentation from Rigid Dynamics of Point Clouds. In *NeurIPS*, 2022. 2
- [49] Carole H Sudre, Wenqi Li, Tom Vercauteren, Sebastien Ourselin, and M Jorge Cardoso. Generalised dice overlap as a deep learning loss function for highly unbalanced segmentations. In *Deep Learning in Medical Image Analysis and Multimodal Learning for Clinical Decision Support: Third International Workshop, DLMIA 2017, and 7th International Workshop, ML-CDS 2017, Held in Conjunction with MICCAI 2017, Québec City, QC, Canada, September 14, Proceedings 3*, pages 240–248. Springer, 2017. 12
- [50] Jiahao Sun, Chunmei Qing, Junpeng Tan, and Xiangmin Xu. Superpoint transformer for 3d scene instance segmentation. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 2393–2401, 2023. 1, 2
- [51] Thang Vu, Kookhoi Kim, Tung M Luu, Thanh Nguyen, Junyeong Kim, and Chang D Yoo. Softgroup++: Scalable 3d instance segmentation with octree pyramid grouping. *arXiv preprint arXiv:2209.08263*, 2022. 2
- [52] Thang Vu, Kookhoi Kim, Tung M. Luu, Xuan Thanh Nguyen, and Chang D. Yoo. Softgroup for 3d instance segmentation on 3d point clouds. In *CVPR*, 2022. 1
- [53] Puzuo Wang, Wei Yao, and Jie Shao. One class one click: Quasi scene-level weakly supervised point cloud semantic

- segmentation with active learning. *ISPRS Journal of Photogrammetry and Remote Sensing*, 204:89–104, 2023. [2](#)
- [54] Weiyue Wang, Ronald Yu, Qiangui Huang, and Ulrich Neumann. Sgpn: Similarity group proposal network for 3d point cloud instance segmentation. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2569–2578, 2018. [1](#)
- [55] Xinlong Wang, Zhiding Yu, Shalini De Mello, Jan Kautz, Anima Anandkumar, Chunhua Shen, and Jose M Alvarez. Freesolo: Learning to segment objects without annotations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14176–14186, 2022. [2](#), [7](#), [12](#), [15](#), [16](#), [17](#)
- [56] Xudong Wang, Rohit Girdhar, Stella X Yu, and Ishan Misra. Cut and learn for unsupervised object detection and instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 3124–3134, 2023. [2](#), [5](#), [6](#), [12](#), [15](#)
- [57] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *Acm Transactions On Graphics (tog)*, 38(5):1–12, 2019. [1](#)
- [58] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pages 9621–9630, 2019. [1](#)
- [59] Z. Wu and R. Leahy. An optimal graph theoretic approach to data clustering: theory and its application to image segmentation. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 15(11):1101–1113, 1993. [2](#)
- [60] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part III 16*, pages 574–591. Springer, 2020. [1](#), [2](#), [12](#)
- [61] Yunhan Yang, Xiaoyang Wu, Tong He, Hengshuang Zhao, and Xihui Liu. Sam3d: Segment anything in 3d scenes. *arXiv preprint arXiv:2306.03908*, 2023. [15](#)
- [62] Zaiwei Zhang, Rohit Girdhar, Armand Joulin, and Ishan Misra. Self-supervised pretraining of 3d features on any point-cloud. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10252–10263, 2021. [2](#), [12](#)

7. Appendix

7.1. UnScene3D as Data Efficient Pretraining

We report additional qualitative details on the data efficient pretraining performance of UnScene3D in Table 5.

We also note that the 3D contrastive pre-training of CSC, similar to other 3D pre-training methods developed for non-transformer backbones [19, 39, 60, 62], was not beneficial for a transformer-based model. A similar observation was also reported in a recent pretraining method [20]. We thus also compare with CSC pretraining on their original 3D backbone (which demonstrated improvement over training from scratch on the same backbone). Our approach can improve notably over both alternatives.

7.2. The effect of noise robust losses.

We adopt DropLoss [56] for our self-training cycles, which is robust to sparse data and missing annotations. In particular, we use a weighted combination of cross-entropy and Dice [49] losses for bipartite-matching with pseudo annotations. We then drop losses for backpropagation which do not have at least τ_{drop} overlap with the annotations from the previous cycle. We evaluate the effect of different noise robust losses for self-training in Table 6. We compare our baseline losses with a 3D extension of the projection loss of [55], and our adaptation of DropLoss from [56]. Our approach does not penalize for missing pseudo masks, which enables more effective self-training to discover previously missed instances.

7.3. Additional Qualitative Results

We show more qualitative results from our method trained on ARKitScenes [2] in Figure 8 and on ScanNet [10] in Figure 9.

7.4. Pseudo Mask Generation Ablations

We also ablate the saliency threshold, oversegmentation parameters, and separation strategy in our pseudo mask generation. If not explicitly stated otherwise in Table 12, we use both 2D and 3D modality features for the pseudo mask generation.

What is the effect of the saliency threshold in pseudo mask generation? We threshold the saliency matrix A with $\tau_{cut} = 0.55$ for geometric-only features and $\tau_{cut} = 0.65$ for combined modalities. Table 7 shows that our approach maintains robust performance across a large range of τ_{cut} thresholds used to estimate salient areas for pseudo masks. In this table we report results using features from combined modalities, but similar behaviour can be observed for the other scenarios as well.

The effect of iterative mask densification. We designed a strategy to leverage a sparse set of relatively clean initial pseudo masks, which are progressively extended with confident self-predictions during later iterations. This leads to a 3x improvement over state of the art in the Average Precision Metric. We could also consider different mask refinement strategies using a mixture of segments, initial masks or self-trained instances. Tab. 8 ablates a mask refinement strategy of discarding previous masks and retaining current predictions. We also consider using Felzenswalb segments directly instead of feature-based pseudo labels. Both these strategies lead to lower performance due to the increased presence of noisy labels, which dominate the training signal.

Robustness to oversegmentation parameters. Table 9 shows that our approach maintains strong robustness to a wide range of oversegmentation parameters for our geometric segments (our used parameters denoted in bold).

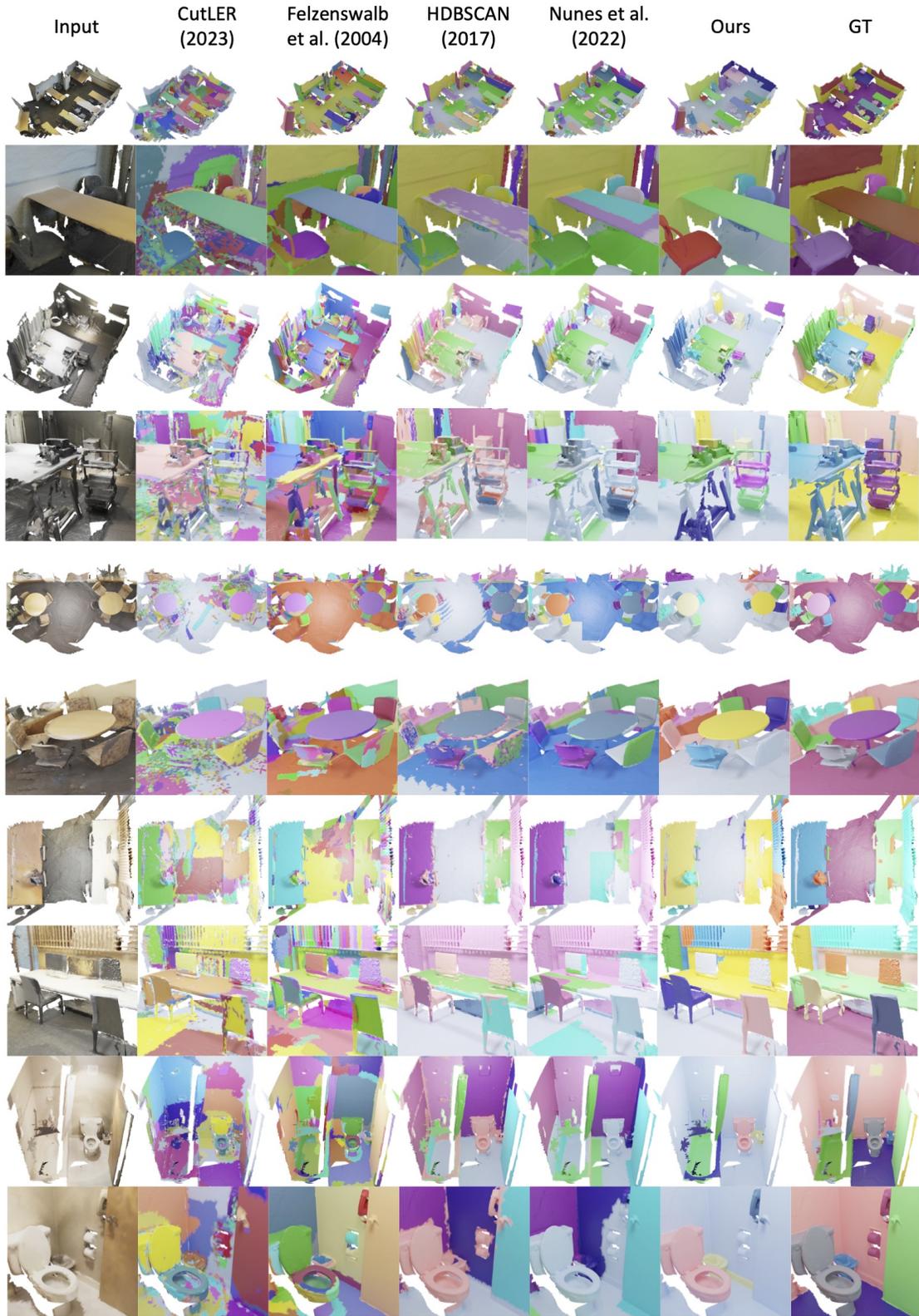
Additional pseudo mask generation hyperparameters. Additionally, we also test the effect of other hyperparameters in our $NCut$ -based pseudo mask generation module, including used distance metrics in the similarity matrix and different methods to separate unconnected patches in the predicted foregrounds. During the foreground separation in the Normalized Cut algorithm, we had an additional condition for the minimum number of foreground segments for the bipartitions. This condition was able to effectively filter out suboptimal partitioning of the full graph leading to separated parts from the full instances. Reducing the size of this parameter can directly lead to a more dense set of initial pseudo masks, with the cost of higher false positive rate. In Table 9 we report a sparser and denser version of the datasets with a minimum number of foreground segments of 8 and 2 accordingly, and show the initial higher scores of the pseudo annotation doesn't necessarily propagate to better downstream self-trained performance.

Finally, we also ablate the effect of our physical connectivity-based foreground separation introduced in Section 3.1. In our main method we separate all set of connected components in the foreground, but only keep the component with the highest eigenvector activation (*Max*). As an alternative we also test a method where we calculate the highest average activation in the connected component (*Avg.*), a method where we keep the component with the largest surface value (*Largest*) and finally, to test the effect of this module, without any kind of connectivity-based separation (*No Sep.*).

7.5. Comparison with methods from the 2D domain

To ensure a fair evaluation of methods operating on different input domains in Table 1. we followed the established procedure of well-known baselines [9, 18, 23]. This involves





Model	Backbone	1%			5%			10%			20%			50%		
		AP@25	AP@50	AP												
Scratch	Bottom-up	22.6	14.1	6.8	45.5	33.3	18.1	54.8	39.2	21.9	61.0	43.4	25.5	67.0	51.4	30.3
CSC [19]	Bottom-up	35.6	22.1	12.5	52.7	39.9	23.3	59.8	43.8	25.0	63.8	48.9	29.6	70.5	56.0	33.6
Scratch	Transformer	24.7	9.3	4.6	48.1	27.6	16.3	59.2	39.1	23.4	66.4	49.6	33.1	78.9	67.5	49.8
CSC	Transformer	17.0	6.8	3.8	44.2	22.7	13.1	55.2	32.3	19.1	62.0	41.2	26.0	73.7	58.2	40.0
Ours	Transformer	43.5	28.4	15.8	63.2	46.8	28.3	70.3	55.7	36.7	72.4	60.7	41.5	78.9	68.0	48.2

Table 5. Unsupervised class-agnostic pretraining with our method can also act as a powerful pretraining strategy, advancing over state of the art. We report pretraining with CSC [19] and UnScene3D, and evaluate the downstream weakly-supervised instance segmentation performance on ScanNet with percentage of limited annotated scenes used denoted in the top row. As we found that CSC degraded performance when using a transformer-based backbone, we also report the performance of training from scratch and CSC on their originally proposed backbone of a sparse UNet with bottom-up voting.

	AP@25	AP@50	AP	AP Final
Initial Pseudo Masks	19.9	10.0	5.9	-
Baseline losses [45]	42.3	16.9	7.2	14.2
Projection loss [55]	35.7	12.1	4.7	7.2
DropLoss [56]	52.9	23.2	10.4	15.9

Table 6. A 3D projection loss struggles with under-determined associations, while DropLoss helps UnScene3D to discover parts of the scene that were missed by the source supervision. We report all metrics after a single iteration and the AP scores after 4 iterations of self-training.

τ_{cut}	AP@25	AP@50	AP
0.40	16.7	9.0	5.2
0.50	20.8	10.7	5.7
0.55	21.0	10.8	5.7
0.60	21.3	11.3	5.8
0.65	19.9	10.0	5.9
0.70	18.2	9.9	5.6
0.80	11.8	5.0	2.6

Table 7. Our pseudo mask generation quality, as measured by AP metrics, maintains robustness to a large range of τ thresholds that extract saliency. Note that this measures the quality of only the pseudo masks; our full approach with self-training produces significantly improved results. In this table we show results and parameters used by our method in bold and report pseudo mask performance generated from both modalities.

using depth information to project 2D predictions into 3D such that all methods are evaluated in the same 3D domain and aggregate multiple predictions through consensus by majority voting or accepting the maximum confidence scores for every voxel location. We also show results evaluated against 2D ScanNet images by projecting our method’s predictions into 2D in Tab. 10, and comparing it to the current state of the art 2D unsupervised segmentation method [56] which demonstrates the usefulness of 3D reasoning.

We also compare to weakly-supervised instance segmentation method SAM3D [61], where powerful class-agnostic

	AP@25	AP@50	AP
Felzenszwalb Masks	35.5	20.6	10.3
Mask Refinement	43.7	24.4	12.4
Mask Addition (Ours)	58.6	32.0	16.0

Table 8. Instead of using masks from previous iteration directly it is the best to keep the initial masks fixed, and iteratively sample plausible predictions to enrich the pseudo dataset during self-training. This method strikes a balance between relatively clean, but sparse labels and increasing number of confident samples. Finally, even though Felzenszwalb oversegmentation yields to higher precision, then our initial mask prediction algorithm, it also includes more background into the training, and this way plateauing at a lower self-training performance.

2D masks are extracted by the powerful SAM model [29]. Here the projected 2D masks are merged into 3D masks iteratively with a bottom-up bidirectional merging approach to achieved cleaner and more view-independent 3D instances. A qualitative comparison on ScanNet can be seen in Table 11, with qualitative comparisons in Figure 10.



Figure 10. While SAM has powerful capabilities in crisp 2D mask generation, when aggregated on 3D, SAM3D tends to over-segment object instances.

SAM3D must resolve view inconsistencies and SAM’s tendency to over-segment objects, which results in SAM3D splitting instances, while UnScene3D is able to achieve complete masks through multi-modal reasoning. We believe integrating SAM or other (weakly-) supervised 2D models into our pipeline to enable multi-modal reasoning is an interesting avenue for future work.

7.6. Additional Implementation Details

Here, we further explain the implementation details of our pseudo mask generation.

Generation Params.				Initial Pseudo Mask			1 Iteration of Self-Training			4 Iterations of Self-Training			
Segment Size	Metric	Separation	Min. # of Foreground	# of Instances	AP@25	AP@50	AP	AP@25	AP@50	AP	AP@25	AP@50	AP
30	Cos	Max	8	2169	21.9	11.5	6.3	53.7	26.2	12.4	55.4	30.3	15.3
50	Cos	Max	8	1414	19.9	10.0	5.9	52.9	23.2	10.4	58.5	32.2	15.9
100	Cos	Max	8	1090	17.4	8.0	4.2	33.1	10.2	3.9	39.6	13.7	5.3
200	Cos	Max	8	584	11.0	3.7	1.8	24.3	8.7	2.1	26.1	9.7	2.4
400	Cos	Max	8	319	6.4	2.5	1.1	19.1	3.9	1.2	19.9	3.2	1.0
50	L2	Max	8	1539	20.1	10.6	5.4	49.0	21.7	9.8	55.3	38.4	14.3
100	L2	Max	8	805	13.3	5.3	2.6	30.8	8.3	2.8	39.0	12.7	5.0
50	Cos	No Sep.	8	125	4.3	0.3	0.1	4.3	0.5	0.2	4.9	0.6	0.2
50	Cos	Largest	8	620	11.5	4.9	2.5	11.5	1.5	0.4	12.9	2.2	12.9
50	Cos	Avg.	8	1078	16.8	9.1	5.1	36.4	12.5	4.9	43.8	17.8	7.5
30	Cos	Max	2	2909	29.0	15.6	8.7	53.6	28.6	14.2	54.2	29.8	15.4
50	Cos	Max	2	2512	24.9	12.4	7.2	56.5	29.8	15.0	51.3	26.2	12.6
100	Cos	Max	2	2317	23.1	12.3	6.8	51.8	24.4	11.6	57.1	31.3	15.6
200	Cos	Max	2	2181	28.4	15.5	8.9	54.6	28.7	13.7	56.6	31.4	15.6
400	Cos	Max	2	1373	20.6	11.1	6.3	51.0	24.8	11.8	55.8	30.3	15.2
50	L2	Max	2	2496	28.6	15.8	9.0	55.8	29.6	14.6	54.8	30.3	15.3
100	L2	Max	2	1668	23.4	12.7	7.3	53.1	25.0	11.3	56.3	27.7	12.9
50	Cos	No Sep.	2	159	0.2	0.5	3.6	5.4	0.6	0.3	3.9	0.4	0.2
50	Cos	Largest	2	1026	14.1	7.2	3.9	11.5	1.8	0.5	14.5	2.5	0.7
50	Cos	Avg.	2	2053	23.3	12.0	6.8	52.5	27.4	12.7	54.9	29.9	14.9

Table 9. We denote the parameters used by our method in bold. We show that our method is robust to a wide range of numbers regarding segments sizes and different similarity metrics, and only degrades somewhat in performance when segments are constrained to be too large. We also show that the separation of physically distant foreground patches is important and it is beneficial to use the activation of the eigenvector for the best results. Finally, we show that denser initial mask predictions lead to quantitatively better initial pseudo annotations, and even better self-training performance after a single iteration, but underperforming in their final scores. This behaviour can be explained by the larger false positive ratio in the denser initial predictions, which is propagating through all iterations, but thanks to the noise robust losses and iterative refinement of predictions the sparse set of labels can be effectively used. In this table we report results using both modalities for the initial pseudo mask generation, and number predicted pseudo instances in the official validation split of the ScanNet dataset.

	AP@25 (2D)	AP@50 (2D)	AP (2D)
CutLER (2D)	7.8	2.8	0.7
Ours (projected)	60.0	38.1	21.1

Table 10. 2D evaluation on ScanNet images.

	AP@25	AP@50	AP
SAM3D	37.2	11.8	3.7
SAM3D with GT Segments	47.6	24.1	10.8
Ours	58.5	32.2	15.9

Table 11. UnScene3D achieves significantly better performance on ScanNet than SAM3D through our strong multi-modal reasoning.

Pseudo code for masked NCut We show the pseudo code-style implementation for the masked normalized cut algorithm generating multiple instances as pseudo masks. The full algorithm can be seen in 1.

3D Adaptation of FreeMask We also evaluate an alternative pseudo mask segmentation algorithm besides the masked *NCut* method. In the 2D domain FreeSOLO [55] also followed a two stage pipeline first generating the pseudo annotations, and then refine those predictions through a series of self-training cycles. We followed their intuition to take a self-supervised pretrained backbone and

extract it’s deep features at multiple levels of the decoder. While in standard pretrained UNet-style models early features represent global context, final features and local semantic meaning, intermediate features can act as an useful proxy to extract self-similar regions in the input samples. In our implementation we used the same backbone features of [3, 19] for the same 2D-3D setup and extracted the penultimate layer features for the self-similarity calculation. Then sampled the feature space with the Furthest Point Sampling [42] strategy to get a more limited set of anchor points, later used to extract self-similar regions. For every seed point we took similarity scores with the other features of the full scene and thresholded it to extract salient regions. Finally, we used the efficient Non Maximum Suppression implementation from [55] to sort the predicted salient areas and filter out overlapping regions. We also used average similarity score combined with the salient region area to get *maskness scores* for every salient region, directly following the original implementation. We report comparative results of the masked *NCut* algorithm and our FreeMask 3D adaptation after self-training in Table 3. of the main paper and in Table 12 of the initial pseudo mask scores.

We also note here that while there is a difference in the initial pseudo mask qualities for the different methods, the

Algorithm 1: Masked NCut on 3D segments

Data: $\mathcal{S} = \{s_1, \dots, s_N\}, \mathcal{F} \in \mathcal{R}^{N \times D}$,
 $\mathcal{C} = \{(s_1, s_k), (s_1, s_l), \dots\}$
Result: $\mathcal{M} = \{m_1, \dots, m_M\}$

```
1  $\mathcal{M} \leftarrow \{\}$ 
2 while  $j \leq \max\_inst\_num$  do
3    $\mathcal{F}' \leftarrow \mathcal{F}$ 
4    $\mathcal{F}'[\mathcal{M}] \leftarrow 0.$  // Mask out previous insts.
5    $\mathcal{W} \leftarrow \mathcal{F} \times \mathcal{F}^T$  // Feature similarity
   // Saliency with connected graph
6    $\mathcal{W}_{i,k} = \begin{cases} 1. & \text{if } \mathcal{W}_{i,k} \geq \tau_{cut} \\ \epsilon & \text{if } \mathcal{W}_{i,k} < \tau_{cut} \end{cases}$ 
7    $\mathcal{D}_{i,i} = \sum_k \mathcal{W}_{i,k}$ 
   // Get 2nd smallest eigenvector
8    $\lambda, \mathbf{v} \leftarrow \text{eigh}(\mathcal{D} - \mathcal{W}, \mathcal{D}, -2)$ 
9    $m_i = \begin{cases} 1 & \text{if } v_i \geq \text{mean}(\mathbf{v}) \\ 0 & \text{if } v_i < \text{mean}(\mathbf{v}) \end{cases}$ 
   // Invert bipartition if too large
10  if  $\text{sum}(\mathbf{m}) > D/2$  then
11     $\mathbf{m} = 1 - \mathbf{m}$ 
12     $\mathbf{v} = -1. * \mathbf{v}$ 
   // Separate unconnected components
13   $v_{max} = \max(\mathbf{v})$ 
14   $\hat{\mathbf{m}} = \text{sep}(\mathbf{v}, v_{max}, \mathcal{C})$ 
15   $M \leftarrow M \cup \{\hat{\mathbf{m}}\}$ 
```

	Modality	AP@25	AP@50	AP
FreeMask	3D	13.7	7.2	3.7
Ours	3D	13.8	4.7	2.0
FreeMask	2D	15.3	6.6	2.9
Ours	2D	15.6	7.2	3.6
FreeMask	both	17.9	7.5	3.7
Ours	both	19.9	10.0	5.9

Table 12. We compare pseudo mask generation from 3D-only features (3D), color-only features (2D), and both color and geometry (both) signal, as well as with pseudo annotation generation algorithm FreeMask. We compare the quality of the initial pseudo mask dataset using our masked *NCut* algorithm and the adaptation of FreeMask [55] to 3D. We see that the normalized cut-based method is superior for both modalities.

downstream performance is way more significant. This can be explained by the nature of the pseudo masks. *NCut* provides a clean and sparse set of annotation, which is easy to densify for following iterations. On the other hand, the more dense, but noisy FreeMask predictions remain in the training for the duration of the whole training, hindering the performance of the self-trained model with noisy supervision.