# MAELi: Masked Autoencoder for Large-Scale LiDAR Point Clouds

Georg Krispel[1]   David Schinagl[1,2]   Christian Fruhwirth-Reisinger[1,2]   Horst Possegger[1]   Horst Bischof[1,2]

{georg.krispel,david.schinagl,christian.reisinger,possegger,bischof}@icg.tugraz.at

[1] Graz University of Technology [2] Christian Doppler Laboratory for Embedded Machine Learning

## Abstract

*The sensing process of large-scale LiDAR point clouds inevitably causes large blind spots,* i.e. *regions not visible to the sensor. We demonstrate how these inherent sampling properties can be effectively utilized for self-supervised representation learning by designing a highly effective pre-training framework that considerably reduces the need for tedious 3D annotations to train state-of-the-art object detectors. Our* Masked AutoEncoder for LiDAR point clouds (MAELi) *intuitively leverages the sparsity of LiDAR point clouds in both the encoder and decoder during reconstruction. This results in more expressive and useful initialization, which can be directly applied to downstream perception tasks, such as 3D object detection or semantic segmentation for autonomous driving. In a novel reconstruction approach, MAELi distinguishes between empty and occluded space and employs a new masking strategy that targets the LiDAR's inherent spherical projection. Thereby, without any ground truth whatsoever and trained on single frames only, MAELi obtains an understanding of the underlying 3D scene geometry and semantics. To demonstrate the potential of MAELi, we pre-train backbones in an end-to-end manner and show the effectiveness of our unsupervised pre-trained weights on the tasks of 3D object detection and semantic segmentation.*

## 1. Introduction

Thanks to recent large-scale and elaborately curated datasets, such as the Waymo Open Dataset [35], we have witnessed tremendous progress in a diverse variety of 3D perception tasks crucial for autonomous driving. However, even with the aid of such cost-intensive datasets, models remain only transferable to other domains while suffering significant performance drops [40].

Self-supervised representation learning (SSRL) provides a technique to reduce the costly labeling effort. The overall idea is to learn a universal feature representation in an unsupervised way, that is later utilized for a specific downstream task, such as object detection. One of the most common
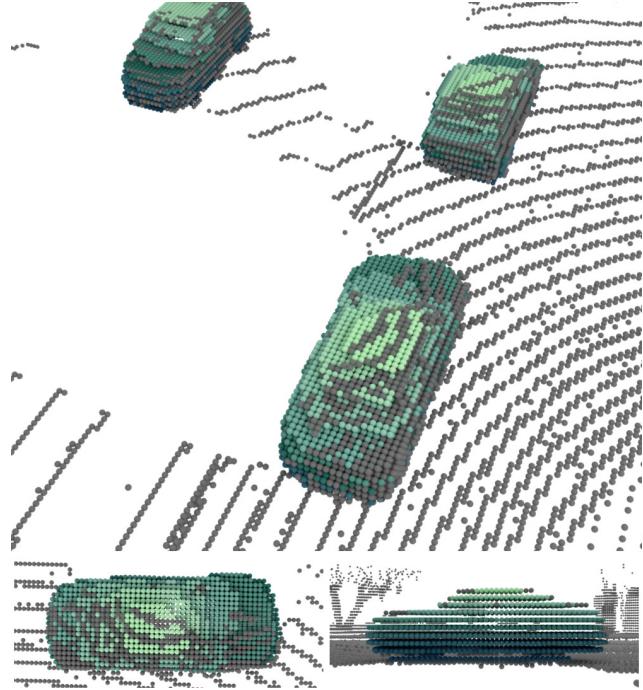


Figure 1. Reconstructed point cloud. In contrast to existing self-supervised representation learning (SSRL) approaches which simply reconstruct the initial point cloud (gray), MAELi learns an expressive feature representation which captures the full geometric object structure of objects, without any ground truth labels. For visualization purposes, we color-coded the points by their $z$-coordinate and removed the reconstructed ground plane.

approaches in 3D is to learn representations via *point cloud reconstruction* [1, 25, 30, 37, 38, 41, 56]. There, the task is to restore removed parts of a point cloud and thereby learning an implicit understanding of the scene and the object's geometric structure. This is especially useful for full 3D point clouds generated from CAD models, such as Model-Net [44] or ShapeNet [6]. Recently, these methods were adapted for large-scale point clouds within the automotive domain [17, 27, 37].

Existing SSRL approaches, however, neglect inherent, but fundamental properties of LiDAR point clouds: i) We

cannot sense beyond a hit surface (*i.e.* we are limited to 2.5D perception), and ii) LiDAR sensors have a limited angular resolution. In this work, we adapt to these properties and propose *MAELi*, a transformer-less masked autoencoder (MAE), which does not simply follow the straightforward methodology of reconstructing the original LiDAR point cloud. Instead, we present a novel reconstruction approach that allows us to go beyond the (visible) points. As a result, MAELi learns how objects look like from any viewing direction, which leads to strong pre-training weights with favorable generalization capabilities. As illustrated in Figure 1, it implicitly learns to reconstruct entire objects whilst training on single frames and in a genuinely unsupervised way without any ground truth whatsoever.

Intuitively, we explicitly distinguish *occupied*, *empty* and *unknown* space. When tracing the path of a LiDAR beam from the sensor to an object (and back), the intervening space is considered *emtpy*, the object itself is *occupying* space at the point of impact and the space behind the object is *unknown* due to occlusion. Moreover, no conclusions can be made about the regions that were not covered due to the limited resolution of the LiDAR. Thus, we task our model with reconstructing removed parts of the point cloud, but we do not penalize it for the completion of structures in *unknown* regions, *i.e.* occluded or unsampled areas.

During training, the model encounters a wide range of objects, differing in pose and sampling density. Despite the missing labels, our unique objective allows the model to reconstruct whole objects, *i.e.* implicitly capturing the entire geometric structure. In this process, MAELi learns a representation that more closely align with the underlying geometric structure, rather than simply mimicking the specific sampling patterns observable in a LiDAR point cloud.

For evaluation, we select the most predominant automotive perception tasks, *i.e.* object detection and semantic segmentation. Our memory-efficient, sparse decoder structure enables efficient pre-training of state-of-the-art object detectors in an end-to-end fashion on a single GPU. In extensive experiments on the Waymo Open Dataset [35], KITTI [2, 12, 24] and ONCE [26], we demonstrate that MAELi is highly effective for pre-training various state-of-the-art 3D detectors and semantic segmentation networks.

In summary, our contributions are threefold:

- We propose a LiDAR-aware SSRL approach to pre-train 3D backbones applicable to various architectures and downstream tasks.

- We introduce a novel masking strategy and reconstruction loss for unsupervised representation learning, especially designed for LiDAR properties.

- We show the effectivity of our pre-trained, visually verifiable representation improving several baselines for 3D object detection and semantic segmentation.

## 2. Related Work

**SSRL for 2D Imagery and 3D Points Clouds:** Self-supervised representation learning strives to learn beneficial representations before introducing any supervision in the form of manually labeled ground truth data. These representations are utilized to improve the results on respective downstream tasks or to reduce the required amount of labeled training data.

*Contrastive learning* approaches task a model to maintain similar embeddings for the same data instance when transformed with different augmentations. Consequently, different data instances should lead to diverging embeddings. Initially applied to 2D imagery [8, 16, 18, 39], these methods were also adapted for point clouds [21, 23, 28, 29, 31, 31, 42, 46, 54, 58]. Naturally, the granularity of the induced consistency loss defines the *semantic level* on which the model agrees upon. In other words, contrastive learning on a global embedding describing an entire image or point cloud is more suitable for downstream tasks like classification. Tasks like object detection or semantic segmentation, however, necessitate a more fine-grained treatment. The causality dilemma is to sample semantically coherent regions with a proper level of detail without knowing what is semantically coherent. For example, Yin *et al.* [54] generate proposals via farthest point sampling and ball queries after removing the ground plane. TARL [29] and STSSL [43] extend the analysis to multiple time frames to cluster objects of interest, requiring the use of globally registered point clouds. This added temporal dimension introduces an additional layer of information, distinguishing them from single-frame based methods.

**SSRL via Reconstruction:** Recently, generative self-supervised representation learning approaches have been on the rise. One of its most successful concepts is the *denoising autoencoders*. Based on the encoder's output embedding, the decoder is tasked to reconstruct the denoised input and if successful, the encoder is forced to learn a useful representation resilient to noise. Especially the reconstruction of a masked input gained huge traction across various domains of application, among others like natural language processing (NLP) [11] and 2D imagery [15] also on point clouds. The predominant research focus started on learning representations on full 3D, synthetic or indoor datasets [7, 13, 19, 25, 30, 38, 49, 52, 56, 57], *e.g.* Model-Net [44] or ScanNet [10].

Recently, a few works consider LiDAR point clouds, such as [17, 27, 37, 45, 48, 51]. Xie *et al.* [45] require full supervision, whereas we do not need any labels for pre-training. MV-JAR [48] tasks a backbone to reconstruct voxel's masked positional encoding and inner point distribution, enforcing the network to reconstruct the exact sampling pattern. In Occupancy-MAE [27], the authors use an MAE-approach to pre-train common 3D voxel backbones
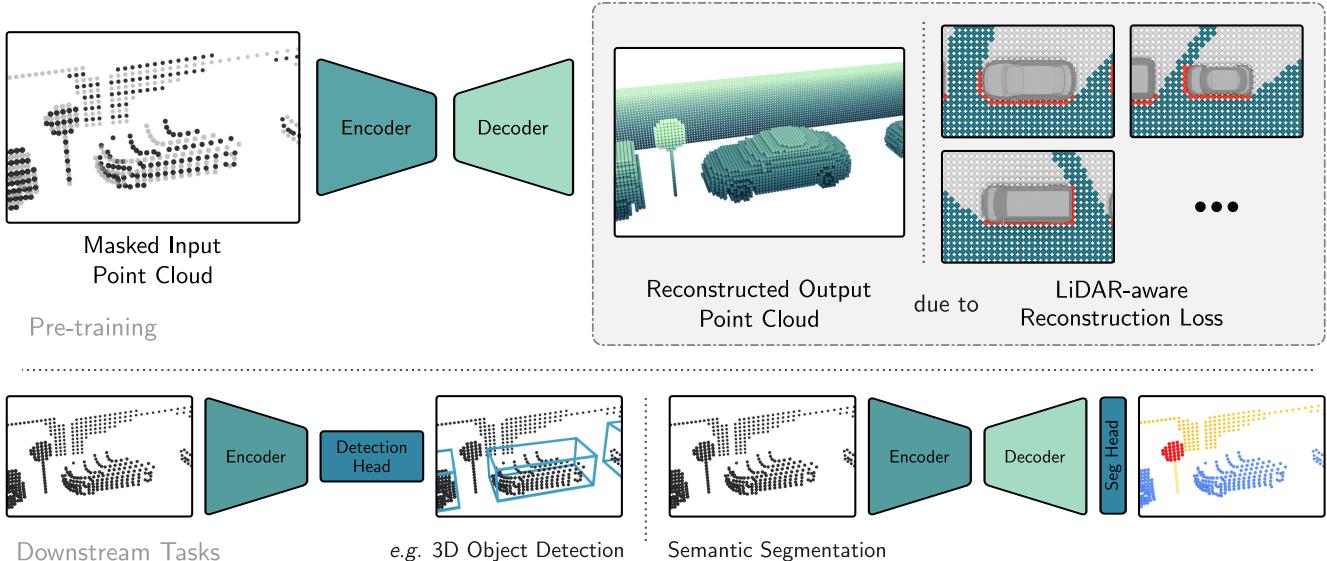
Figure 2. Schematic overview of our MAELi pre-training. The task of our sparse decoder is to reconstruct the missing parts of the masked input point cloud. Thereby, the encoder is forced to learn a reasonable representation usable for a downstream task, *e.g.* 3D object detection. Since we do not penalize our network for reconstructing voxels at areas not visible to the LiDAR, after seeing numerous different samples, it learns to reconstruct occluded parts without any ground truth labels, leading to a more expressive feature representation. The ground plane has been removed for visualization purposes.

without the bird's-eye view (BEV) encoder, while GD-MAE [51] introduces a newly developed elaborate multi-stage transformer. Both approaches employ a dense decoder and do not differentiate between the empty and occluded spaces inherent in a LiDAR point cloud. In contrast, our reconstruction strategy combined with our sparse decoder enables us to pre-train the entire encoder of the most common 3D object detectors without penalizing reconstruction in unsampled areas on a single GPU.

Only a few studies consider the inherent properties of a LiDAR point cloud, *i.e.* its sampling resolution and 2.5D perception. In this context, Wang *et al.* [38] synthesize occlusions on small-scale full 3D point scans and leverage this to pre-train an encoder. Hu *et al.* [20] generate visibility maps via raycasting to mitigate inconsistent object augmentation [50], serving as an additional input to a detection network. Xu *et al.* [47] enhance ground truth object points by grouping similar instances and utilize them to train an auxiliary task, estimating the likelihood of an occluded area being occupied by an object. GeoMAE [37] proposes a masked autoencoder, which is trained to reconstruct underlying point statistics and surface properties, *i.e.* an estimation of normal and curvature based on neighboring voxels. ALSO [3], most closely related to our approach, employs SSRL via surface interpolation. It generates query points along LiDAR rays situated in front and behind the detected hits and instructs the network to predict the occupancy of these selected points, treating occluded areas as occupied. While this strategy is straightforward, it may en-

counter limitations in a multi-LiDAR setup on the ego vehicle, a configuration notably utilized in the Waymo Open Dataset [35]. In comparison, our approach is capable to handle such complex scenarios robustly and can implicitly learn geometric structure beyond sampled query points. The self-supervision tasks of both, GeoMAE and ALSO, are deemed achieved once the network is able to interpolate the underlying surface. While the resulting feature representation already shows promising results, our LiDAR-aware reconstruction demonstrates further improvements as our model inherently captures the whole object geometry.

## 3. MAE for Large-scale LiDAR Point Clouds

We aim to significantly reduce the expensive labeling effort for large-scale LiDAR point clouds via self-supervised representation learning (SSRL). We build upon the successful masking and reconstruction paradigm, but address the fundamental limitation that existing approaches primarily focus on reconstructing the original input point cloud. While this strategy is effective for full 3D models, such as those generated by CAD renderings [6, 44], it limits the usefulness of the learned representation for LiDAR point clouds in two principal aspects.

First, the limited angular resolution of a LiDAR sensor induces gaps between the LiDAR beams. Simply reconstructing the original point cloud would mean to penalize the model for (correctly) reconstructed points in these gaps. Second, a single LiDAR sweep cannot capture objects fully. Once a beam is reflected by a surface, the sensor is unable

to capture any spatial information from objects situated behind that surface. Thus, models based on standard reconstruction are penalized for completing occluded parts and, as a consequence, hindered from inherently understanding the underlying objects and learning implicit contextual information. With MAELi, we address these challenges by introducing our *LiDAR-aware loss*. Rather than penalizing uniformly, this loss specifically targets known regions sampled by the LiDAR.

While our approach is versatile and can be applied to a wide range of tasks, we chose to demonstrate its efficacy specifically on 3D detection and semantic segmentation due to their paramount importance and widespread application in the field. For 3D detection, we pre-train the encoder by attaching a reconstruction decoder to its final layer. Once pre-training is complete, we discard the decoder and utilize the encoder's weights as an initialization for the subsequent detection task. For semantic segmentation, we use the weights of both the encoder and decoder as an initialization for downstream fine-tuning. This is illustrated in Figure 2.

In the following, we briefly explain our sparse decoder (Sec. 3.1), before describing our reconstruction objective (Sec. 3.2) and masking strategy (Sec. 3.3).

**Definitions & Notations:** Utilizing sparse operations, a voxel position is considered *active* if any of its corresponding features deviates from zero and is consequently included in the computation process. Only the active sites are actually stored in memory. Furthermore, we follow the common notation, *e.g.* [9], for sparse convolutions, where the *voxel/tensor stride* **s** refers to the distance between two voxels w.r.t. the highest voxel resolution along each axis. For example, applying two (or three) 3D convolutions with a stride of **2** leads to a feature map with a tensor stride of **4** (or **8**). A downsampling step increases the stride, while an upsampling step decreases it.

So, let $\mathcal{V}^{E,\mathbf{s}} = \{\mathbf{v}_i^{E,\mathbf{s}}\}_{i=1...M^{E,\mathbf{s}}}$ be the $M^{E,\mathbf{s}}$ *active* voxel positions of a certain *tensor stride* **s** in the sparse encoder and $\mathcal{V}^{D,\mathbf{s}} = \{\mathbf{v}_i^{D,\mathbf{s}}\}_{i=1...M^{D,\mathbf{s}}}$ the $M^{D,\mathbf{s}}$ active voxel positions in the sparse decoder, where $\mathbf{v}_i \in \mathbb{R}^3$. To avoid cluttering the notation, unless stated otherwise, $\mathcal{V}^{\mathbf{s}}$ and $\mathbf{v}_i^{\mathbf{s}}$ imply the decoder and default to $\mathcal{V}^{D,\mathbf{s}}$ and $\mathbf{v}_i^{D,\mathbf{s}}$, respectively.

### 3.1. Sparse Reconstruction Decoder

Our goal is to obtain more expressive feature representations by pre-training the entire backbone for the respective downstream task. Contrary to existing sparse, voxel-based encoder/decoder structures for LiDAR point clouds, such as Part-A$^2$ [33], we must address a key difference. The typical approach involves voxelizing and processing the sparse data via dedicated sparse convolutions (SCs). Unlike its dense counterpart, a SC is only applied if the kernel covers any

active sites. Even with small $3 \times 3$ kernels, these active sites dilute rapidly, escalating computational effort. Thus, methods like [32, 50, 55] use submanifold sparse convolutions (SSCs) [14], where the kernel center is placed only on active sites, considering only active sites covered by the kernel while maintaining favorable memory consumption. The decoder then re-uses the active sites of the respective encoder layer during upsampling. This, however, renders the common upsampling schema ineffective for reconstructing voxels not present in the encoder, making it unsuitable for a reconstruction task.

To perform reconstruction on a point cloud, we require a decoder that can also restore the removed voxels and expand beyond the active sites present in the encoder. Dense upsampling to the original spatial resolution is a possibility, but it is rather impractical for pre-training larger architectures on limited hardware setups due to excessive memory and computational demands. Thus, we allow each upsampling layer to grow (cubically), but add a subsequent pruning layer that learns to remove redundant voxels, as illustrated in Figure 3. For this, we leverage the idea of small-scale point cloud reconstruction [9] and apply a $1 \times 1$ convolution to the sparse feature map that learns the pruning. This way, the decoder is able to reconstruct and complete parts of the point cloud, while the pruning layer removes superfluous voxels. To obtain stronger feature representations, we propose an extended reconstruction objective in the following.

### 3.2. Reconstruction Objective

We formulate our reconstruction objective to incorporate a deeper understanding of the underlying environmental structure. Intuitively, it should be more straightforward for the model to grasp the complete appearance of a car rather than merely reconstructing specific points captured by LiDAR beams. Moreover, the pre-trained representation should be more closely aligned with the objectives of downstream tasks, *e.g.* fitting a 3D bounding box around an entire car. Consequently, during the reconstruction within the decoder, we distinguish three categories of voxels: *occupied*, *empty*, and *unknown*. These categories are depicted in Figure 4.

Occupied voxels contain surface points from the original point cloud (before masking) and should thus be part of the reconstruction. Empty voxels are the ones traversed by the LiDAR beam without hitting any surface and thus, should remain empty. Finally, we categorize a voxel as *unknown* if neither of the first two cases applies, *i.e.* these are either occluded or were not sensed by the beam due to the limited angular resolution. This categorization enables the reconstruction to grow beyond the initially sensed point cloud. Intuitively this makes sense, since it is counter-productive to punish the network for reconstructing points not sampled
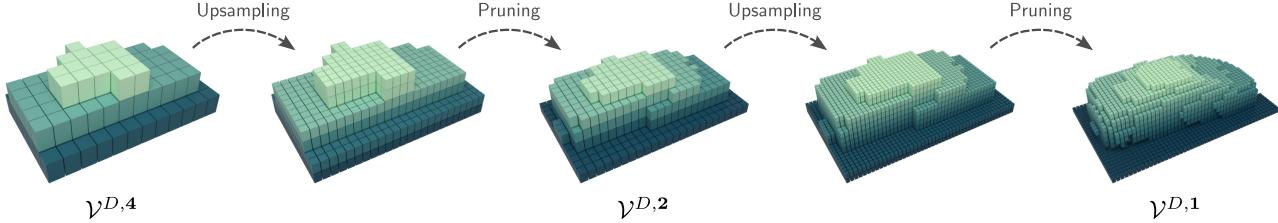
Figure 3. Reconstruction of an actual car within our decoder. An upsampling step halves the *voxel stride* and voxel size and increases the number of voxels. Subsequently, superfluous voxels are removed during the pruning step. We cropped the car and color coded the $z$-coordinate for visualization purposes.

by the LiDAR, even if they are part of the underlying object.

We observe that the discrete voxelization can cause inaccuracies, especially at low resolution voxel grids. To mitigate this, we reduce the loss for an empty voxel via the perpendicular distance $d_i^{\mathbf{s}}$ from the voxel center $\mathbf{v}_i^{\mathbf{s}}$ to the closest LiDAR beam. More formally, we define the weight $w_i^{\mathbf{s}}$ for a voxel as

$$w_i^{\mathbf{s}} = \begin{cases} 0 & \text{if } \mathbf{v}_i^{\mathbf{s}} \text{ is } \textit{unknown}, \\ 1 & \text{if } \mathbf{v}_i^{\mathbf{s}} \text{ is } \textit{occupied}, \\ 1 - \frac{2\, d_i^{\mathbf{s}}}{d_v^{\mathbf{s}}} & \text{if } \mathbf{v}_i^{\mathbf{s}} \text{ is } \textit{empty}, \end{cases} \quad (1)$$

where $d_v^{\mathbf{s}}$ denotes the length of a voxel diagonal at stride $\mathbf{s}$. Note that a setup with multiple LiDARs (*e.g.* Waymo [35]) can simply be handled by iterating the categorization process for each LiDAR. There, the number of *unknown* voxels would be reduced due to the additional LiDAR sensors.

We require the same categorization for voxels of lower resolution and larger strides, respectively. For example, a voxel of stride $\mathbf{s}$ spatially covers eight voxels of stride $\mathbf{s}/2$ after upsampling (see Figure 3). Here, however, we need to consider that if a low-resolution voxel is pruned, it can no longer generate a self-supervision signal for its higher-resolution voxels during upsampling. Thus, a low-resolution voxel is considered *occupied* if it incorporates any high-resolution *occupied* voxel (red square in Figure 4). Similarly, it is *unknown* (gray square) if it incorporates any *unknown* voxel but no *occupied* ones. Otherwise, the low-resolution voxel is categorized as *empty*.

In summary, the objective during self-supervised pre-training is to correctly classify whether a voxel is *occupied* or *empty*, while not penalizing *unknown* ones. For this, we use a weighted binary cross-entropy loss

$$\mathcal{L} = -\frac{1}{\widetilde{M}} \sum_{\mathbf{s} \in \mathcal{S}} \sum_{i=1}^{M^{\mathbf{s}}} w_i^{\mathbf{s}}\, l_i^{\mathbf{s}}, \text{ with} \quad (2)$$

$$l_i^{\mathbf{s}} = \left[ y_i^{\mathbf{s}} \log x_i^{\mathbf{s}} + (1 - y_i^{\mathbf{s}}) \log(1 - x_i^{\mathbf{s}}) \right], \quad (3)$$

where $\mathcal{S}$ denotes the set of all strides; $\widetilde{M}$ is the amount of all *occupied* and *empty* voxels in the decoder; $y_i^{\mathbf{s}}$ is the actual occupancy of the voxelized input point cloud, *i.e.* 1 if $\mathbf{v}_i^{\mathbf{s}}$ is *occupied* and 0 otherwise; and $x_i^{\mathbf{s}}$ is our model's predicted occupancy probability. In other words, we penalize each pruning step for removing *occupied* and maintaining *empty* voxels, but we do not induce any loss for *unknown* voxels.

### 3.3. Masking Strategy

In methods applied to small-scale full 3D point clouds, a common masking strategy is to apply Farthest Point Sampling (FPS) and the K-Nearest Neighbor algorithm to create overlapping patches [30, 56], which are then randomly removed. Applying FPS to LiDAR point clouds would, however, predominantly pick isolated points, which are usually single outliers far off. Masking these would not be benefi-



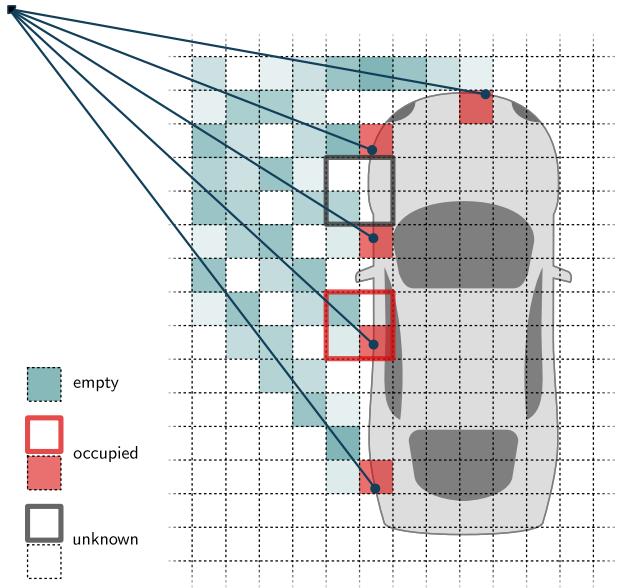Figure 4. A loss should only be induced for voxels, where we *actually know* whether the space is empty or not. LiDAR beams traverse *empty* voxels (blue) until they hit a surface, *i.e.* an *occupied* voxel (red). All other voxels are considered *unknown*. The loss for an *empty* voxel is weighted by the proximity of its center to the nearest beam (different shades of blue) to counter inaccuracies due to discrete voxelization.

5

cial to learn a strong feature representation, as such outliers usually do not correspond to real scene structures and have only few neighboring points. Instead, we exploit the voxelization step already in place and randomly mask voxels.

Furthermore, consider the sampling resolution of a standard spinning LiDAR sensor: It sends out multiple beams, waits for their return and then derives the distance to a hit obstacle, *i.e.* the *range*, from the elapsed time. Every beam has an inclination and rotates around a common vertical axis, defining an azimuth. Since two specific beams enclose a constant angle, the spatial resolution and thus, the number of points decreases with the distance to the hit object. Considering SSRL via point cloud reconstruction, there is less self-supervision with increasing distance to the sensor.

Intuitively, if we remove points from nearby regions, such that they look similar to sparser regions further away, the model should be able to better generalize to these. We incorporate this *spherical masking* idea by reducing the angular resolution in both, azimuth and inclination. To do this efficiently, we subsample the LiDAR's range image. In particular, we randomly sample two integers $1 \leq m_r, m_c \leq 4$ and filter all rows $r$ and columns $c$ of the range image for which $r \mod m_r \neq 0$ or $c \mod m_c \neq 0$, respectively. For example, using only every second row and column halves the angular resolution. This way, an object is sampled as if it would have been further away, but we are able to induce a stronger self-supervision signal during our reconstruction task.

# 4. Experiments

We demonstrate our MAELi pre-training by focusing on key tasks in automotive perception, namely object detection (Section 4.1) and semantic segmentation (Section 4.2). To validate our approach, we employ a range of architectures and pre-train them on widely-recognized datasets in these domains. Specifically, we use the Waymo Open Dataset [35], KITTI [2, 12, 24] and ONCE [26] datasets.

We compare against state-of-the-art approaches and, unless stated otherwise, report their official results, taken either from official benchmarks or the corresponding publications. An additional ablation study to isolate the contributions of various components in our methodology is provided in the supplementary material.

**Implementation Details:** We integrate MAELi into the OpenPCDet [36] framework (v0.5.2) and employ the Minkowski Engine [9] to design our sparse decoder. We use a voxel size of $[0.05, 0.05, 0.1]$, $[0.1, 0.1, 0.15]$ and $[0.1, 0.1, 0.2]$ for KITTI, Waymo and ONCE, respectively. Given the respective architecture and dataset, we pre-train for 30 epochs without any labeled ground truth. We use Adam [22] and a one-cycle policy [34] with maximum learning rate of 0.003. We conducted all experiments on a single NVIDIA® Quadro RTX™ 8000.

| Method | Pre-train | mAP | Car | Pedestrian | Cyclist |
|---|---|---|---|---|---|
| SECOND [50] | - | 65.35 | 81.50 | 48.82 | 65.72 |
| + ALSO [3] | nuScenes | 68.07 | 81.78 | 54.24 | 68.19 |
| + ALSO [3] | KITTI 3D | 67.68 | 81.97 | 51.93 | 69.14 |
| + ALSO [3] | KITTI360 | 68.31 | 81.79 | 52.45 | 70.68 |
| **+ MAELi** | Waymo | 69.05 | 81.70 | 54.34 | 71.11 |
| **+ MAELi** | KITTI 3D | 68.18 | 81.37 | 54.37 | 68.80 |
| **+ MAELi** | KITTI360 | 69.51 | 81.51 | 54.74 | 72.28 |
| PV-RCNN [32] | - | 70.57 | 84.50 | 57.06 | 70.14 |
| + STRL [21] | KITTI 3D | 71.46 | 84.70 | 57.80 | 71.88 |
| + GCC-3D [23] | Waymo | 71.26 | - | - | - |
| + PropCont [54] | Waymo | 72.92 | 84.72 | 60.36 | 73.69 |
| + ALSO [3] | nuScenes | 72.53 | 84.86 | 57.76 | 74.98 |
| + ALSO [3] | KITTI 3D | 72.76 | 84.72 | 58.49 | 75.06 |
| + ALSO [3] | KITTI360 | 72.96 | 84.68 | 60.16 | 74.04 |
| **+ MAELi** | Waymo | 72.15 | 84.80 | 57.46 | 74.18 |
| **+ MAELi** | KITTI 3D | 72.58 | 82.94 | 59.51 | 75.29 |
| **+ MAELi** | KITTI360 | 73.43 | 84.71 | 62.29 | 73.27 |

Table 1. Performance comparison of pre-trained weights fine-tuned on the full KITTI 3D training set for SECOND and PV-RCNN. Results are reported on the KITTI 3D *val* set using the standard $R_{40}$ metric.

## 4.1. 3D Object Detection

Our pre-training methodology MAELi is well-suited for 3D object detectors as our loss formulation lets the model learn how an object should look like. These pre-trained weights generalize well across datasets and allow for data-efficient fine-tuning of a detector on a target dataset.

For 3D detection, we design a sparse decoder with four blocks, each inverting one downsampling operation of the encoder (exact architecture in the supplementary material). We pre-train the entire encoder including the dense BEV backbone [50]. Therefore, to utilize the sparse processing capabilities of the decoder, we take the active voxels from the last layer of the sparse 3D encoder and sample them from the BEV feature map. To fine-tune the different 3D detectors with our pre-trained weights, we warm-up the detection head by freezing the encoder for one epoch and train end-to-end according to the default OpenPCDet configuration.

OpenPCDet provides several augmentation techniques, *i.e.* random flip/rotation/scaling and *ground truth sampling* [50]. During pre-training, we utilize random flip/rotation/scaling in addition to our masking strategy (see Section 3.3). For fine-tuning, we adopt all of them as-is, except *ground truth sampling* [50] during the data efficiency experiments. We cannot apply its vanilla form there, as it samples ground truth objects across the entire dataset and copy-pastes these randomly into frames which contain only few objects. Without change, this would add ground truth objects from other than our subsampled frames and thus, invalidate these studies. Thus, we filter the original ground truth database and ensure that it only contains objects from the actually used frames.

6

**Detection Results:** The **KITTI 3D** dataset and benchmark [12] is one of the first publicly available datasets for 3D object detection and consists of $\sim$ 7.5k LiDAR frames, which were labeled in the front camera view. KITTI360 [24] is an extension of the original KITTI dataset, containing $\sim$ 80k LiDAR frames, offering additional modalities and more exhaustive annotations. For a comparison to other pre-training approaches, we show detection results on the KITTI 3D set while using pre-trained weights from KITTI 3D, KITTI360 and Waymo. For our experiments, we report evaluation metrics based on the moderate difficulty level, utilizing the official $R_{40}$ metric with 40 recall points.

Table 1 shows the detection performance when using different pre-trained weights for the widely-used SECOND [50] and PV-RCNN [32] detectors. Throughout all datasets used for pre-training, MAELi enables strong improvements over detection models trained from scratch and yields top performing detection results. In particular, we outperform the current state-of-the-art, *i.e.* ALSO, in overall mAP, while the performance difference in the car category remains within a very narrow margin, indicating a saturation in performance for this class.

The large-scale **Waymo Open Dataset** [35] contains 798 training sequences and 202 validation sequences. Labels are provided for vehicles, pedestrians and cyclists. We report our results using the official Waymo evaluation protocol, reporting APH at the more challenging LEVEL 2 difficulty. Detailed results including the AP scores are provided in the supplemental material.

We pre-train our weights on the Waymo *train* set and use them to initialize the backbones of SECOND [50], CenterPoint [55] and PV-RCNN [32]. For a fair comparison to [27, 54], we follow the common protocol of using 20% of the entire Waymo *train* set, including vanilla ground truth sampling, to fine-tune the detectors.

Table 2 shows the Waymo results in comparison to the pre-training approaches Occupancy-MAE [27], GCC-3D [23] and ProposalContrast [54]. Pre-training with MAELi demonstrates strong improvements across all detectors. While MAELi pre-trained weights are clearly better suited for SECOND than Occupancy-MAE's, the results for CenterPoint and PV-RCNN are on par with Occupancy-MAE, clearly outperforming GCC-3D or ProposalContrast.

**ONCE** [26] is a comprehensive dataset designed for tasks such as 3D object detection, tracking and motion forecasting. This dataset includes 1 million frames, the majority of which are unlabeled. A primary objective of ONCE is to serve as a foundation for research that leverages on large-scale unlabeled data. In our work, we utilize the official $U_{small}$ subset for pre-training purposes, subsequently fine-tuning our models on the training set and conducting evaluations on the validation set.

| Method | | 3D APH (LEVEL 2) | | | |
|---|---|---|---|---|---|
| | Gain | Overall | Vehicle | Pedestrian | Cyclist |
| SECOND [50] | - | 54.35 | 62.02 | 47.49 | 53.53 |
| + Occ-MAE [27] | +0.75 | 55.10 | 62.34 | 48.79 | 54.17 |
| **+ MAELi** | +2.35 | 56.69 | 63.20 | 50.93 | 55.95 |
| CenterPoint [55] | - | 61.92 | 62.65 | 58.23 | 64.87 |
| + Occ-MAE [27] | +1.31 | 63.23 | 63.53 | 59.62 | 66.53 |
| **+ MAELi** | +1.08 | 63.00 | 63.70 | 59.79 | 65.52 |
| PV-RCNN [32] | - | 56.23 | 64.38 | 45.14 | 59.18 |
| + GCC-3D [23] | +1.95 | 58.18 | 65.10 | 48.02 | 61.43 |
| + PropCont [54] | +3.05 | 59.28 | 65.47 | 49.51 | 62.86 |
| + Occ-MAE [27] | +5.74 | 61.98 | 67.34 | 55.57 | 63.02 |
| **+ MAELi** | +5.92 | 62.15 | 67.34 | 56.32 | 62.79 |

Table 2. Performance comparison on the Waymo *val* set trained on 20% of the Waymo *train* set. We compare different detectors trained from scratch with their pendants utilizing pre-trained weights from GCC-3D, ProposalContrast, Occupancy-MAE and the proposed MAELi.

| Method | Pre-train | Orientation-Aware AP | | | |
|---|---|---|---|---|---|
| | | mAP | Car | Pedestrian | Cyclist |
| SECOND [50] | - | 51.89 | 71.19 | 26.44 | 58.04 |
| + SwAV [5] | $U_{small}$ | 51.96 | 72.71 | 25.13 | 58.05 |
| + DeepCluster [4] | $U_{small}$ | 52.06 | 73.19 | 24.00 | 58.99 |
| + ALSO [3] | $U_{small}$ | 52.68 | 71.73 | 28.16 | 58.13 |
| **+ MAELi** | $U_{small}$ | 57.39 | 75.73 | 34.83 | 61.62 |
| **+ MAELi** | Waymo | 55.84 | 75.86 | 31.03 | 60.65 |
| CenterPoint [55] | - | 64.24 | 75.26 | 51.65 | 65.79 |
| + PropCont [54] | Waymo | 66.24 | 78.00 | 52.56 | 68.17 |
| **+ MAELi** | Waymo | 66.72 | 80.09 | 51.87 | 68.21 |

Table 3. Performance comparison on the ONCE validation set. Our initialization, even when pre-trained on a different dataset, helps outperforming state-of-the-art methods.

Table 3 shows the results comparing MAELi with other pre-training methods on ONCE. Besides our favorable results, this evaluation also demonstrates the strong cross-domain generalization capabilities of our pre-training approach: for SECOND [50], we outperform the state-of-the-art approaches even when using pre-trained weights from Waymo, *i.e.* 55.84 mAP (MAELi pre-trained on Waymo) versus 52.68 (ALSO pre-trained on $U_{small}$). Naturally, pre-training on ONCE itself further improves the results.

**Data Efficiency:** Pre-trained weights play a crucial role during detector initialization. The goal of self-supervised representation learning approaches is to reduce the costly labeling effort. If properly pre-trained, a detector should achieve strong performance with only few annotated samples. We conduct the following experiments to demonstrate the benefits of MAELi for low-data regimes:

On **Waymo**, we follow the evaluation of Proficient-Teachers [53], a state-of-the-art (but semi-supervised) approach which is specifically tailored for data-efficient 3D object detection. The protocol is to pre-train SECOND using the first 399 Waymo *train* sequences and then using different fractions of labeled data from the latter 399 se-

| Fraction | Method | 3D APH (LEVEL 2) | | | | |
|---|---|---|---|---|---|---|
| | | Gain | Overall | Vehicle | Ped. | Cyc. |
| 1% (791 frames) | SECOND [50] | - | 22.25 | 40.02 | 17.45 | 9.29 |
| | + MAELi | +10.79 | 33.05 | 50.11 | 24.65 | 24.38 |
| 5% (3952 frames) | SECOND [50] | - | 34.05 | 52.37 | 22.80 | 26.97 |
| | + ProfTeach [53] | +11.70 | 45.75 | 52.54 | 38.67 | 46.03 |
| | + MAELi | +13.94 | 47.99 | 56.75 | 41.27 | 45.94 |
| 10% (7904 frames) | SECOND [50] | - | 38.23 | 57.46 | 28.15 | 29.07 |
| | + ProfTeach [53] | +12.20 | 50.43 | 56.92 | 43.19 | 51.18 |
| | + MAELi | +13.61 | 51.84 | 59.47 | 45.52 | 50.52 |
| 20% (15808 frames) | SECOND [50] | - | 51.26 | 59.54 | 43.30 | 50.93 |
| | + ProfTeach [53] | +2.90 | 54.16 | 59.36 | 46.97 | 56.15 |
| | + MAELi | +2.75 | 54.01 | 61.21 | 47.63 | 53.18 |

Table 4. Data efficiency comparison for SECOND on the Waymo *val* set. The first 399 Waymo *train* sequences are used for pre-training and different fractions of the latter 399 sequences are used for fine-tuning.

| Fraction | Method | mAP | Car | Ped. | Cyc. |
|---|---|---|---|---|---|
| 20% (743 frames) | PV-RCNN [32] | 66.71 | 82.52 | 53.33 | 64.28 |
| | + PropCont [54] | 68.13 | 82.65 | 55.05 | 66.68 |
| | + MAELi | 69.41 | 82.21 | 56.71 | 69.30 |
| 50% (1856 frames) | PV-RCNN [32] | 69.63 | 82.68 | 57.10 | 69.12 |
| | + PropCont [54] | 71.76 | 82.92 | 59.92 | 72.45 |
| | + MAELi | 70.13 | 83.89 | 56.48 | 70.02 |

Table 5. Data efficiency comparison for PV-RCNN. Following ProposalContrast, we pre-train on Waymo and fine-tune on different fractions of KITTI 3D.

quences to fine-tune the detection heads. For the evaluation on **KITTI 3D**, we adhere to the methodology presented in ProposalContrast [54]. Thus, we pre-train PV-RCNN on the Waymo dataset and subsequently fine-tune it on various fractions of labeled KITTI 3D frames.

Tables 4 and 5 show the data efficiency evaluations for Waymo and KITTI, respectively. Comparisons for other detectors are provided in the supplementary material. Overall, MAELi performs favorably across the different low-data setups. As with any SSRL approach, the initial benefits diminish as the downstream model is trained on larger quantities of annotated samples. However, MAELi performs favorably with fewer annotations, where SSRL is most important since such techniques are designed to provide strong pretrained weights for the crucial initialization phase.

### 4.2. Semantic Segmentation

Since our pre-training methodology does not depend on the downstream task, we further demonstrate its efficacy in the context of semantic segmentation. To this end, we employ regular sparse transpose convolutions in the decoder, which are capable of expanding beyond the active voxels in the encoder, as opposed to using submanifold convolutions (see Section 3.1). Importantly, this maintains the weight dimensionality, enabling the direct utilization of pre-trained weights from both the encoder and decoder during the fine-tuning phase.

| Method | Fraction of labeled samples | | | | |
|---|---|---|---|---|---|
| | 0.1% (17 fr.) | 1% (188) | 10% (1912) | 50% (9560) | 100% (19130) |
| MinkUNet [28] | 30.0 | 46.2 | 57.6 | 61.8 | 62.7 |
| + PointContrast [46] | 32.4 | 47.9 | 59.7 | 62.7 | 63.4 |
| + DepthContrast [58] | 32.5 | 49.0 | 60.3 | 62.9 | 63.9 |
| + SegContrast [28] | 32.3 | 48.9 | 58.7 | 62.1 | 62.3 |
| + ALSO [3] | 35.0 | 50.0 | 60.5 | 63.4 | 63.6 |
| + MAELi | 34.6 | 50.7 | 61.3 | 63.6 | 64.2 |

Table 6. Performance comparison on the SemanticKITTI dataset using a MinkUNet backbone. Results show the mean Intersection-over-Union (mIoU) averaged over 5 runs.

We initialize with the pre-trained weights and employ the same framework and configurations as the state-of-the-art ALSO [3] for subsequent fine-tuning and evaluation procedures. Following ALSO, we use the MinkUNet [9] variant from [28] and report the average performance over 5 runs on the **SemanticKITTI** dataset [2]. This dataset provides semantic labels for each point cloud in the odometry task of the KITTI dataset [12]. SemanticKITTI comprises 22 sequences, 19 classes and ~23k frames for training and validation purposes. Our performance is assessed based on the official evaluation protocol, reported as the mean Intersection-over-Union (mIoU) across all classes.

We follow the evaluation setup of ALSO: pre-training is conducted on the complete dataset and fine-tuning is done on varying fractions of the training set using the training and validation splits from [28]. Our results, stated in Table 6, reveal performance enhancements across almost all considered fractions of the training set (except the highly unreliable setting with merely 17 frames), thereby affirming the general applicability of our MAELi approach in semantic segmentation tasks.

## 5. Conclusion

We proposed MAELi, a self-supervised pre-training approach, carefully designed to adapt to the inherent but subtle properties of large-scale LiDAR point clouds. We were the first to put aspects like occlusion and intrinsic spherical sampling of LiDAR data into the context of SSRL. Our learned representation not only leads to significant improvements in various tasks, but can also be visually verified. Moreover, it can be easily applied to other datasets, with minimal data requirements for fine-tuning. Additionally, our method offers the potential for further investigation on a multi-frame basis. With MAELi, we offer a new method to sustainably reduce the amount of tedious and costly annotation tasks for LiDAR point clouds.

# References

[1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning Representations and Generative Models for 3D Point Clouds. In *Proc. ICML*, 2018. 1

[2] Jens Behley, Martin Garbade, Andres Milioto, Jan Quenzel, Sven Behnke, Cyrill Stachniss, and Jurgen Gall. SemanticKITTI: A Dataset for Semantic Scene Understanding of LiDAR Sequences. In *Proc. ICCV*, 2019. 2, 6, 8

[3] Alexandre Boulch, Corentin Sautier, Björn Michele, Gilles Puy, and Renaud Marlet. ALSO: Automotive Lidar Self-supervision by Occupancy estimation. In *Proc. CVPR*, 2023. 3, 6, 7, 8, 11, 12

[4] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep Clustering for Unsupervised Learning of Visual Features. In *Proc. ECCV*, 2018. 7

[5] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised Learning of Visual Features by Contrasting Cluster Assignments. In *Proc. NeurIPS*, 2020. 7

[6] Angel X. Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. ShapeNet: An Information-Rich 3D Model Repository. *arXiv:1512.03012*, 2015. 1, 3

[7] Anthony Chen, Kevin Zhang, Renrui Zhang, Zihan Wang, Yuheng Lu, Yandong Guo, and Shanghang Zhang. PiMAE: Point Cloud and Image Interactive Masked Autoencoders for 3D Object Detection. In *Proc. CVPR*, 2023. 2

[8] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A Simple Framework for Contrastive Learning of Visual Representations. In *Proc. ICML*, 2020. 2

[9] Christopher Choy, JunYoung Gwak, and Silvio Savarese. 4D Spatio-Temporal ConvNets: Minkowski Convolutional Neural Networks. In *Proc. CVPR*, 2019. 4, 6, 8

[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. CVPR*, 2017. 2

[11] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proc. NAACL: Human Language Technologies*, 2019. 2

[12] Andreas Geiger, Philip Lenz, and Raquel Urtasun. Are we ready for autonomous driving? The KITTI vision benchmark suite. In *Proc. CVPR*, 2012. 2, 6, 7, 8

[13] Rohit Girdhar, David F. Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a Predictable and Generative Vector Representation for Objects. In *Proc. ECCV*, 2016. 2

[14] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3D Semantic Segmentation With Submanifold Sparse Convolutional Networks. In *Proc. CVPR*, 2018. 4

[15] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked Autoencoders Are Scalable Vision Learners. In *Proc. CVPR*, 2022. 2

[16] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum Contrast for Unsupervised Visual Representation Learning. In *Proc. CVPR*, 2020. 2

[17] Georg Hess, Johan Jaxing, Elias Svensson, David Hagerman, Christoffer Petersson, and Lennart Svensson. Masked Autoencoder for Self-Supervised Pre-Training on Lidar Point Clouds. In *WACV Workshop*, 2023. 1, 2

[18] R. Devon Hjelm, Alex Fedorov, Samuel Lavoie-Marchildon, Karan Grewal, Phil Bachman, Adam Trischler, and Yoshua Bengio. Learning deep representations by mutual information estimation and maximization. In *Proc. ICLR*, 2019. 2

[19] Ji Hou, Xiaoliang Dai, Zijian He, Angela Dai, and Matthias Nießner. Mask3D: Pre-Training 2D Vision Transformers by Learning Masked 3D Priors. In *Proc. CVPR*, 2023. 2

[20] Peiyun Hu, Jason Ziglar, David Held, and Deva Ramanan. What You See is What You Get: Exploiting Visibility for 3D Object Detection. In *Proc. CVPR*, 2020. 3

[21] Siyuan Huang, Yichen Xie, Song-Chun Zhu, and Yixin Zhu. Spatio-Temporal Self-Supervised Representation Learning for 3D Point Clouds. In *Proc. ICCV*, 2021. 2, 6

[22] Diederik P. Kingma and Jimmy Ba. Adam: A Method for Stochastic Optimization. In *Proc. ICLR*, 2015. 6

[23] Hanxue Liang, Chenhan Jiang, Dapeng Feng, Xin Chen, Hang Xu, Xiaodan Liang, Wei Zhang, Zhenguo Li, and Luc Van Gool. Exploring Geometry-Aware Contrast and Clustering Harmonization for Self-Supervised 3D Object Detection. In *Proc. ICCV*, 2021. 2, 6, 7, 12

[24] Yiyi Liao, Jun Xie, and Andreas Geiger. KITTI-360: A Novel Dataset and Benchmarks for Urban Scene Understanding in 2D and 3D. *IEEE TPAMI*, 45(3):3292–3310, 2023. 2, 6, 7

[25] Haotian Liu, Mu Cai, and Yong Jae Lee. Masked Discrimination for Self-Supervised Learning on Point Clouds. In *Proc. ECCV*, 2022. 1, 2

[26] Jiageng Mao, Niu Minzhe, ChenHan Jiang, Hanxue Liang, Jingheng Chen, Xiaodan Liang, Yamin Li, Chaoqiang Ye, Wei Zhang, Zhenguo Li, Jie Yu, Chunjing Xu, and Hang Xu. One Million Scenes for Autonomous Driving: ONCE Dataset. In *Proc. NeurIPS*, 2021. 2, 6, 7

[27] Chen Min, Liang Xiao, Dawei Zhao, Yiming Nie, and Bin Dai. Occupancy-MAE: Self-Supervised Pre-Training Large-Scale LiDAR Point Clouds With Masked Occupancy Autoencoders. *Transactions on Intelligent Vehicles*, pages 1–13, 2023. 1, 2, 7, 11, 12

[28] Lucas Nunes, Rodrigo Marcuzzi, Xieyuanli Chen, Jens Behley, and Cyrill Stachniss. SegContrast: 3D Point Cloud Feature Representation Learning Through Self-Supervised Segment Discrimination. *Robotics and Automation Letters*, 7(2):2116–2123, 2022. 2, 8

[29] Lucas Nunes, Louis Wiesmann, Rodrigo Marcuzzi, Xieyuanli Chen, Jens Behley, and Cyrill Stachniss. Temporal Consistent 3D LiDAR Representation Learning for Semantic Perception in Autonomous Driving. In *Proc. CVPR*, 2023. 2

[30] Yatian Pang, Wenxiao Wang, Francis E. H. Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked Autoencoders for Point Cloud Self-supervised Learning. In *Proc. ECCV*, 2022. 1, 2, 5

[31] Xiaoxiao Sheng, Zhiqiang Shen, and Gang Xiao. Contrastive Predictive Autoencoders for Dynamic Point Cloud Self-Supervised Learning. In *Proc. AAAI*, 2023. 2

9

[32] Shaoshuai Shi, Chaoxu Guo, Li Jiang, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. PV-RCNN: Point-Voxel Feature Set Abstraction for 3D Object Detection. In *Proc. CVPR*, 2020. 4, 6, 7, 8, 11, 12

[33] Shaoshuai Shi, Zhe Wang, Jianping Shi, Xiaogang Wang, and Hongsheng Li. From Points to Parts: 3D Object Detection From Point Cloud With Part-Aware and Part-Aggregation Network. *IEEE TPAMI*, 43(08):2647–2664, 2021. 4

[34] Leslie N. Smith. A disciplined approach to neural network hyper-parameters: Part 1 – learning rate, batch size, momentum, and weight decay. *arXiv:1803.09820*, 2018. 6

[35] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, Vijay Vasudevan, Wei Han, Jiquan Ngiam, Hang Zhao, Aleksei Timofeev, Scott Ettinger, Maxim Krivokon, Amy Gao, Aditya Joshi, Yu Zhang, Jonathon Shlens, Zhifeng Chen, and Dragomir Anguelov. Scalability in Perception for Autonomous Driving: Waymo Open Dataset. In *Proc. CVPR*, 2020. 1, 2, 3, 5, 6, 7, 11, 13

[36] OpenPCDet Development Team. OpenPCDet: An Open-source Toolbox for 3D Object Detection from Point Clouds. https://github.com/open-mmlab/OpenPCDet, 2020. 6

[37] Xiaoyu Tian, Haoxi Ran, Yue Wang, and Hang Zhao. GeoMAE: Masked Geometric Target Prediction for Self-Supervised Point Cloud Pre-Training. In *Proc. CVPR*, 2023. 1, 2, 3

[38] Hanchen Wang, Qi Liu, Xiangyu Yue, Joan Lasenby, and Matt J. Kusner. Unsupervised Point Cloud Pre-Training via Occlusion Completion. In *Proc. ICCV*, 2021. 1, 2, 3

[39] Wenguan Wang, Tianfei Zhou, Fisher Yu, Jifeng Dai, Ender Konukoglu, and Luc Van Gool. Exploring Cross-Image Pixel Contrast for Semantic Segmentation. In *Proc. ICCV*, 2021. 2

[40] Yan Wang, Xiangyu Chen, Yurong You, Li Erran Li, Bharath Hariharan, Mark Campbell, Kilian Q. Weinberger, and Wei-Lun Chao. Train in Germany, Test in the USA: Making 3D Object Detectors Generalize. In *Proc. CVPR*, 2020. 1

[41] Xin Wen, Tianyang Li, Zhizhong Han, and Yu-Shen Liu. Point Cloud Completion by Skip-Attention Network With Hierarchical Folding. In *Proc. CVPR*, 2020. 1

[42] Xiaoyang Wu, Xin Wen, Xihui Liu, and Hengshuang Zhao. Masked Scene Contrast: A Scalable Framework for Unsupervised 3D Representation Learning. In *Proc. CVPR*, 2023. 2

[43] Yanhao Wu, Tong Zhang, Wei Ke, Sabine Süsstrunk, and Mathieu Salzmann. Spatiotemporal Self-Supervised Learning for Point Clouds in the Wild. In *Proc. CVPR*, 2023. 2

[44] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3D ShapeNets: A Deep Representation for Volumetric Shapes. In *Proc. CVPR*, 2015. 1, 2, 3

[45] Guangda Xie, Yang Li, Hongquan Qu, and Zaiming Sun. Masked Autoencoder for Pre-Training on 3D Point Cloud Object Detection. *Mathematics*, 10(19):3549, 2022. 2

[46] Saining Xie, Jiatao Gu, Demi Guo, Charles R. Qi, Leonidas Guibas, and Or Litany. PointContrast: Unsupervised Pre-training for 3D Point Cloud Understanding. In *Proc. ECCV*, 2020. 2, 8

[47] Qiangeng Xu, Yiqi Zhong, and Ulrich Neumann. Behind the Curtain: Learning Occluded Shapes for 3D Object Detection. In *Proc. AAAI*, 2022. 3

[48] Runsen Xu, Tai Wang, Wenwei Zhang, Runjian Chen, Jinkun Cao, Jiangmiao Pang, and Dahua Lin. MV-JAR: Masked Voxel Jigsaw and Reconstruction for LiDAR-Based Self-Supervised Pre-Training. In *Proc. CVPR*, 2023. 2

[49] Siming Yan, Zhenpei Yang, Haoxiang Li, Chen Song, Li Guan, Hao Kang, Gang Hua, and Qixing Huang. Implicit Autoencoder for Point-Cloud Self-Supervised Representation Learning. In *Proc. ICCV*, 2023. 2

[50] Yan Yan, Yuxing Mao, and Bo Li. SECOND: Sparsely Embedded Convolutional Detection. *Sensors*, 18(10):3337, 2018. 3, 4, 6, 7, 8, 11, 12, 13, 14, 16

[51] Honghui Yang, Tong He, Jiaheng Liu, Hua Chen, Boxi Wu, Binbin Lin, Xiaofei He, and Wanli Ouyang. GD-MAE: Generative Decoder for MAE Pre-Training on LiDAR Point Clouds. In *Proc. CVPR*, 2023. 2, 3

[52] Juyoung Yang, Pyunghwan Ahn, Doyeon Kim, Haeil Lee, and Junmo Kim. Progressive Seed Generation Auto-Encoder for Unsupervised Point Cloud Learning. In *Proc. ICCV*, 2021. 2

[53] Junbo Yin, Jin Fang, Dingfu Zhou, Liangjun Zhang, Cheng-Zhong Xu, Jianbing Shen, and Wenguan Wang. Semi-supervised 3D Object Detection with Proficient Teachers. In *Proc. ECCV*, 2022. 7, 8

[54] Junbo Yin, Dingfu Zhou, Liangjun Zhang, Jin Fang, Cheng-Zhong Xu, Jianbing Shen, and Wenguan Wang. Proposal-Contrast: Unsupervised Pre-training for LiDAR-based 3D Object Detection. In *Proc. ECCV*, 2022. 2, 6, 7, 8, 12

[55] Tianwei Yin, Xingyi Zhou, and Philipp Krähenbühl. Center-Based 3D Object Detection and Tracking. In *Proc. CVPR*, 2021. 4, 7, 11, 12

[56] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-BERT: Pre-Training 3D Point Cloud Transformers With Masked Point Modeling. In *Proc. CVPR*, 2022. 1, 2, 5

[57] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Point-M2AE: Multi-scale Masked Autoencoders for Hierarchical Point Cloud Pre-training. In *Proc. NeurIPS*, 2022. 2

[58] Zaiwei Zhang, Rohit Girdhar, Armand Joulin, and Ishan Misra. Self-Supervised Pretraining of 3D Features on Any Point-Cloud. In *Proc. ICCV*, 2021. 2, 8

10

# Supplementary Material

This supplementary material presents further details, results and insights into MAELi. We state further results in Appendix A, describe the detailed architecture of our decoder in Appendix B and illustrate the motivation behind spherical masking in Appendix C. Furthermore, we evaluate the impact of different amounts of masked voxels in Appendix D and discuss potential limitations in Appendix E. Finally, we discuss additional insights on reconstruction results and data efficiency in Appendix F.

## A. Additional Results

We provide results for CenterPoint [55] and PV-RCNN [32] to illustrate that our pre-trained initialization significantly enhances these baseline models. Following insights from the main manuscript (Section 4.1), we observe in Table S1 that our MAELi pre-training effectively improves detection performance in a low-data regime where only a limited number of annotated samples are available for fine-tuning. In Table S2, we report AP scores for Waymo Open Dataset (WOD) [35], extending the results from Table 2 in the main manuscript. Additionally, in Table S3, we present our findings on the KITTI 3D dataset using the $R_{11}$ metric, and make comparisons with ALSO [3] and Occupancy-MAE [27].

## B. Sparse Reconstruction Decoder for 3D Object Detection

To describe the architecture of our decoder in detail, we group operations with the same *voxel/tensor stride* into a *block*. In Table S4, we list the different decoder blocks in addition to the preceding BEV encoder (summarized as single entry) and the required *reshaping+sampling* step to transform the dense feature representation back to a sparse 3D tensor.

Each block comprises an upsampling step using *generative transposed convolution* and a pruning step via $1 \times 1 \times 1$ *submanifold sparse convolution*. The operations are listed in Table S5.

## C. Spherical Masking - Illustration

As discussed in the main manuscript (Section 3.3), spherical masking reduces the angular resolution in azimuth and inclination by subsampling the LiDAR's range image. Figure S1 illustrates the effect of this sampling on the Li-DAR's range image. We sample objects as if they were located at a larger distance. Since nearby objects are more densely sensed by the LiDAR, we have more knowledge about the actual occupancy and thus, can induce a stronger

self-supervision signal. This helps to improve the model's ability to generalize to objects located farther away.

## D. Ablation Study

**Analyzing Pipeline Components:** We perform various experiments to investigate specific aspects of our pipeline, such as assessing the influence of our *distance weighting* for empty voxels, evaluating our *LiDAR-aware reconstruction* objective, and comparing the effectiveness of the *voxel-* and *spherical masking* strategies.

Therefore, we pre-train according to the *Data Efficiency* protocol depicted in the main manuscript (Section 4.1) and fine-tune a SECOND [50] model on 1% of the latter 399 sequences of the Waymo *train* set. To disable *distance weighting*, we set $w_j^{D,\mathbf{s}} = 1$ for all *empty* voxels $\mathbf{v}_j^{D,\mathbf{s}}$. To disable our *LiDAR-aware reconstruction*, we additionally consider all *unknown* voxels as *empty*.

We state the results in Table S6 on *Vehicle* LEVEL 2 across the distance ranges [0m, 30m], [30m, 50m) and [50m, $+\infty$). Our *LiDAR-aware reconstruction* objective improves the overall results across all ranges. While *voxel masking* naturally has a nearly equal impact over all ranges, our *spherical masking* is especially beneficial for the range [30m, 50m) with a gain of 1.87AP and 2.01APH. The overall lower impact on the far distance range (above 50m) is also reasonable, since at this distance only very few points are sampled on the same object.

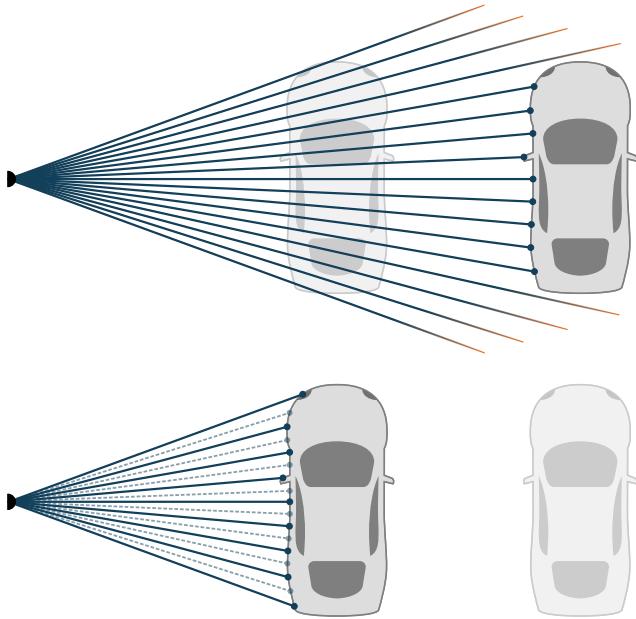**Masking:** We evaluated our *voxel masking* for different



Figure S1. Spherical masking reduces the angular resolution of the LiDAR (bottom). The resulting sampling is thus similar to objects that are farther away (top).

| Fraction | Method | 3D AP/APH (LEVEL 2) | | | | | | | | | |
| | | Gain | | Overall | | Vehicle | | Pedestrian | | Cyclist | |
| | | AP | APH | AP | APH | AP | APH | AP | APH | AP | APH |
| 1%<br>(791 frames) | Centerpoint [55] | - | - | 39.64 | 36.50 | 41.01 | 40.32 | 40.01 | 32.63 | 37.90 | 36.55 |
| | + MAELi | +9.29 | +8.75 | 48.93 | 45.25 | 49.99 | 49.24 | 51.92 | 43.07 | 44.89 | 43.43 |
| | PV-RCNN [32] | - | - | 43.93 | 30.72 | 51.34 | 48.70 | 41.59 | 20.35 | 38.86 | 23.11 |
| | + MAELi | +7.53 | +4.89 | 51.46 | 35.61 | 56.24 | 55.38 | 49.41 | 25.32 | 48.73 | 26.14 |
| 5%<br>(3952 frames) | Centerpoint [55] | - | - | 53.91 | 51.16 | 53.04 | 52.45 | 52.73 | 46.51 | 55.96 | 54.53 |
| | + MAELi | +4.49 | +4.21 | 58.40 | 55.37 | 57.62 | 57.01 | 59.01 | 51.83 | 58.57 | 57.27 |
| | PV-RCNN [32] | - | - | 56.98 | 38.98 | 61.66 | 60.86 | 53.28 | 27.15 | 56.00 | 28.92 |
| | + MAELi | +1.64 | +1.39 | 58.62 | 40.37 | 62.77 | 62.04 | 57.07 | 29.05 | 56.02 | 30.02 |
| 10%<br>(7904 frames) | Centerpoint [55] | - | - | 58.09 | 55.41 | 56.95 | 56.40 | 56.97 | 50.90 | 60.35 | 58.94 |
| | + MAELi | +3.26 | +3.07 | 61.35 | 58.48 | 59.93 | 59.36 | 62.06 | 55.30 | 62.06 | 60.78 |
| | PV-RCNN [32] | - | - | 60.09 | 41.89 | 63.73 | 63.05 | 57.32 | 30.09 | 59.23 | 32.53 |
| | + MAELi | +1.19 | -0.05 | 61.28 | 41.84 | 64.63 | 63.99 | 59.82 | 30.90 | 59.40 | 30.62 |
| 20%<br>(15808 frames) | Centerpoint [55] | - | - | 61.81 | 59.15 | 60.59 | 60.06 | 61.19 | 55.03 | 63.64 | 62.36 |
| | + MAELi | +0.98 | +0.90 | 62.79 | 60.05 | 61.79 | 61.23 | 63.47 | 57.04 | 63.11 | 61.87 |
| | PV-RCNN [32] | - | - | 62.15 | 42.99 | 65.01 | 64.35 | 60.40 | 30.30 | 61.05 | 34.32 |
| | + MAELi | +0.49 | +7.21 | 62.65 | 50.20 | 65.45 | 64.85 | 61.54 | 35.52 | 60.95 | 50.24 |

Table S1. Quantitative results of our pre-training on Centerpoint and PV-RCNN on the Waymo *val* set. For each detector, we report the results of training from scratch (upper row) and the improved results utilizing a MAELi-pre-trained initialization (lower row), respectively. We use the first 399 sequences of the Waymo *train* set for pre-training and different fractions of the second 399 sequences for fine-tuning.

| Method | 3D AP/APH (LEVEL 2) | | | | | | | | | |
| | Gain | | Overall | | Vehicle | | Pedestrian | | Cyclist | |
| | AP | APH | AP | APH | AP | APH | AP | APH | AP | APH |
| SECOND [50] | - | - | 58.26 | 54.35 | 62.58 | 62.02 | 57.22 | 47.49 | 54.97 | 53.53 |
| + Occ-MAE [27] | +0.85 | +0.75 | 59.11 | 55.10 | 62.67 | 62.34 | 59.03 | 48.79 | 55.62 | 54.17 |
| + MAELi | +2.32 | +2.35 | 60.57 | 56.69 | 63.75 | 63.20 | 60.71 | 50.93 | 57.26 | 55.95 |
| CenterPoint [55] | - | - | 64.51 | 61.92 | 63.16 | 62.65 | 64.27 | 58.23 | 66.11 | 64.87 |
| + Occ-MAE [27] | +1.35 | +1.31 | 65.86 | 63.23 | 64.05 | 63.53 | 65.78 | 59.62 | 67.76 | 66.53 |
| + MAELi | +1.09 | +1.08 | 65.60 | 63.00 | 64.22 | 63.70 | 65.93 | 59.79 | 66.66 | 65.52 |
| PV-RCNN [32] | - | - | 59.84 | 56.23 | 64.99 | 64.38 | 53.80 | 45.14 | 60.72 | 59.18 |
| + GCC-3D [23] | +1.46 | +1.95 | 61.30 | 58.18 | 65.65 | 65.10 | 55.54 | 48.02 | 62.72 | 61.43 |
| + PropCont [54] | +2.78 | +3.05 | 62.62 | 59.28 | 66.04 | 65.47 | 57.58 | 49.51 | 64.23 | 62.86 |
| + Occ-MAE [27] | +5.99 | +5.74 | 65.82 | 61.98 | 67.94 | 67.34 | 64.91 | 55.57 | 64.62 | 63.02 |
| + MAELi | +5.88 | +5.92 | 65.72 | 62.15 | 67.90 | 67.34 | 65.14 | 56.32 | 64.13 | 62.79 |

Table S2. Performance comparison on the Waymo *val* set trained on 20% of the Waymo *train* set including AP scores. We compare different detectors trained from scratch with their pendants utilizing pre-trained weights from GCC-3D, ProposalContrast, Occupancy-MAE and the proposed MAELi.

| Method | Pre-train | mAP | Car | Pedestrian | Cyclist |
| --- | --- | --- | --- | --- | --- |
| SECOND [50] | - | 66.25 | 78.62 | 52.98 | 67.15 |
| + Occ-MAE [27] | KITTI 3D | 66.71 | 78.90 | 53.14 | 68.08 |
| + ALSO [3] | nuScenes | 67.29 | 78.65 | 55.17 | 68.05 |
| + ALSO [3] | KITTI 3D | 66.86 | 78.78 | 53.57 | 68.22 |
| + ALSO [3] | KITTI360 | 67.40 | 78.63 | 54.23 | 69.35 |
| + MAELi | Waymo | 68.31 | 78.44 | 55.72 | 70.78 |
| + MAELi | KITTI 3D | 67.51 | 78.20 | 55.48 | 68.86 |
| + MAELi | KITTI360 | 68.74 | 78.44 | 56.00 | 71.79 |
| PV-RCNN [32] | - | 70.66 | 83.61 | 57.90 | 70.47 |
| + Occ-MAE [27] | KITTI 3D | 71.73 | 83.82 | 59.37 | 71.99 |
| + ALSO [3] | nuScenes | 72.20 | 83.77 | 58.49 | 74.35 |
| + ALSO [3] | KITTI 3D | 71.96 | 83.67 | 58.48 | 73.74 |
| + ALSO [3] | KITTI360 | 72.69 | 83.39 | 60.83 | 73.85 |
| + MAELi | Waymo | 71.79 | 83.38 | 58.53 | 73.45 |
| + MAELi | KITTI 3D | 70.70 | 79.22 | 60.02 | 72.87 |
| + MAELi | KITTI360 | 73.03 | 83.99 | 62.43 | 72.67 |

Table S3. Quantitative results of our pre-training on SECOND and PV-RCNN on the KITTI 3D *val* set using the $R_{11}$ metric.

| Description | # Channels | Voxel/Tensor Stride | Spatial Dimension |
| --- | --- | --- | --- |
| Output BEV Encoder | 512 | - | 188 × 188 |
| Reshaping + Sampling | 256 | 8 × 8 × 16 | 188 × 188 × 2 |
| DBlock 1 | 64 | 8 × 8 × 8 | 188 × 188 × 5 |
| DBlock 2 | 64 | 4 × 4 × 4 | 376 × 376 × 11 |
| DBlock 3 | 32 | 2 × 2 × 2 | 752 × 752 × 21 |
| DBlock 4 | 16 | 1 × 1 × 1 | 1504 × 1504 × 41 |

Table S4. Architecture of our decoder. We state the number of channels, the voxel stride and the maximum spatial dimension for the Waymo Open Dataset *after* each block. Stride and spatial dimensions are depicted in the format $x \times y \times z$. Each decoder block *inverts* one downsampling step from the sparse 3D encoder, eventually resulting in the original voxel stride.

amounts of voxels. We maintain the training and evaluation scheme from above and vary the amount of kept voxels.

The results are shown in Table S7. Keeping $60\%$ of the voxels leads to the best overall results, eventually used for all other experiments with MAELi. However, in combination with *spherical masking* significantly fewer points than this

| Operation | Kernel Size | Stride |
|---|---|---|
| Generative Transposed Convolution | $2 \times 2 \times 2^{\dagger}$ | $2 \times 2 \times 2^{\dagger}$ |
| Batch Norm | - | - |
| ReLU | - | - |
| Submanifold Sparse Convolution | $3 \times 3 \times 3$ | $1 \times 1 \times 1$ |
| Batch Norm | - | - |
| ReLU | - | - |
| Submanifold Sparse Convolution | $1 \times 1 \times 1$ | $1 \times 1 \times 1$ |
| Pruning | - | - |

Table S5. Structure of each decoder block. We additionally state the operation's kernel size and stride, each in the format $x \times y \times z$. The upper part depicts the upsampling and feature transformation. The lower part uses the final feature representation from above and decides via classification whether a voxels is pruned or not. $^{\dagger}$These values deviate for DBlock 1, where it has a kernel size of $1 \times 1 \times 3$ and a stride of $1 \times 1 \times 2$ to invert the encoder's respective downsampling step.

fraction actually remain. To get an estimate, we evaluated the amount of effectively used points over 10000 iterations, resulting in a fraction of $15.57\%$ on average.

## E. Limitations

Even though our sparse decoder allows for a memory efficient reconstruction, the amount of reconstructed voxels is obviously constrained by the available compute infrastructure. Especially during the first iterations of our pre-training, while not sufficiently trained, some samples may lead to an uncontrolled reconstruction. In order to regulate the amount of reconstructed voxels and to avoid training breakdowns, we introduce two limiting factors. First, we estimate an average ground plane for each dataset and prune all reconstructed voxels that are $0.1m$ below this plane, as these generally do not contribute valuable information. Second, we introduce a threshold for the maximum amount of reconstructed voxels to ensure that we do not run into memory issues. If an upsampling step would generate more voxels than this limit, we randomly prune before the upsampling. For these pruned voxels, simply no loss is induced, which only slightly *delays* the training effect for these rare cases. For the detection experiments, we set the maximum number of total voxels to 6 million, which are easily processable, *e.g.* on an NVIDIA® GeForce® RTX 3090 GPU. We counted only 62 limit exceedances within the first 10k iterations of a random experiment.

## F. Additional Insights

**Reconstruction Capabilities:** In Figure S2, we visualize the reconstruction capabilities of our *LiDAR-aware loss* on a full point cloud. It encourages the network to fill up gaps in the wall and reconstruct the occluded areas of cars.

Figure S3 highlights that reconstruction outcomes can vary across different objects, with some objects such as the less frequent trucks presenting greater challenges. However, the MAELi model demonstrates an advanced understanding of object semantics, enabling it to complete objects beyond the visible LiDAR input point cloud.

In Figure S4, we show the individual layers of two different reconstructed cars. We can see that our pre-training approach indeed encourages the model to go beyond the sampled LiDAR surface, reconstructing the entire car, also showing some hints for the correct placement of tires. Furthermore, especially parts of the interior, which are often surrounded by glass and thus, sometimes traversed by LiDAR beams, are seemingly left out during the reconstruction.

**AP vs APH Data Efficiency:** In Figure S5, we plot the data efficiency results on Waymo [35] using our MAELi-based pre-training on the SECOND [50] detector. All three classes benefit from our pre-training (solid lines) compared to the vanilla version trained from scratch (dotted lines). With little data, the vanilla detector especially struggles to estimate the heading, which can be clearly seen for the smaller, less represented classes *Pedestrian* and *Cyclist*. However, utilizing our initialization, a proper estimation and significant detection improvements are possible already, early on, with only few annotated data samples.
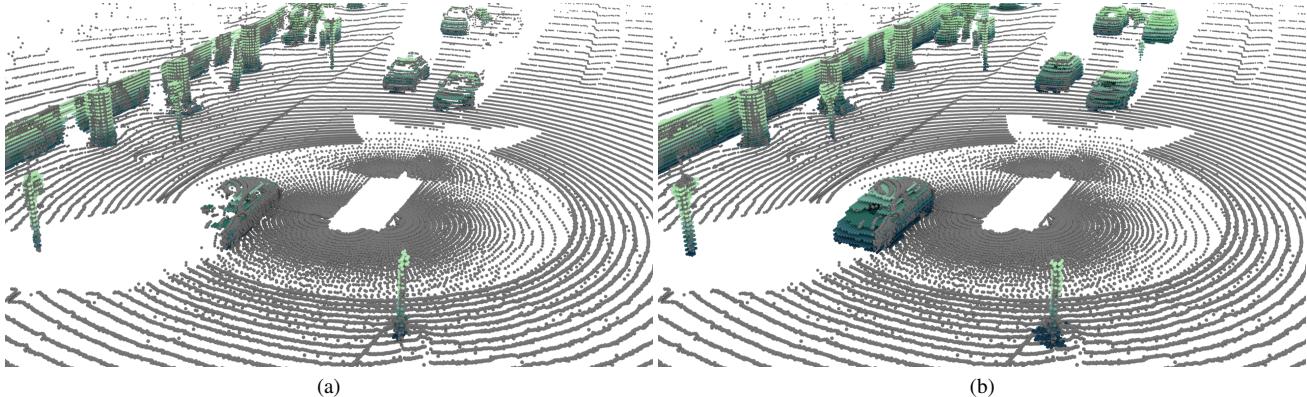
Figure S2. Completion results (a) without and (b) with our *LiDAR-aware reconstruction* on a full point cloud (gray). For visualization purposes, we color-coded the output points by their $z$-coordinate and removed the reconstructed ground plane.

| Method | 3D AP/APH (LEVEL 2) | | | | | | | |
| | Overall | | [0m, 30m) | | [30m, 50m) | | [50m, +inf) | |
| | AP | APH | AP | APH | AP | APH | AP | APH |
|---|---|---|---|---|---|---|---|---|
| **MAELi** | 51.05 | 50.11 | 80.22 | 79.37 | 48.86 | 47.64 | 21.09 | 19.98 |
| w/o DW | 50.91 | 49.85 | 79.95 | 79.04 | 48.62 | 47.20 | 21.47 | 20.31 |
| w/o LAR | 50.05 | 48.87 | 79.33 | 78.20 | 47.77 | 46.27 | 20.66 | 19.56 |
| w/o VM | 48.99 | 47.86 | 78.45 | 77.46 | 46.43 | 45.05 | 19.47 | 18.23 |
| w/o SM | 49.90 | 48.85 | 79.50 | 78.54 | 46.99 | 45.63 | 20.53 | 19.37 |
| Baseline | 41.64 | 40.02 | 72.59 | 70.79 | 37.09 | 34.86 | 13.14 | 11.96 |

Table S6. Impact of the components of MAELi evaluated on the Waymo *val* set for *Vehicle*. We disable *distance weighting* (w/o DW), *LiDAR-aware reconstruction* (w/o LAR), *voxel masking* (w/o VM) and *spherical masking* (w/o SM). We use the first 399 sequences of the Waymo *train* set for pre-training and 1% of the second 399 sequences for fine-tuning. We utilize a SECOND [50] model and state a version trained from scratch as baseline.

| Fraction Voxels | 3D AP/APH (LEVEL 2) | | | | | | | |
| | Overall | | Vehicle | | Pedestrian | | Cyclist | |
| | AP | APH | AP | APH | AP | APH | AP | APH |
|---|---|---|---|---|---|---|---|---|
| 0.8 | 45.34 | 32.84 | 50.26 | 49.10 | 48.03 | 24.33 | 37.74 | 25.09 |
| 0.7 | 44.91 | 32.10 | 50.35 | 49.19 | 47.54 | 24.29 | 36.83 | 22.83 |
| 0.6 | 46.01 | 33.05 | 51.05 | 50.11 | 48.13 | 24.65 | 38.86 | 24.38 |
| 0.5 | 45.60 | 32.27 | 50.87 | 49.79 | 47.08 | 23.72 | 38.86 | 23.31 |
| 0.4 | 45.79 | 31.85 | 50.36 | 49.24 | 48.27 | 24.50 | 38.74 | 21.81 |
| Baseline | 31.09 | 22.25 | 41.64 | 40.02 | 33.39 | 17.45 | 18.24 | 9.29 |

Table S7. Impact of different amounts of voxels kept during *voxel masking* evaluated on the Waymo *val* set for *Vehicle*. We use the first 399 sequences of the Waymo *train* set for pre-training and 1% of the second 399 sequences for fine-tuning. We utilize a SECOND [50] model and state a version trained from scratch as baseline.
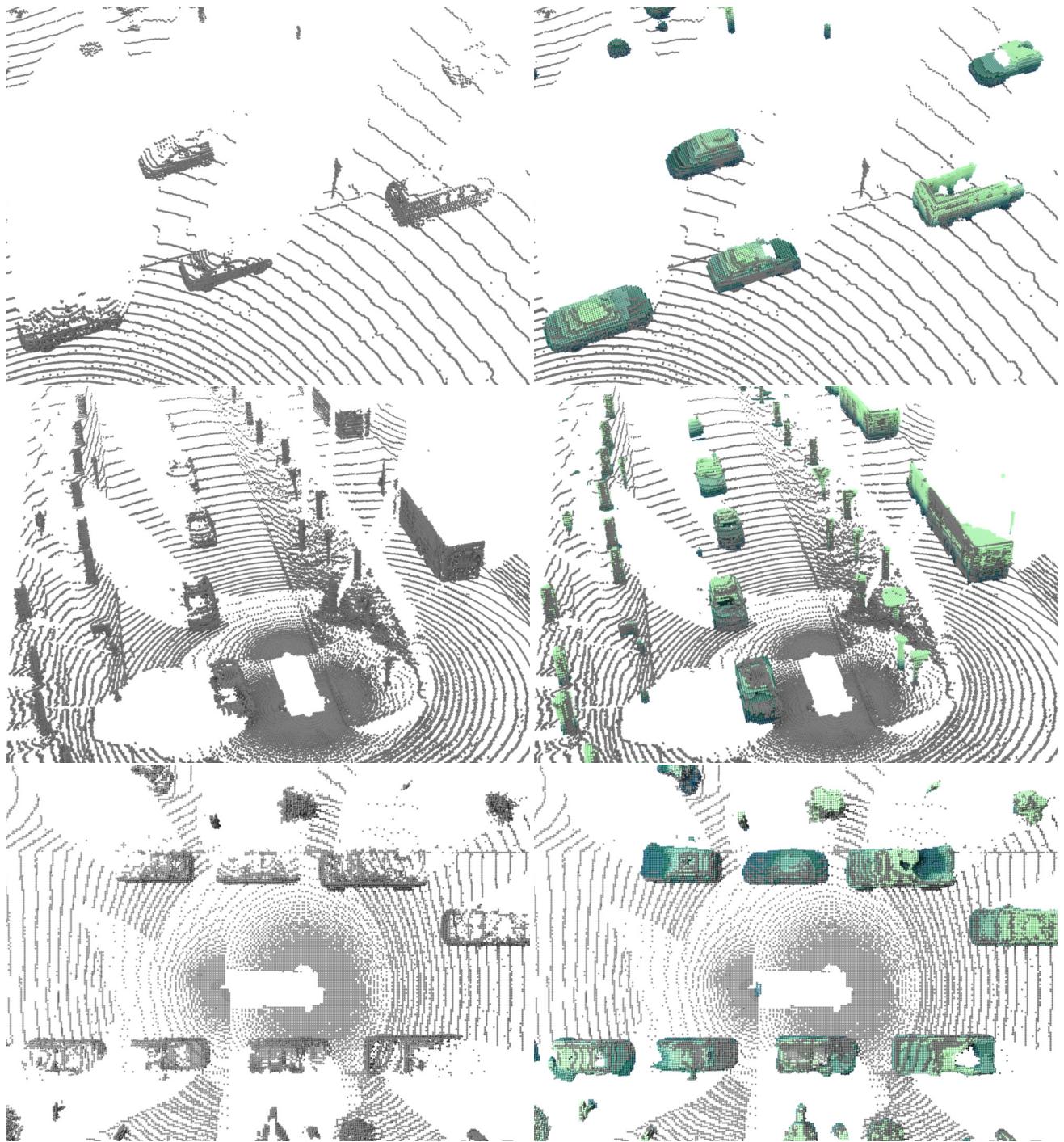
Figure S3. Further reconstruction results. We show the input point cloud on the left and the completed point cloud on the right. MAELi's reconstruction exhibits imperfections when reconstructing objects that are sparsely sampled or less frequent, such as trucks. However, it shows an apparent understanding of traffic scenes beyond a LiDAR's 2.5D sampling, *e.g.* by symmetrically completing occluded parts of cars and poles. For visualization purposes, we color-coded the output points by their $z$-coordinate and removed the reconstructed ground plane.
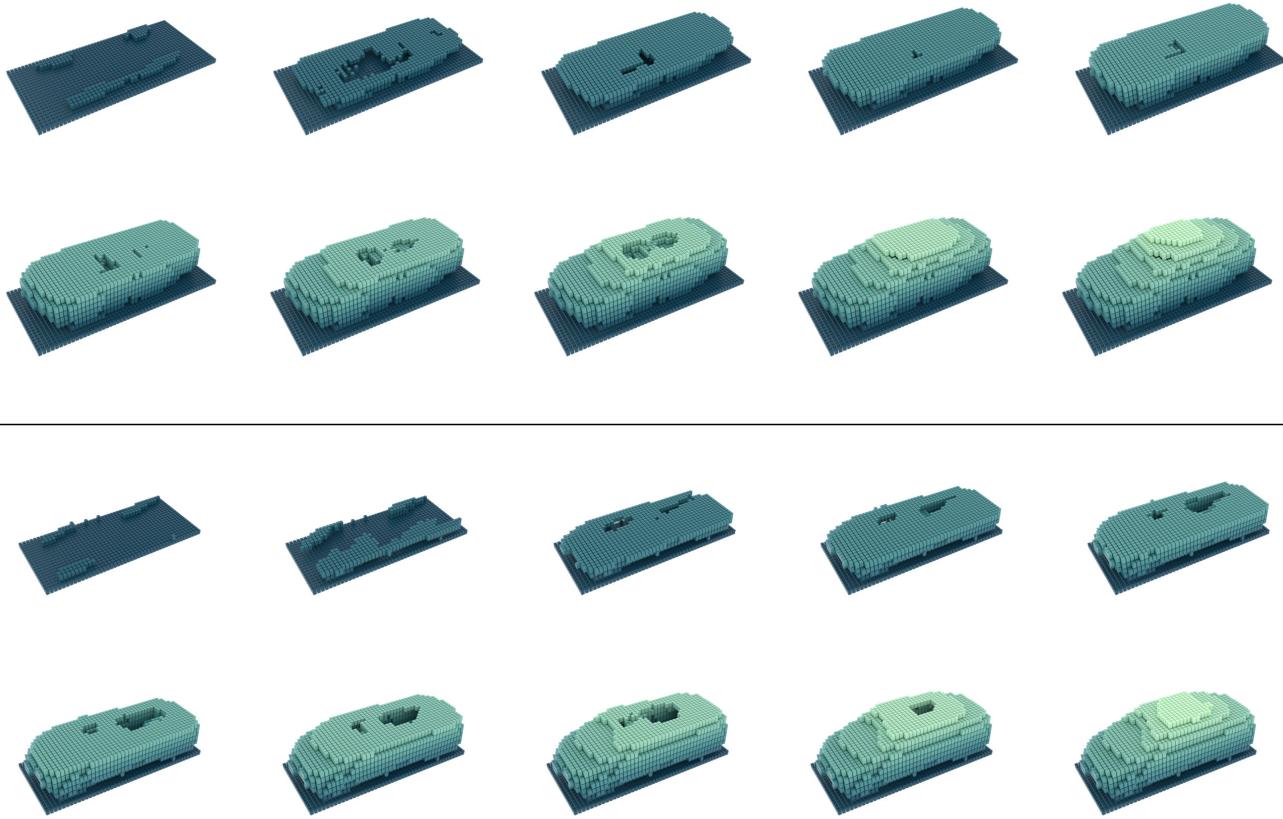
Figure S4. Different layers of two reconstructed cars. We observed that the reconstructed cars are often hollow. There are visible tendencies to leave parts of the interior free.
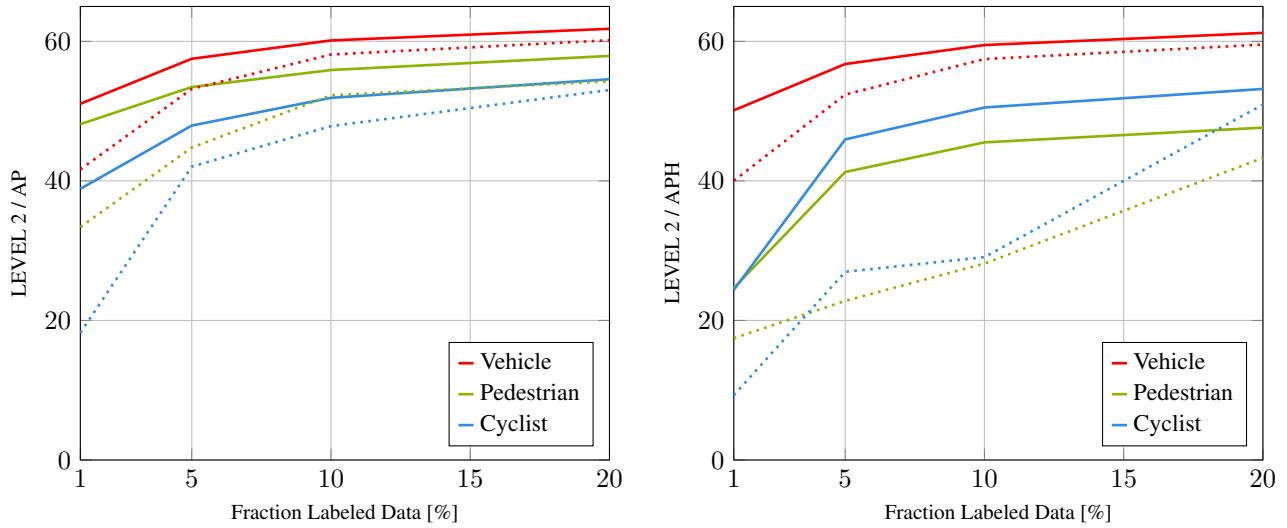


Figure S5. Results of our pre-training on SECOND [50] on the Waymo *val* set, using different amounts of labeled data. We use the first 399 sequences of the Waymo *train* set without any labels for pre-training and different fractions of the latter 399 sequences for fine-tuning. The *solid lines* are the results utilizing our pre-training with MAELi. The *dotted lines* denote the version trained from scratch.