

Self-Supervised Pretraining of 3D Features on *any* Point-Cloud

Zaiwei Zhang^{1,2*}Rohit Girdhar¹Armand Joulin¹Ishan Misra¹¹Facebook AI Research²The University of Texas at Austin

Abstract

Pretraining on large labeled datasets is a prerequisite to achieve good performance in many computer vision tasks like 2D object recognition, video classification etc. However, pretraining is not widely used for 3D recognition tasks where state-of-the-art methods train models from scratch. A primary reason is the lack of large annotated datasets because 3D data is both difficult to acquire and time consuming to label. We present a simple self-supervised pretraining method that can work with any 3D data — single or multi-view, indoor or outdoor, acquired by varied sensors, without 3D registration. We pretrain standard point cloud and voxel based model architectures, and show that joint pretraining further improves performance. We evaluate our models on 9 benchmarks for object detection, semantic segmentation, and object classification, where they achieve state-of-the-art results and can outperform supervised pretraining. We set a new state-of-the-art for object detection on ScanNet (69.0% mAP) and SUNRGBD (63.5% mAP). Our pretrained models are label efficient and improve performance for classes with few examples. Code and models are available at <https://github.com/facebookresearch/DepthContrast>

1. Introduction

Pretraining visual features on large labeled datasets is a pre-requisite to achieve good performance when access to annotations is limited [27, 48, 55, 90]. More recently, self-supervised pretraining has become a popular alternative to supervised pretraining especially for tasks where annotations are time-consuming, such as detection and segmentation in images [9, 36, 37, 59, 97] or tracking in videos [41]. In 3D computer vision, single-view depth scans are easy to acquire while reconstructed 3D scenes and annotations are difficult to obtain. Reconstructing a 3D scene requires registering and aligning multiple depth maps captured in a static environment, which may fail easily if there exists fast camera motion or odometry drift [16]. The resulting 3D scene is a point cloud composed of thousands of 3D

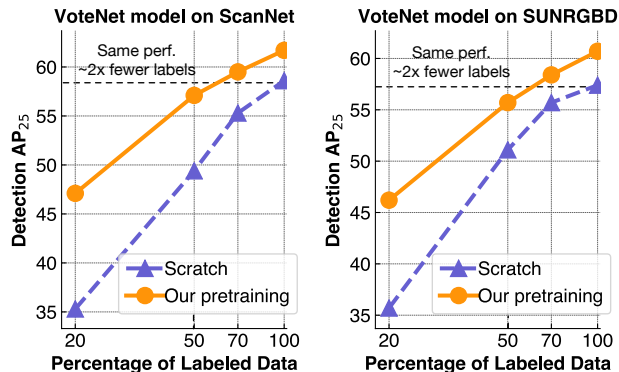


Figure 1: Label-efficiency of our self-supervised pretraining. We finetune detection models from scratch or using our pretraining as initialization. Our pretraining which uses unlabeled single-view 3D data, outperforms training from scratch, and achieves the same detection performance with about half the detection labels.

points that need to be annotated individually in the case of segmentation or by groups for detection, which is time-consuming and laborious. For example, it takes around 22 minutes to annotate a single scene in ScanNet [18]. Thus, 3D data acquisition and annotation both require significantly more effort than images or videos, and there is no existing weak supervision, like image tags, to shortcut the process [45, 55, 90]. This cumbersome annotation process results in a lack of large annotated 3D datasets. However, consumer-grade depth sensors have become more easily accessible and simpler to use, *e.g.*, in phones [24, 76, 86], leading to large quantities of single-view raw depth maps. These depth maps can be leveraged to pretrain 3D features, but, there is surprisingly little work that can be applied.

Pretraining on depth maps faces several challenges: first, the absence of any supervision at scale requires the use of a self-supervised signal. Second, the absence of alignments between depth maps excludes methods using multi-view constraints, such as finding correspondences between views [109]. Finally, the representation of a 3D scene varies depending on the application—segmentation typically uses a voxel representation [17], whereas detection uses point-clouds [67]—and therefore the ideal pretraining method should be easily applicable to any 3D representation.

*Work done during an internship at Facebook.

In this paper, we introduce a simple contrastive framework, **DepthContrast**, to simultaneously learn different representations of any 3D data. Our approach is based on the Instance Discrimination method by Wu *et al.* [107], applying to depth maps. We side-step the need of registered point clouds or correspondences, by considering each depth map as an instance and discriminating between them, even if they come from the same scene. For 3D representations with different architectures, we jointly learn features by considering the different representations as data augmentations that are processed with their associated networks [97].

Our contributions can be summarized as follows:

- We show that single view 3D depth scans can be used to learn powerful feature representations using self-supervised learning.
- We show that *joint* training of different input representations like points and voxels is important for learning good representations, and a naive application of contrastive learning may not yield good results.
- Our method is applicable across different model architectures, indoor/outdoor 3D data, single/multi-view 3D data. We also show that it can be used to pretrain high capacity 3D architectures which otherwise overfit on tasks like detection and segmentation.
- We show performance improvements over *nine* downstream tasks, and set a new state-of-the-art for *two* object detection tasks (ScanNet and SUNRGBD). Our models are efficient few-shot learners.

2. Related Work

Our method builds on the work from the self-supervised learning literature, with 3D data as an application. In this section, we give an overview of the recent advances in both self-supervision and 3D representations.

Self-supervised learning for images. Self-supervised learning is a well studied problem in machine learning and computer vision [56, 63, 72, 75, 99]. There are many classes of methods for learning representations - clustering [7, 8, 43], GANs [20, 58], pretext tasks [19, 62, 101] *etc.* Recent advances [9, 13, 30, 36, 37, 49, 59, 98] have shown that self-supervised pretraining is a viable alternative to supervised pretraining for 2D recognition tasks. Our work builds upon contrastive learning [34, 64] where models are trained to discriminate between each instance [21] with no explicit classifier [107]. These instance discrimination methods can be extended to multiple modalities [60, 65, 97]. Our method extends the work of Wu *et al.* [107] to multiple 3D input formats following Tian *et al.* [97] using a momentum encoder [36] instead of a memory bank.

Self-supervised learning for 3D data. Most methods on self-supervised learning focus on single 3D object representation with different applications to reconstruction, classifi-

cation or part segmentation [1, 2, 25, 33, 35, 44, 51, 77, 103, 114]. Recently, Xie *et al.* [109] proposed a self-supervised method to build representations of scene level point clouds. Their method relies on the complete 3D reconstruction of a scene with point-wise correspondences between the different views of a point cloud. These point-wise correspondences requires post-processing the data by registering the different depth maps into a single 3D scene. Their method can only be applied to static scenes that have been registered, which greatly limits the applications of their work. We show a simple self-supervised method that learns state-of-the-art representations from *single-view* 3D data, and can also be applied to multi-view data.

Representations of 3D scenes. There are multiple ways to represent 3D information in different vectorized forms such as point-clouds, voxels or meshes. Point-cloud based models [69, 71] are widely used in classification and segmentation tasks [6, 40, 47, 69, 71, 94, 102, 104, 110, 111], 3D reconstruction [23, 93, 114] and 3D object detection [66, 67, 80, 102, 116, 116, 118]. Since many 3D sensors acquire data in terms of 3D points, point clouds are a convenient input for deep networks. However, since using convolution operations on point-clouds directly is difficult [31, 69], voxelized data is another popular input representation. 3D convolutional models [3, 17, 28, 31, 38, 52, 73, 88, 95] are widely used in 3D scene understanding [32, 85, 113, 120]. There are also efforts to combine different 3D input representations [29, 79, 100, 117–119]. In this work, we propose to jointly pretrain two architectures for points and voxels, that are PointNet++ [71] for points and Sparse Convolution based U-Net [17] for voxels.

3D transfer tasks and datasets. We use shape classification, scene segmentation, and object detection as the recognition tasks for transfer learning. Shape classification techniques [11, 57, 69–71, 89] are widely evaluated on the ModelNet [106] dataset, which we use. It contains synthetic 3D data and each sample contains exactly one object. We also evaluate on complete 3D scenes using the more general 3D scene understanding task. Scene-centric datasets can be broadly divided into indoor scenes [5, 10, 18, 39, 61, 66, 78, 84, 87, 108], and outdoor (self-driving focussed) scenes [26, 74, 91]. We use these datasets and evaluate the performance of our methods on the indoor detection [12, 22, 67, 68, 118], scene segmentation [17, 71, 95, 105, 112], and outdoor detection tasks [15, 50, 79–81, 113, 115].

3. Approach

We develop a scalable pretraining method, DepthContrast, for 3D representations that uses unprocessed single-view or multi-view depth maps without human annotations. Our method, illustrated in Fig 2, is based on the instance discrimination framework of Wu *et al.* [107] with a momentum encoder [36]. DepthContrast learns 3D representations across multiple 3D input formats like points and voxels, and

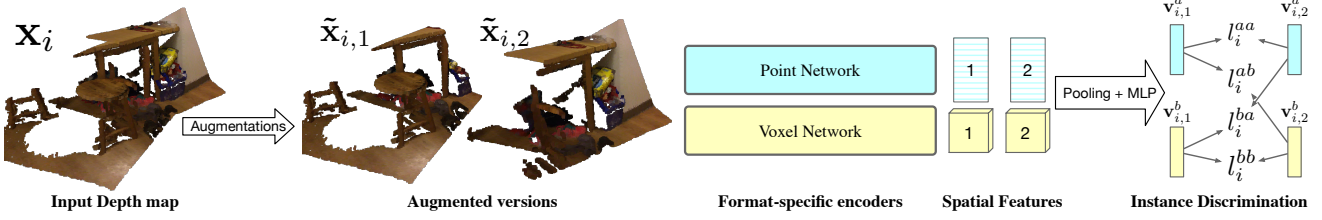


Figure 2: Approach Overview. We propose DepthContrast - a simple 3D representation learning method that uses large amounts of unprocessed single/multi-view depth maps. Given a depth map we construct two augmented versions using data augmentation and represent them with different input formats (point coordinates and voxels). We use format-specific encoders to get spatial features which are pooled and projected to obtain global features \mathbf{v} . The global features are used to setup an instance discrimination task and pretrain the encoders.

across different 3D architectures by using an extension of contrastive learning to multiple data formats [97].

3.1. Instance Discrimination

Given a dataset $\mathcal{D} = \{\mathbf{X}\}_{i=1}^N$ containing N samples \mathbf{X} , we wish to learn a function $g(\mathbf{X})$ that produces useful representations $\mathbf{v} = g(\mathbf{X})$ of the input sample. Our method uses 3D data where \mathbf{X} can be represented by point coordinates or voxels¹. We apply a data augmentation t sampled randomly from a large set of augmentations \mathcal{T} , to obtain an augmented sample $\tilde{\mathbf{X}} = t(\mathbf{X})$. The augmented sample is input to a deep network g that extracts unit-norm global features $\mathbf{v} = g(\tilde{\mathbf{X}})$ by pooling over the 3D spatial coordinates. We setup an instance discrimination problem where the features $\mathbf{v}_{i,1}$ and $\mathbf{v}_{i,2}$ obtained from two data augmented versions of sample i must be similar to each other, and different from features \mathbf{v}_j obtained using K other (negative) samples j in the dataset. We use a contrastive loss [34, 64, 83] to achieve this goal:

$$l_i = -\log \frac{\exp(\mathbf{v}_{i,1}^\top \mathbf{v}_{i,2} / \tau)}{\exp(\mathbf{v}_{i,1}^\top \mathbf{v}_{i,2} / \tau) + \sum_{j \neq i}^K \exp(\mathbf{v}_{i,1}^\top \mathbf{v}_j / \tau)}, \quad (1)$$

where τ is the temperature that controls the smoothness of the softmax distribution. This loss encourages features from different augmentations of the same scene to be similar, while being dissimilar to features of other scenes. Thus, it learns features that focus on discriminative regions of a scene that make it different from other scenes in the dataset. **Minimal assumptions on input data.** Our method, by design, makes minimal assumptions about the input \mathbf{X} , *i.e.*, it is an unprocessed single-view depth map. It does not require careful sampling of overlapping multi-view 3D inputs [109] or object centric depth maps [35, 44]. These minimal assumptions enable us to learn from large scale single-view 3D depth maps in § 4, indoor and outdoor 3D depth maps obtained from a variety of sensors without relying on 3D calibration in § 5.3.

Momentum encoder. As using a large number of negatives is important for contrastive learning [13, 36, 59, 107], we

¹Points in a depth map are a set, but for simplicity we denote them as a matrix. Our method does not rely on any specific ordering of the points.

use the method of He *et al.* [36] where the features of the other augmentation $\mathbf{v}_{i,2}$ and negative samples \mathbf{v}_j in Eq 1 are obtained using a momentum encoder and a queue respectively. This allows us to use a large number K of negative samples without increasing the training batch size.

3.2. Extension to Multiple 3D Input Formats

Multiple input formats are commonly used to represent 3D data - point clouds, voxels, meshes *etc.* and have their specific deep learning architectures and applications. Our self-supervised method can be naturally extended to accommodate these input formats and architectures. For each input format f , we denote the corresponding input sample as \mathbf{X}^f , the format-specific encoder network as g^f , and the extracted feature as \mathbf{v}^f . Extending Eq 1, we can minimize a single objective that performs instance discrimination within and across input formats a, b :

$$l_i^{ab} = -\log \frac{\exp(\mathbf{v}_{i,1}^{a\top} \mathbf{v}_{i,2}^b / \tau)}{\exp(\mathbf{v}_{i,1}^{a\top} \mathbf{v}_{i,2}^b / \tau) + \sum_{j \neq i}^K \exp(\mathbf{v}_{i,1}^{a\top} \mathbf{v}_j^b / \tau)}. \quad (2)$$

When the input formats a, b are identical, this objective reduces to the *within format* loss of Eq 1, and when $a \neq b$ this objective aligns the feature representations $\mathbf{v}^f = g^f(\mathbf{X}^f)$ obtained *across formats* f using different network architectures g^f . As illustrated in Fig 2, we use two popular input formats - point clouds and voxels, and train these format-specific models with a single joint loss function

$$L_i = \underbrace{l_i^{ab} + l_i^{ba}}_{\text{across format}} + \underbrace{l_i^{aa} + l_i^{bb}}_{\text{within format}}. \quad (3)$$

Similar techniques have been explored in the context of different modalities of data, *e.g.*, color and grayscale images [97], audio and video [60, 65] *etc.* While these methods use different modalities, our extension uses the *same 3D data* and only changes the input format.

3.3. Model Architecture

We describe the model architecture used for our input format-specific encoders. Once we obtain the augmented data sample $\tilde{\mathbf{X}}$, we use the XYZ coordinates for point input. We follow [17, 109] to voxelize the augmented sample

$\tilde{\mathbf{X}}$ and use the 3D voxel occupancy grid and RGB values for the voxel model. Thus, both encoders operate on the same input 3D data, and differ only in the way the input is represented. We provide the full, layer-wise architecture details in the supplemental material.

Point input. We use PointNet++ [67] as the backbone network which takes as input the XYZ coordinates of the 3D data. PointNet++ employs a U-Net structure which has four layers of feature extraction and down-sampling, and two layers of feature aggregation and up-sampling. Our network takes as input 20K points from the scene represented as a $20K \times 3$ matrix of XYZ coordinates. The network’s final layer produces C dimensional per-point features for 1024 points after aggregation. We obtain the scene level 256 dimensional feature \mathbf{v} in Eq 2 by global max pooling to these last layer features, followed by a two layer MLP as in [13] and L2 normalization. We increase the size C of the output feature in the last two upsampling (fp) layers from 256 in [67] to 512. This improves the performance of the baselines, *e.g.*, training from scratch, as well as our method and we use this improved architecture throughout the paper.

Voxel input. We use a sparse convolution U-Net model [17] as the backbone for the voxel 3D input. The network takes a 3D occupancy grid as the input representation of the 3D data. Our U-Net consists of four layers of feature extraction and pooling and four layers of feature aggregation and up-sampling. We use a voxel size of $5cm$ to voxelize the input data and following [17], use the voxel occupancy grid and the RGB values as input to the model. To obtain the scene level 256 dimensional feature \mathbf{v} (Eq 2), similar to the point input, we apply global max-pooling to the last layer feature, followed by a two layer MLP and L2 normalization.

Using global feature representations. Our instance discrimination formulation learns global feature representations \mathbf{v} from the input 3D data \mathbf{X} , *i.e.*, the feature representation \mathbf{v} is obtained by pooling over the 3D coordinates of the input. On the other hand many downstream 3D tasks require local per-point predictions like segmentation. We show in § 4 that despite using global features, our pretraining benefits such prediction tasks.

3.4. Data Augmentation for 3D

Data augmentation is as an essential component of our framework. We first adopt standard pointcloud data augmentation methods proposed in [67], which are random point up/down sampling, random flip in xy axis, and random rotation. However, after adding these methods, it is still easy for the network to distinguish different training instances. Thus, we add two new data augmentation methods: random cuboid and random drop patches. Inspired by the random crop in 2D images [92], we define a random cuboid augmentation that extracts random cuboids from the input point cloud. Cuboids are sampled using a random scale $[0.5, 1.0]$ of the original scene, and a random aspect

Dataset	Stats	Task	Gain of DepthContrast
Self-supervised Pretraining			
ScanNet-vid [18]	190K single-view depth maps (Indoor)		
Redwood-vid [16]	370K single-view depth maps (Indoor/Outdoor)		
Transfer tasks			
ScanNet [18]	1.2K train, 312 val (Indoor)	Det.	+3.6% mAP
		Seg.	+0.9% mIOU[†]
SUNRGBD [84]	5.2K train, 5K val (Indoor)	Det.	+3.3% mAP
S3DIS [4]	199 train, 67 val (Indoor)	Det.	+12.1% mAP
		Seg.	+2.4% mIOU
Synthia [74]	19.8K train, 1.8K val (Synth.)	Seg.	+2.4% mIOU
Matterport3D [10]	1.4K train, 232 val (Indoor)	Det.	+3.9% mAP
ModelNet [106]	9.8K train, 2.4K val (Synth.)	Cls.	+3.1% Acc[†]

Det.: Object Detection, Seg: Semantic Segmentation

Cls: Classification, Synth.: Synthetic, [†]Results in supplemental.

Table 1: Pretraining datasets and transfer tasks used in this paper. We use two different pretraining datasets without post-processing like 3D registration, camera calibration. We use 8 different transfer tasks for evaluation where our DepthContrast pretraining gives consistent gains (last column) over scratch pretraining. Additionally, we show evaluation results on LiDAR data in § 5.3.

ratio $[0.75, 1.0]$. We also drop (erase) cuboids to force the network learn local geometric features. The dropped cuboid is randomly cropped with 0.2 of the scene scale. The performance boost from each augmentation is analyzed in § 5. For voxelized inputs, in addition to all the point augmentations, we use the augmentations from [17].

3.5. Implementation Details

We use 130K negatives for contrastive learning in Eq 3 and a momentum of 0.9 for the momentum encoder following [36]. As noted in § 3.3, we follow Chen *et al.* [13] and use an additional non-linear projection (two layer MLP) and L2 normalization after the network g to obtain the features \mathbf{v} . The features \mathbf{v} are 128 dimensional and we use a temperature value of 0.1 while computing the non-parametric softmax in Eq 1. We use a standard SGD optimizer with momentum 0.9, cosine learning rate scheduler [53] starting from 0.12 to 0.00012 and train the model for 1000 epochs with a batch size of 1024.

4. Experiments

We evaluate DepthContrast pretraining by transfer learning, *i.e.*, fine-tuning on downstream tasks and datasets. As Table 1 shows, we use a diverse set of 3D understanding tasks like object classification, semantic segmentation, and object detection. We first study a single input 3D format and a single network architecture in § 4.2 showing that DepthContrast’s performance improves with large data and higher capacity models, and benefits finetuning with limited labels. Finally, in § 4.3, we show the benefits of our pre-training across different 3D input formats (points and voxels).

Pretraining Details. We use single-view depth map videos

Initialization	ScanNet	SUNRGBD	Matterport3D	S3DIS
Scratch	58.6	57.4	38.8	31.2
Supervised	-	59.1 (+1.7)	41.7 (+2.9)	48.5 (+17.3)
DepthContrast (Ours)	61.3 (+2.7)	60.4 (+3.0)	41.9 (+3.1)	43.3 (+12.1)
PointContrast [109]	59.2	57.5	-	-

Table 2: Detection AP₂₅ using VoteNet [67]. We evaluate different pretrained models - random initialization, supervised VoteNet on ScanNet, and our self-supervised DepthContrast using the point input format. DepthContrast outperforms the scratch model on all benchmarks and is better than the detection-specific supervised pretraining on two datasets.

Method	ScanNet		SUNRGBD	
	AP ₂₅	AP ₅₀	AP ₂₅	AP ₅₀
F-PointNet [68]	54.0	-	-	-
VoteNet [67]	58.6	33.5	57.7	32.9
H3DNet [118]	67.2	48.1	60.1	39.0
HGNet [12]	61.3	34.4	61.6	34.4
3D-MPA [22]	64.2	49.2	-	-
PointContrast (VoteNet) [109]	59.2	38.0	57.5	34.8
DepthContrast (VoteNet)	64.0	42.9	61.6	35.5
DepthContrast (H3DNet)	69.0	50.0	63.5	43.4

Table 3: Transfer using state-of-the-art detection frameworks. We use our pretrained model (PointNet++ 3× on Redwood-vid +ScanNet-vid) and transfer it using two state-of-the-art detection frameworks - H3DNet [118] and VoteNet [67]. Our DepthContrast pretraining outperforms all prior work and sets a new state-of-the-art on both ScanNet and SUNRGBD detection datasets.

from the popular ScanNet [18] dataset and term it as ScanNet-vid. ScanNet-vid contains about 2.5 million RGB-D scans for more than 1500 indoor scenes. Following the train/val split from [67], we extract around 190K RGB-D scans (one frame every 15 frames) from about 1200 video sequences in the train set. We do not use camera calibration or 3D registration methods and operate directly on single-view depth maps. We use our data augmentation described in § 3.4 and use the training objectives from § 3. Additional details are provided in § 3.5 and the supplemental material.

4.1. Transfer Datasets and Tasks

We evaluate our pretrained model by transfer learning and finetune it on different downstream datasets and tasks summarized in Table 1. We use diverse downstream datasets - full scenes/object centric; using different 3D sensors; single/multi-view; real/synthetic; indoor/outdoor. On these diverse datasets, we use three major tasks, which are classification, semantic segmentation, and object detection. These tasks test different aspects of the pretrained model - while object detection and semantic segmentation use local features, classification is performed on global features.

4.2. Pretraining with Point Input Format

We pretrain a PointNet++ model using the instance discrimination objective in Eq 1 on the single-view depth maps from ScanNet-vid. We study the transfer performance of the pretrained model on object detection using the VoteNet [67] framework that uses a PointNet++ backbone network. In Table 2 we report the detection results by finetuning the VoteNet model with different backbone initializations. We use the implementation of [67] for finetuning and report the detection performance using the mean Average Precision at IoU=0.25 (AP₂₅) metric. Training from scratch or random initialization is standard practice in VoteNet [67] and serves as a baseline for comparing other pretraining methods. As a supervised pretraining baseline, we use the VoteNet model trained on ScanNet detection. Since the supervised baseline is pretrained specifically on object detection, it serves as a strong baseline.

Scratch training provides competitive results on the larger detection datasets like ScanNet and SUNRGBD [42, 82, 84, 108], however, its performance on the smaller S3DIS dataset is low. In comparison, supervised pretraining provides large gains in the detection performance across all datasets. DepthContrast outperforms training from scratch on all the four datasets, and **improves performance by 12.1% mAP** on the small S3DIS dataset that has only 200 labeled training samples. We further analyze label efficiency of our model in § 4.2.4. Interestingly, despite using no labels during pretraining, DepthContrast is better than the detection-specific supervised pretraining for two datasets (SUNRGBD and Matterport3D). Our method also outperforms recent work [109] by significant margins.

4.2.1 Training Higher Capacity Models

We use DepthContrast for training higher capacity models. We follow the practice in 2D self-supervised learning [49] and increase the capacity of the PointNet++ model by multiplying the channel width of all the layers by {2, 3, 4}. We pretrain all models on the ScanNet-vid dataset and measure their transfer performance in Fig 3. Training large models from scratch provides some benefit, but quickly leads to reduced or plateauing performance. We observe overfitting on the small datasets like S3DIS where increasing the model capacity does not improve performance. On the other hand, our self-supervised pretraining on ScanNet-vid reduces this overfitting and performance improves or stays the same for larger models. This suggests that *pretraining is crucial* for training large 3D detection models.

4.2.2 Using More Pretraining Data

We increase the pretraining data by using readily available single-view 3D data from the Redwood-vid dataset [16]. Redwood-vid contains over 23 million depth scans from RGB-D videos taken in both indoor and outdoor settings.

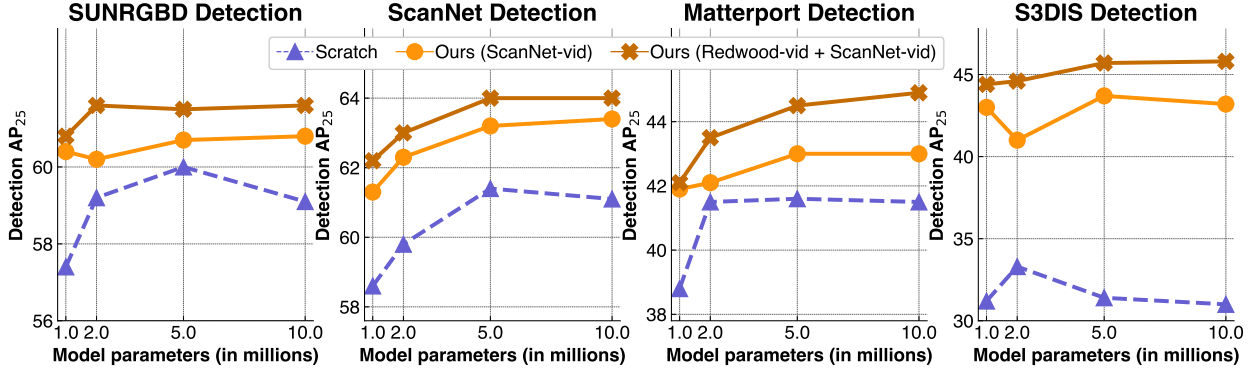


Figure 3: Scaling the model size and pretraining data. We increase the model capacity of the PointNet++ model by increasing the width by $\{2\times, 3\times, 4\times\}$. When training from scratch, increasing the model capacity increases the performance but ultimately leads to overfitting. Overfitting is more pronounced on small datasets like S3DIS. Our DepthContrast pretraining on ScanNet-vid improves the performance for larger models and reduces overfitting. We increase the pretraining data by combining the readily available single-view depth maps from ScanNet-vid and Redwood-vid. DepthContrast’s performance improves significantly when using both large data and large models.

As this dataset is extremely large, we use a subset of 2500 video sequences consisting of 10 categories and extract 370K RGB-D scans. More importantly, the Redwood-vid dataset does *not contain* camera extrinsic parameters and thus cannot be registered to get a multi-view dataset which is a necessity for prior self-supervised methods [109].

Combining the Redwood-vid and ScanNet-vid datasets allows us to triple our pretraining data. We pretrain all models on this combined dataset and report their performance (AP_{25}) in Fig 3. DepthContrast’s performance improves with both model capacity and number of pretraining samples across all four detection datasets. The higher capacity models show a larger improvement in performance particularly on the smaller S3DIS dataset. This suggests that DepthContrast can leverage large amounts of single-view 3D data to obtain better and higher capacity 3D models.

4.2.3 State-of-the-art Detection Frameworks

We use two state-of-the-art detection frameworks - H3DNet [118] and VoteNet [67] and study the benefit of using our pretrained model. We use our PointNet++ $3\times$ model pretrained on the combined Redwood-vid and ScanNet-vid dataset and transfer it using these detection frameworks. The detection results in Table 3 show that our pretrained model achieves state-of-the-art performance on SUNRGBD and ScanNet. In particular, as the gains are larger on stricter mAP at IoU=0.5, our pretrained models result in detection models that are better at localization.

4.2.4 Label Efficiency of Pretrained Models

Pretraining allows models to be finetuned with small amount of labeled data. In Table 2, we observe that small labeled datasets benefit more from pretraining. We study the label efficiency of DepthContrast pretrained models by varying the amount of labeled data used for finetuning.

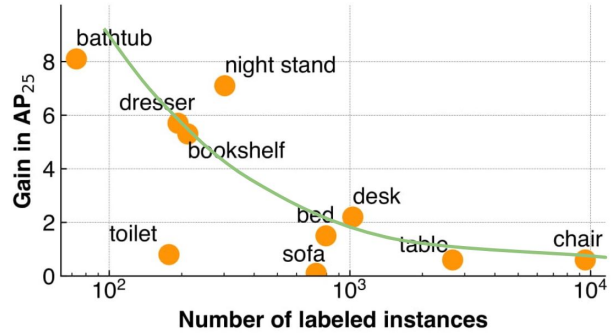


Figure 4: Pretraining benefits long tail classes. We analyze the gain of our pretraining across different classes for SUNRGBD object detection. The training data has a long tailed distribution where the least frequent classes occur $50\times$ less than the most frequent classes. Our pretraining improves performance for classes with fewer labeled instances by 4 – 8%. (Trending line in green.)

While varying the data, we draw 3 independent samples and report average results. We use the PointNet++ models pretrained on ScanNet-vid (§ 4.2) and report the detection performance in Fig 1. DepthContrast pretraining provides large gains in performance at every setting. On both the ScanNet and SUNRGBD datasets, our model with just 50% *samples gets the same performance* as training from scratch with the full dataset. When using 20% samples for finetuning, our pretrained models provide a **gain of over 10% mAP**. This shows that our pretraining is label efficient and can improve performance especially on tasks with limited supervision.

Does pretraining benefit tail classes? 3D detection datasets like SUNRGBD and ScanNet exhibit a long tailed distribution where many ‘tail’ classes have few training instances. In the SUNRGBD dataset, the ‘tail’ classes like bathtub, toilet, dresser have less than 200 training

Loss	Point Transfer		Voxel Transfer	
	SUNRGBD	ScanNet	S3DIS	Synthia
Scratch	57.4	58.6	68.2	78.9
Within Format only	60.4 (+3.0)	61.3 (+1.7)	66.5 (-2.7)	80.1 (+1.2)
Across format only	60.0 (+2.6)	61.1 (+2.5)	69.9 (+1.7)	81.2 (+2.3)
Both (Ours)	60.7 (+3.3)	62.2 (+3.6)	70.6 (+2.4)	81.3 (+2.4)
PointContrast [109]	57.5	59.2	70.9	83.1

Table 4: Multiple input formats. We study the importance of training 3D representations jointly using multiple input formats - points and voxels. We vary the within format and across format loss terms in Eq 3. We report detection mAP@0.25 on the point transfer tasks and segmentation mIOU for the voxel transfer tasks. A naive within format instance discrimination loss (second row) can give worse performance than finetuning from scratch. We observe that performing instance discrimination across the input formats (third row) greatly improves over the within format loss term. Our joint loss gives the best performance.

instances, while classes like chair have over 9000 instances. Fig 4 shows the gain of our pretrained model over the scratch model across object classes on the SUNRGBD dataset. Our pretraining improves the performance of classes with fewer instances, *i.e.*, the tail classes, by 4 – 8% AP. This suggests that DepthContrast pretraining can partially address the long tailed label distributions of current 3D scene understanding benchmarks.

4.3. Pretraining with Multiple Input Formats

We pretrain DepthContrast using both the point and voxel input formats and use two format-specific encoders - PointNet++ for points and UNet for voxels. As explained in Eq 3, when using multiple 3D input formats, we can define two loss terms - a within format loss and an across format loss. To analyze which of these loss terms matter for pretraining, we consider three variants - (1) *Within format* which independently trains format-specific models for each input format and is a straightforward application of instance discrimination to 3D; (2) *Across format* which trains the format-specific models jointly using the second term of Eq 3; (3) *Ours* which trains the format-specific models jointly using our combined loss function. We evaluate the pretrained models by transfer learning. As in § 4.2, we finetune the pretrained point input format PointNet++ models on SUNRGBD and ScanNet detection using the VoteNet framework. We finetune the voxel UNet models on segmentation using the framework from Spatio-Temporal Segmentation [17] which uses a UNet backbone network. The results are summarized in Table 4.

Compared to training from scratch, the within format pretraining only provides a benefit for the point input format PointNet++ models. For the voxel models, this pretraining does not improve consistently over training from scratch, which is in line with observations from recent work [109]. This shows that a naive application of instance discrimination to 3D representation learning *may not* yield good pretrained models. The across format loss improves perfor-

Task	VoteNet [67]	+Rand. Cuboid	+Rand. Drop
ModelNet Linear (Accuracy)	80.6	85.4	85.0
SUNRGBD Detection (mAP)	58.6	59.5	60.7

Table 5: Data augmentation. We vary the data augmentation used for pretraining DepthContrast point models and report their transfer performance. The standard data augmentation used in supervised learning (VoteNet) is not sufficient to learn good self-supervised models. Our proposed Random Cuboid and Random Drop augmentations improve performance.

mance for both the point and voxel models, suggesting the benefit of using multiple input formats. Our proposed loss function that combines both the within and across format losses provides the best transfer performance. The gains are particularly significant on the voxel format model which improves by 4% **over the within format loss**. In the supplemental material, we show that this benefit of joint training over the within format loss also holds across different pretraining data and architectures. Although our method only uses single-view unprocessed depth scans, our results on the voxel transfer tasks are comparable to the recent PointContrast method [109] that uses multi-view point clouds and pointwise correspondences. We note that our UNet architecture is different from [109] since their architecture underfit on our self-supervised pretraining task. We provide results with their architecture in the supplement.

5. Analysis

In this section we present a series of experiments designed to understand DepthContrast better. We first pretrain point format (PointNet++) models on the ScanNet-vid dataset following the settings from § 4.2. We use two transfer tasks for evaluation - (1) object detection on SUNRGBD using VoteNet [67] where we finetune the full model and test the quality of the pretraining; (2) object classification on ModelNet where we keep the model fixed and only train linear classifiers on fixed features, thus testing the quality of the learned representations [30, 49]. Finally, we also evaluate DepthContrast’s generalizability to outdoor 3D data.

5.1. Importance of Data Augmentation

Data augmentations play an important role for self-supervised representation learning and have been studied extensively in the case of 2D images [9, 14, 59, 97, 98]. However, the impact of data augmentation for 3D representation learning is less well understood. Thus, we analyze the effect of our proposed augmentations from § 3.4 on transfer performance. We train different DepthContrast point models with the same training setup and only vary the data augmentation used. Our results are summarized in Table 5.

The widely used VoteNet [67] augmentations from su-

Task	Scratch	Pretraining		
		ScanNet (Multi-view)	ScanNet-vid (Single-view)	Redwood-vid (Single-view)
ModelNet Linear (Accuracy)	50.7	85.1	85.0	86.4
SUNRGBD Detection (mAP)	57.4	60.5	60.7	60.4

Table 6: Single-view or multi-view 3D data. We study whether our pretraining is sensitive to single-view or multi-view data. We use ScanNet and ScanNet-vid which are multi-view and single-view versions of the same dataset [18], and Redwood-vid [16] which is a single-view only 3D dataset. Our pretrained model is robust to 3D preprocessing, and using single-view or multi-view data gives similar performance.

pervised learning perform worse than our proposed augmentations. Our augmentations lead to both a better feature representation: a gain of 5% accuracy on ModelNet classification, and a better pretrained model: 2% mAP on SUNRGBD detection. In our experiments, we consistently observe gains from our improved data augmentation on all the downstream tasks from § 4 which underscores the importance of designing good data augmentation.

5.2. Impact of Single-view or Multi-view 3D Data

Our self-supervised method does not make assumptions on the input data and can use single-view depth maps without 3D preprocessing as input. We study whether pretraining on reconstructed multi-view 3D scenes impacts the downstream performance. We use the ScanNet [18] dataset which contains multi-view 3D data obtained by 3D registration of the ScanNet-vid depth maps. As another single-view dataset, we pretrain on the Redwood-vid dataset from § 4.2.2. We pretrain DepthContrast point models on these datasets and compare their performance by transfer learning in Table 6.

The transfer performance is similar when models are pretrained on ScanNet-vid or ScanNet. Since ScanNet-vid and ScanNet only differ in the 3D preprocessing involved, the result suggests that DepthContrast is not sensitive to single-view or multi-view input data. This is not surprising given that our objective does not rely on multi-view information. Pretraining on the single-view Redwood-vid dataset also gives similar performance suggesting that DepthContrast is robust to different data distributions during pretraining. All the DepthContrast models outperform the scratch model.

5.3. Generalization to Outdoor LiDAR data

We test DepthContrast’s generalization to outdoor LiDAR data by pretraining on the Waymo Open Dataset [91] where we extract 79K single-view scans from the videos. We use the same data augmentation parameters from § 3.4 and only modify random cuboid to work on the full scale of the Z (depth) dimension of the scene. We use the standard LiDAR-specific model architectures as our format-specific encoders - PointnetMSG [80] for point clouds and Spconv-

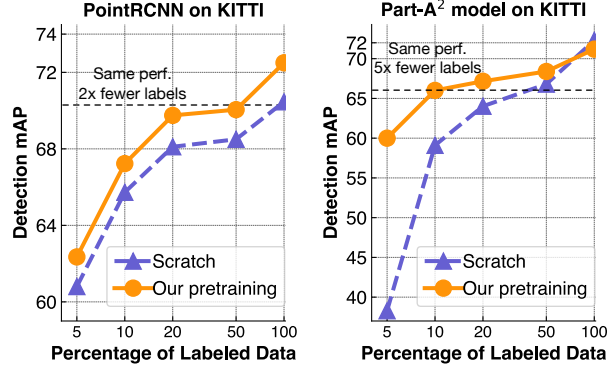


Figure 5: Using outdoor LiDAR data. We finetune detection models from scratch or using our pretraining and report mAP (with 40 recall positions) on the cyclist class at moderate difficulty level of the KITTI val split. Our models are pretrained using unlabeled outdoor data from the Waymo dataset and outperform scratch training using either point (left) or voxel (right) inputs.

UNet [81] for voxels. Similar to § 3.5, we obtain features from these models after global max pooling and a two layer MLP. The models are optimized jointly with Eq 3 using both within and across-format losses. For transfer learning, we use the standard KITTI [26] object detection benchmark, and PointRCNN [80] and Part-A² [81] for down-stream models. We report results on the cyclist class since it has fewer examples in the training set compared to the other classes. We provide results for other classes and finetuning details in the supplemental material. Similar to § 4.2.4, while varying the fraction of pretraining data, we report average performance across 3 independent samplings of the data. Fig 5 shows that our pretrained models outperform training from scratch especially when finetuning on fewer training samples. For Spconv-Unet, we achieve a **20% gain with 5% of labeled data**. This suggests that DepthContrast pretraining generalizes across multiple input formats, and our proposed data augmentation generalizes to different depth sensors and scene types.

6. Conclusion

We propose DepthContrast- an easy to implement self-supervised method that works across model architectures, input data formats, indoor/outdoor 3D, single/multi-view 3D data. DepthContrast pretrains high capacity models for 3D recognition tasks, and leverages large scale 3D data that may not have multi-view information. We show state-of-the-art performance on detection and segmentation benchmarks, outperforming all prior work on detection. We provide crucial insights that make our simple implementation work well - training jointly with multiple input data formats, and designing a generalizable data augmentation scheme. We hope DepthContrast helps future work in 3D self-supervised learning.

References

- [1] Idan Achituve, Haggai Maron, and Gal Chechik. Self-supervised learning for domain adaptation on point-clouds. *arXiv preprint arXiv:2003.12641*, 2020.
- [2] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas Guibas. Learning representations and generative models for 3d point clouds. In *Proceedings of the International Conference on Machine Learning (ICML)*, pages 40–49. PMLR, 2018.
- [3] Andrew Adams, Jongmin Baek, and Myers Abraham Davis. Fast high-dimensional filtering using the permutohedral lattice. In *Computer Graphics Forum*, volume 29, pages 753–762. Wiley Online Library, 2010.
- [4] Iro Armeni, Sasha Sax, Amir Roshan Zamir, and Silvio Savarese. Joint 2d-3d-semantic data for indoor scene understanding. *CoRR*, abs/1702.01105, 2017.
- [5] 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 *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1534–1543, 2016.
- [6] Matan Atzmon, Haggai Maron, and Yaron Lipman. Point convolutional neural networks by extension operators. *arXiv preprint arXiv:1803.10091*, 2018.
- [7] Piotr Bojanowski and Armand Joulin. Unsupervised learning by predicting noise. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [8] Mathilde Caron, Piotr Bojanowski, Armand Joulin, and Matthijs Douze. Deep clustering for unsupervised learning of visual features. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018.
- [9] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. In *NeurIPS*, 2020.
- [10] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Niessner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments. *International Conference on 3D Vision (3DV)*, 2017. Matterport license available at <http://kaldir.vc.in.tum.de/matterport/MP.TOS.pdf>.
- [11] Angel X Chang, Thomas Funkhouser, Leonidas Guibas, Pat Hanrahan, Qixing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, et al. Shapenet: An information-rich 3d model repository. *arXiv preprint arXiv:1512.03012*, 2015.
- [12] Jintai Chen, Biwen Lei, Qingyu Song, Haochao Ying, Danny Z Chen, and Jian Wu. A hierarchical graph network for 3d object detection on point clouds. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 392–401, 2020.
- [13] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. *arXiv preprint arXiv:2002.05709*, 2020.
- [14] Xinlei Chen, Haoqi Fan, Ross Girshick, and Kaiming He. Improved baselines with momentum contrastive learning. *arXiv preprint arXiv:2003.04297*, 2020.
- [15] Yilun Chen, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Fast point r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9775–9784, 2019.
- [16] Sungjoon Choi, Qian-Yi Zhou, Stephen Miller, and Vladlen Koltun. A large dataset of object scans. *arXiv:1602.02481*, 2016.
- [17] 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.
- [18] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Niessner. ScanNet: Richly-annotated 3d reconstructions of indoor scenes. In *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, July 2017.
- [19] Carl Doersch, Abhinav Gupta, and Alexei A Efros. Unsupervised visual representation learning by context prediction. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015.
- [20] J. Donahue, P. Krahenbühl, and T. Darrell. Adversarial feature learning. In *International Conference on Learning Representations (ICLR)*, 2016.
- [21] Alexey Dosovitskiy, Philipp Fischer, Jost Tobias Springenberg, Martin Riedmiller, and Thomas Brox. Discriminative unsupervised feature learning with exemplar convolutional neural networks. *IEEE Transactions on Pattern Analysis and Machine Intelligence (TPAMI)*, 38(9):1734–1747, 2016.
- [22] Francis Engelmann, Martin Bokeloh, Alireza Fathi, Bastian Leibe, and Matthias Nießner. 3d-mpa: Multi-proposal aggregation for 3d semantic instance segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 9031–9040, 2020.
- [23] Haoqiang Fan, Hao Su, and Leonidas J Guibas. A point set generation network for 3d object reconstruction from a single image. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 605–613, 2017.
- [24] Cat Franklin. Apple unveils new ipad pro with breakthrough lidar scanner and brings trackpad support to ipados. <https://www.apple.com/>, 2020.
- [25] Matheus Gadelha, Rui Wang, and Subhansu Maji. Multiresolution tree networks for 3d point cloud processing. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 103–118, 2018.
- [26] Andreas Geiger, Philip Lenz, Christoph Stiller, and Raquel Urtasun. Vision meets robotics: The kitti dataset. *International Journal of Robotics Research (IJRR)*, 2013.
- [27] Deepti Ghadiyaram, Matt Feiszli, Du Tran, Xueting Yan, Heng Wang, and Dhruv Mahajan. Large-scale weakly-supervised pre-training for video action recognition. *arXiv preprint arXiv:1905.00561*, 2019.
- [28] Rohit Girdhar, David Ford Fouhey, Mikel Rodriguez, and Abhinav Gupta. Learning a predictable and generative vector representation for objects. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [29] Georgia Gkioxari, Jitendra Malik, and Justin Johnson. Mesh r-cnn. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 9785–9795, 2019.
- [30] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan

- Misra. Scaling and benchmarking self-supervised visual representation learning. *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- [31] Ben Graham. Sparse 3d convolutional neural networks. *arXiv preprint arXiv:1505.02890*, 2015.
- [32] Benjamin Graham, Martin Engelcke, and Laurens van der Maaten. 3d semantic segmentation with submanifold sparse convolutional networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9224–9232, 2018.
- [33] Thibault Groueix, Matthew Fisher, Vladimir G Kim, Bryan C Russell, and Mathieu Aubry. A papier-mâché approach to learning 3d surface generation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 216–224, 2018.
- [34] Raia Hadsell, Sumit Chopra, and Yann LeCun. Dimensionality reduction by learning an invariant mapping. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2006.
- [35] Kaveh Hassani and Mike Haley. Unsupervised multi-task feature learning on point clouds. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 8160–8171, 2019.
- [36] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. *arXiv preprint arXiv:1911.05722*, 2019.
- [37] Olivier J Hénaff, Ali Razavi, Carl Doersch, SM Eslami, and Aaron van den Oord. Data-efficient image recognition with contrastive predictive coding. *arXiv preprint arXiv:1905.09272*, 2019.
- [38] Pedro Hermosilla, Tobias Ritschel, Pere-Pau Vázquez, Àlvar Vinacua, and Timo Ropinski. Monte carlo convolution for learning on non-uniformly sampled point clouds. *ACM Transactions on Graphics (TOG)*, 37(6):1–12, 2018.
- [39] Binh-Son Hua, Quang-Hieu Pham, Duc Thanh Nguyen, Minh-Khoi Tran, Lap-Fai Yu, and Sai-Kit Yeung. Scenenn: A scene meshes dataset with annotations. In *International Conference on 3D Vision (3DV)*, 2016.
- [40] Binh-Son Hua, Minh-Khoi Tran, and Sai-Kit Yeung. Point-wise convolutional neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 984–993, 2018.
- [41] Allan Jabri, Andrew Owens, and Alexei A Efros. Space-time correspondence as a contrastive random walk. *arXiv preprint arXiv:2006.14613*, 2020.
- [42] Allison Janoch, Sergey Karayev, Yangqing Jia, Jonathan T Barron, Mario Fritz, Kate Saenko, and Trevor Darrell. A category-level 3d object dataset: Putting the kinect to work. In *Consumer depth cameras for computer vision*, pages 141–165. Springer, 2013.
- [43] Xu Ji, João F Henriques, and Andrea Vedaldi. Invariant information clustering for unsupervised image classification and segmentation. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2019.
- [44] Longlong Jing, Yucheng Chen, Ling Zhang, Mingyi He, and Yingli Tian. Self-supervised modal and view invariant feature learning. *arXiv preprint arXiv:2005.14169*, 2020.
- [45] Armand Joulin, Laurens Van Der Maaten, Allan Jabri, and Nicolas Vasilache. Learning visual features from large weakly supervised data. In *European Conference on Computer Vision*, pages 67–84. Springer, 2016.
- [46] Diederik P. Kingma and Jimmy Ba. Adam: A method for stochastic optimization. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, 2015.
- [47] Roman Klokov and Victor Lempitsky. Escape from cells: Deep kd-networks for the recognition of 3d point cloud models. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 863–872, 2017.
- [48] Alexander Kolesnikov, Lucas Beyer, Xiaohua Zhai, Joan Puigcerver, Jessica Yung, Sylvain Gelly, and Neil Houlsby. Big transfer (bit): General visual representation learning. *arXiv preprint arXiv:1912.11370*, 6, 2019.
- [49] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. *arXiv preprint arXiv:1901.09005*, 2019.
- [50] Buyu Li, Wanli Ouyang, Lu Sheng, Xingyu Zeng, and Xiaogang Wang. Gs3d: An efficient 3d object detection framework for autonomous driving. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1019–1028, 2019.
- [51] Jiaxin Li, Ben M Chen, and Gim Hee Lee. So-net: Self-organizing network for point cloud analysis. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 9397–9406, 2018.
- [52] Yangyan Li, Rui Bu, Mingchao Sun, Wei Wu, Xinhan Di, and Baoquan Chen. Pointcnn: Convolution on x-transformed points. In *Advances in Neural Information Processing Systems*, pages 820–830, 2018.
- [53] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. *arXiv preprint arXiv:1608.03983*, 2016.
- [54] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. *arXiv preprint arXiv:1711.05101*, 2017.
- [55] Dhruv Mahajan, Ross Girshick, Vignesh Ramanathan, Kaiming He, Manohar Paluri, Yixuan Li, Ashwin Bharambe, and Laurens van der Maaten. Exploring the limits of weakly supervised pretraining. In *ECCV*, 2018.
- [56] J. Masci, U. Meier, D. Cires, and J. Schmidhuber. Stacked convolutional auto-encoders for hierarchical feature extraction. In *ICANN*, pages 52–59, 2011.
- [57] Daniel Maturana and Sebastian Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *2015 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 922–928. IEEE, 2015.
- [58] L. Mescheder, S. Nowozin, and A. Geiger. Adversarial variational Bayes: Unifying variational autoencoders and generative adversarial networks. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2017.
- [59] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. *arXiv preprint arXiv:1912.01991*, 2019.
- [60] Pedro Morgado, Nuno Vasconcelos, and Ishan Misra. Audio-visual instance discrimination with cross-modal agreement. <https://arxiv.org/abs/2004.12943>, 2020.
- [61] Pushmeet Kohli Nathan Silberman, Derek Hoiem and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2012.

- [62] Mehdi Noroozi and Paolo Favaro. Unsupervised learning of visual representations by solving jigsaw puzzles. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2016.
- [63] B. A. Olshausen and D. J. Field. Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381(6583):607, 1996.
- [64] Aaron van den Oord, Yazhe Li, and Oriol Vinyals. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*, 2018.
- [65] Mandela Patrick, Yuki M Asano, Ruth Fong, João F Henriques, Geoffrey Zweig, and Andrea Vedaldi. Multi-modal self-supervision from generalized data transformations. *arXiv preprint arXiv:2003.04298*, 2020.
- [66] Quang-Hieu Pham, Thanh Nguyen, Binh-Son Hua, Gemma Roig, and Sai-Kit Yeung. Jsis3d: Joint semantic-instance segmentation of 3d point clouds with multi-task pointwise networks and multi-value conditional random fields. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8827–8836, 2019.
- [67] Charles R Qi, Or Litany, Kaiming He, and Leonidas J Guibas. Deep hough voting for 3d object detection in point clouds. *arXiv preprint arXiv:1904.09664*, 2019.
- [68] Charles R Qi, Wei Liu, Chenxia Wu, Hao Su, and Leonidas J Guibas. Frustum pointnets for 3d object detection from rgb-d data. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 918–927, 2018.
- [69] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 652–660, 2017.
- [70] Charles R Qi, Hao Su, Matthias Nießner, Angela Dai, Mengyuan Yan, and Leonidas J Guibas. Volumetric and multi-view cnns for object classification on 3d data. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 5648–5656, 2016.
- [71] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Advances in neural information processing systems*, pages 5099–5108, 2017.
- [72] Marc’Aurelio Ranzato, Fu-Jie Huang, Y-Lan Boureau, and Yann LeCun. Unsupervised learning of invariant feature hierarchies with applications to object recognition. In *CVPR*, 2007.
- [73] Gernot Riegler, Ali Osman Ulusoy, and Andreas Geiger. Octnet: Learning deep 3d representations at high resolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3577–3586, 2017.
- [74] German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio Lopez. The SYNTHIA Dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [75] R. Salakhutdinov and G. Hinton. Deep Boltzmann machines. In *AI-STATS*, pages 448–455, 2009.
- [76] Samsung. What is depthvision camera on galaxy s20+ and s20 ultra? <https://www.samsung.com/>, 2020.
- [77] Jonathan Sauder and Bjarne Sievers. Self-supervised deep learning on point clouds by reconstructing space. In *Advances in Neural Information Processing Systems (NeurIPS)*, pages 12962–12972, 2019.
- [78] Manolis Savva, Angel X Chang, Pat Hanrahan, Matthew Fisher, and Matthias Nießner. Pigraphs: learning interaction snapshots from observations. *ACM Transactions on Graphics (TOG)*, 35(4):1–12, 2016.
- [79] 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 *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10529–10538, 2020.
- [80] Shaoshuai Shi, Xiaogang Wang, and Hongsheng Li. Point-rcnn: 3d object proposal generation and detection from point cloud. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 770–779, 2019.
- [81] 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. *arXiv preprint arXiv:1907.03670*, 2019.
- [82] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *European conference on computer vision*, pages 746–760. Springer, 2012.
- [83] Kihyuk Sohn. Improved deep metric learning with multi-class n-pair loss objective. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2016.
- [84] Shuran Song, Samuel P Lichtenberg, and Jianxiong Xiao. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 567–576, 2015.
- [85] Shuran Song and Jianxiong Xiao. Deep sliding shapes for amodal 3d object detection in rgb-d images. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 808–816, 2016.
- [86] Scott Stein. Lidar on the iphone 12 pro. <https://www.cnet.com/>, 2020.
- [87] Jürgen Sturm, Nikolas Engelhard, Felix Endres, Wolfram Burgard, and Daniel Cremers. A benchmark for the evaluation of rgb-d slam systems. In *2012 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 573–580. IEEE, 2012.
- [88] Hang Su, Varun Jampani, Deqing Sun, Subhransu Maji, Evangelos Kalogerakis, Ming-Hsuan Yang, and Jan Kautz. Splatnet: Sparse lattice networks for point cloud processing. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2530–2539, 2018.
- [89] Hang Su, Subhransu Maji, Evangelos Kalogerakis, and Erik Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Proceedings of the IEEE international conference on computer vision*, pages 945–953, 2015.
- [90] Chen Sun, Abhinav Shrivastava, Saurabh Singh, and Abhinav Gupta. Revisiting unreasonable effectiveness of data in deep learning era. In *ICCV*, 2017.
- [91] Pei Sun, Henrik Kretzschmar, Xerxes Dotiwalla, Aurelien Chouard, Vijaysai Patnaik, Paul Tsui, James Guo, Yin Zhou, Yuning Chai, Benjamin Caine, et al. Scalability in perception for autonomous driving: Waymo open dataset. In *Proceedings of the IEEE/CVF Conference on Computer*

- Vision and Pattern Recognition*, pages 2446–2454, 2020. This publication was made using the Waymo Open Dataset, provided by Waymo LLC under license terms available at waymo.com/open.
- [92] Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2015.
 - [93] Maxim Tatarchenko, Alexey Dosovitskiy, and Thomas Brox. Octree generating networks: Efficient convolutional architectures for high-resolution 3d outputs. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2088–2096, 2017.
 - [94] Maxim Tatarchenko, Jaesik Park, Vladlen Koltun, and Qian-Yi Zhou. Tangent convolutions for dense prediction in 3d. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3887–3896, 2018.
 - [95] Lyne Tchapmi, Christopher Choy, Iro Armeni, JunYoung Gwak, and Silvio Savarese. Segcloud: Semantic segmentation of 3d point clouds. In *2017 international conference on 3D vision (3DV)*, pages 537–547. IEEE, 2017.
 - [96] OpenPCDet Development Team. Openpcdet: An open-source toolbox for 3d object detection from point clouds. <https://github.com/open-mmlab/OpenPCDet>, 2020.
 - [97] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. *arXiv preprint arXiv:1906.05849*, 2019.
 - [98] Yonglong Tian, Chen Sun, Ben Poole, Dilip Krishnan, Cordelia Schmid, and Phillip Isola. What makes for good views for contrastive learning? *arXiv preprint arXiv:2005.10243*, 2020.
 - [99] P. Vincent, H. Larochelle, Y. Bengio, and P.-A. Manzagol. Extracting and composing robust features with denoising autoencoders. In *ICML*, 2008.
 - [100] Nanyang Wang, Yinda Zhang, Zhuwen Li, Yanwei Fu, Wei Liu, and Yu-Gang Jiang. Pixel2mesh: Generating 3d mesh models from single rgb images. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 52–67, 2018.
 - [101] Xiaolong Wang and Abhinav Gupta. Unsupervised learning of visual representations using videos. In *Proceedings of the International Conference on Computer Vision (ICCV)*, 2015.
 - [102] Xinlong Wang, Shu Liu, Xiaoyong Shen, Chunhua Shen, and Jiaya Jia. Associatively segmenting instances and semantics in point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4096–4105, 2019.
 - [103] Yue Wang and Justin M Solomon. Deep closest point: Learning representations for point cloud registration. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages 3523–3532, 2019.
 - [104] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E Sarma, Michael M Bronstein, and Justin M Solomon. Dynamic graph cnn for learning on point clouds. *arXiv preprint arXiv:1801.07829*, 2018.
 - [105] Wenxuan Wu, Zhongang Qi, and Li Fuxin. Pointconv: Deep convolutional networks on 3d point clouds. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 9621–9630, 2019.
 - [106] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *IEEE Conference on Computer Vision and Pattern Recognition, CVPR 2015, Boston, MA, USA, June 7-12, 2015*, pages 1912–1920, 2015.
 - [107] Zhirong Wu, Yuanjun Xiong, Stella X Yu, and Dahua Lin. Unsupervised feature learning via non-parametric instance discrimination. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018.
 - [108] Jianxiong Xiao, Andrew Owens, and Antonio Torralba. Sun3d: A database of big spaces reconstructed using sfm and object labels. In *Proceedings of the IEEE international conference on computer vision*, pages 1625–1632, 2013.
 - [109] Saining Xie, Jiatao Gu, Demi Guo, Charles R Qi, Leonidas J Guibas, and Or Litany. Pointcontrast: Unsupervised pre-training for 3d point cloud understanding. In *Proceedings of the European Conference on Computer Vision (ECCV)*, 2020.
 - [110] Saining Xie, Sainan Liu, Zeyu Chen, and Zhuowen Tu. Attentional shapecontextnet for point cloud recognition. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4606–4615, 2018.
 - [111] Yifan Xu, Tianqi Fan, Mingye Xu, Long Zeng, and Yu Qiao. Spidercnn: Deep learning on point sets with parameterized convolutional filters. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 87–102, 2018.
 - [112] Xu Yan, Chaoda Zheng, Zhen Li, Sheng Wang, and Shuguang Cui. Pointasnl: Robust point clouds processing using nonlocal neural networks with adaptive sampling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5589–5598, 2020.
 - [113] Yan Yan, Yuxing Mao, and Bo Li. Second: Sparsely embedded convolutional detection. *Sensors*, 18(10):3337, 2018.
 - [114] Yaoqing Yang, Chen Feng, Yiru Shen, and Dong Tian. Foldingnet: Point cloud auto-encoder via deep grid deformation. In *Proceedings of the Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 206–215, 2018.
 - [115] Zetong Yang, Yanan Sun, Shu Liu, Xiaoyong Shen, and Jiaya Jia. Std: Sparse-to-dense 3d object detector for point cloud. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1951–1960, 2019.
 - [116] Li Yi, Wang Zhao, He Wang, Minhyuk Sung, and Leonidas J Guibas. Gspn: Generative shape proposal network for 3d instance segmentation in point cloud. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 3947–3956, 2019.
 - [117] Zaiwei Zhang, Zhenxiao Liang, Lemeng Wu, Xiaowei Zhou, and Qixing Huang. Path-invariant map networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 11084–11094, 2019.
 - [118] Zaiwei Zhang, Bo Sun, Haitao Yang, and Qixing Huang. H3dnet: 3d object detection using hybrid geometric primitives. *arXiv preprint arXiv:2006.05682*, 2020.
 - [119] Zaiwei Zhang, Zhenpei Yang, Chongyang Ma, Linjie Luo, Alexander Huth, Etienne Vouga, and Qixing Huang. Deep

generative modeling for scene synthesis via hybrid representations. *ACM Transactions on Graphics (TOG)*, 39(2):1–21, 2020.

- [120] Yin Zhou and Oncel Tuzel. Voxelnet: End-to-end learning for point cloud based 3d object detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4490–4499, 2018.

Supplemental Material

We provide the model architecture details in appendix A. Hyper-parameters used in training and fine-tuning for PointNet++ and UNet and additional results are shown in appendix B. We also show hyper-parameters used in training and fine-tuning for PointnetMSG and Spconv-UNet and results for other categories of detection in KITTI in appendix C.

A. Architecture Details

PointNet++ model used in § 4 and §§ 5.1 and 5.2. As shown in Table 7, PointNet++ contains four set abstraction layers and two feature up-sampling layers, designed in [67]. Each SA layer is specified by $(n, r, [c_1, \dots, c_k])$, where n represents number of output points, r represents the ball-region radius of the reception field, c_i represents the feature channel size of the i -th layer in the MLP. Each feature up-sampling (FP) layer upsamples the point features by interpolating the features on input points to output points, as designed in [71]. Each FP layer is specified by $[c_1, \dots, c_k]$ where c_i is the output of the i -th layer in the MLP. In § 4.2.1 of the main paper, we create higher capacity versions of the PointNet++ model by increasing the channel width. For PointNet++ 2 \times , 3 \times and 4 \times , we multiply the feature size by $\{2, 3, 4\}$ of each layer in the MLP of the four set abstraction layers. When applying PointNet++ in VoteNet and H3DNet, we adjust the rest of the model accordingly based on different point feature size.

PointnetMSG model on LiDAR data used in § 5.3. Table 8 shows the architecture details for PointnetMSG, which processes lidar point cloud for PointRCNN detection model. PointnetMSG contains four multi-scale set abstraction layers and four feature up-sampling layers. Multi-scale set abstraction samples points in different scales and pro-

layer name	input layer	type	output size	layer params
sa1	point cloud (xyz)	SA	(2048,3+128)	(2048,0.2,[64, 64, 128])
sa2	sa1	SA	(1024,3+256)	(1024,0.4,[128, 128, 256])
sa3	sa2	SA	(512,3+256)	(512,0.8,[128, 128, 256])
sa4	sa3	SA	(256,3+256)	(256,1.2,[128, 128, 256])
fp1	sa3,sa4	FP	(512,3+512)	[512, 512]
fp2	sa2,sa3	FP	(1024,3+512)	[512, 512]

Table 7: PointNet++ Network Architecture used in § 4 and §§ 5.1 and 5.2

layer name	input layer	type	output size	layer params
sa1	point cloud (xyz)	SA	(4096,3+96)	(0.1, [16, 16, 32], 0.5, [32, 32, 64])
sa2	sa1	SA	(1024,3+256)	(0.5, [64, 64, 128], 1.0, [64, 96, 128])
sa3	sa2	SA	(256,3+512)	(0.1, [128, 196, 256], 0.5, [128, 196, 256])
sa4	sa3	SA	(64,3+1024)	(0.1, [256, 256, 512], 0.5, [256, 384, 512])
fp1	sa3,sa4	FP	(256,3+1024)	[512, 512]
fp2	sa2,sa3	FP	(1024,3+1024)	[512, 512]
fp3	sa1,sa2	FP	(4096,3+512)	[256, 256]
fp4	point cloud,sa1	FP	(16384,3+256)	[128, 128]

Table 8: PointnetMSG Network Architecture used in § 5.3.

cess them with different MLPs. In here, we adopt the architecture design in [80], in which each set abstraction layer contains two point features produced with two ball-region radius and two MLPs. As shown in Table 8, each SA layer is specified by $(r^1, [c_1^1, \dots, c_k^1], r^2, [c_1^2, \dots, c_k^2])$, where r^1, r^2 indicates the ball-region radius for each scale and c_i^1, c_i^2 indicates the feature channel size of the i -th layer in the MLP of each scale. We directly apply the learnt PointnetMSG in PointRCNN for detection evaluation in KITTI.

UNet model used in § 4.3. Fig 6 shows the network architecture for UNet. It mainly contains four encoding resblock and four decoding resblock. We use sparse convolution and sparse resblock designed in [17]. For each sparse resblock, we first apply sparse convolution or deconvolution, depending on encoding or decoding, with kernel size 2 and stride 2. Then, we apply N number of sparse convolution layers with kernel size 3 and stride 1. D represents the output feature dimension. Since sparse convolution takes variable sized input and output, we do not specify the number of voxels in each layer here. We directly apply the learnt UNet backbone for different scene segmentation tasks.

Spconv-UNet model on LiDAR data used in § 5.3. Fig 7 shows the network architecture for Spconv-UNet used for Part A² detection model. It mainly contains four sparse blocks for encoding and four sparse upblocks for decoding. We use sparse convolution and sparse resblock designed in [81]. For each sparse block/upblock, we show the number of convolution layers N and output feature di-

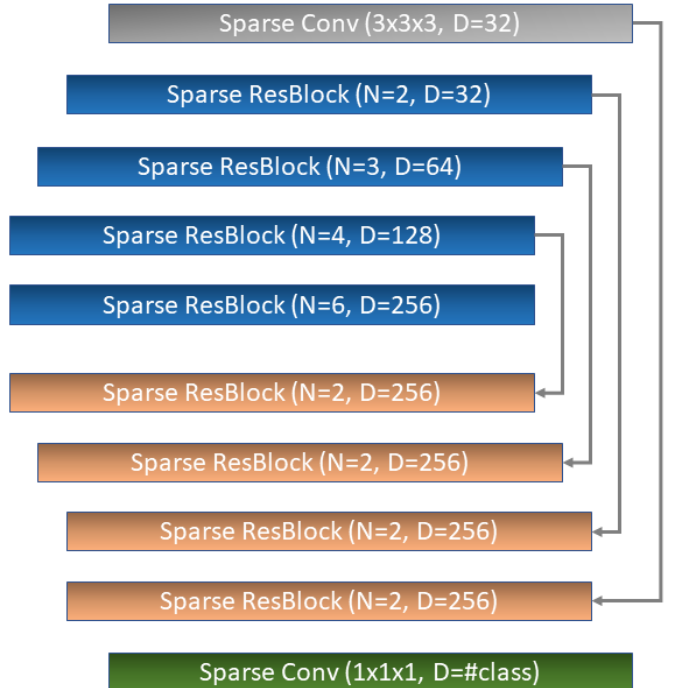


Figure 6: UNet Network Architecture used in § 4.3.

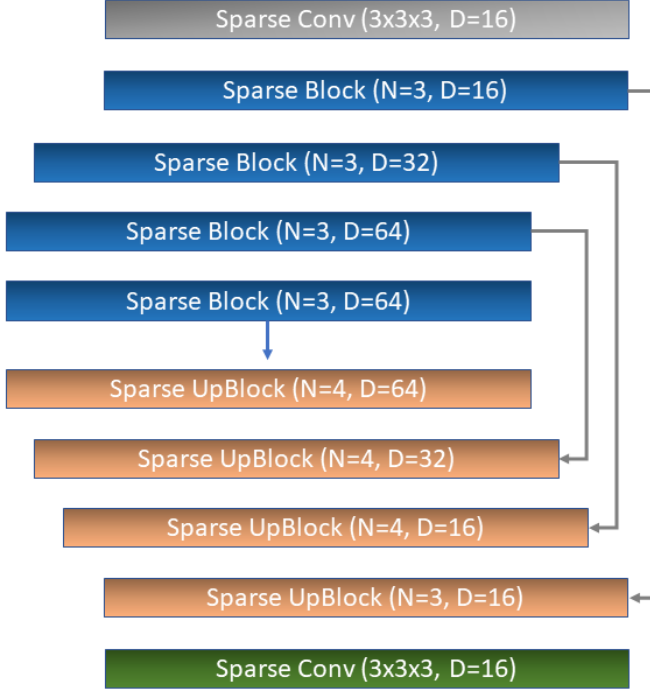


Figure 7: Sconv-UNet Network Architecture used in Section 5.3 of the main paper.

mension D . We directly adopted the architecture design in Part A^2 . For KITTI detection evaluation, we directly load the learnt backbone for fine-tuning.

B. Training Details

B.1. Pretraining Details

As mentioned in the main paper, we use a standard SGD optimizer with momentum 0.9, cosine learning rate scheduler starting from 0.12 to 0.00012 and train the model for 1000 epochs with a batch size of 1024. We observed that pretraining for 400 epochs already gives good results, and 1000 epochs for pretraining only slightly improve the model.

B.2. Experimental Details for PointNet++

For all the VoteNet fine-tuning evaluations, we use the original configurations [67], where we apply Adam optimizer [46] and use a base learning rate 0.001 with a $10\times$ weight decrease at 80, 120 and 160 epochs. The model is trained for 180 epochs in total. Since the training set of S3DIS [4] only contains 200 training instances, we train for 360 epochs. We use a batch size of 8 for ScanNet, Matterport3D, and S3DIS, and a batch size of 16 for SUNRGBD. We use the same configuration for training from scratch and fine-tuning, and we only load the pretrained PointNet++ backbone during fine-tuning.

For the H3DNet fine-tuning, we only use one backbone network instead of the original four [118]. For initial learning rate and decays, we use the original configurations. We found that with $3\times$ PointNet++ backbone, we are able to re-produce the previous results reported in the paper with one backbone network. We use the same configuration for training from scratch and fine-tuning, and we only load the pretrained PointNet++ backbone during fine-tuning.

B.2.1 ModelNet Classification

In § 5 of the main paper, we used the ModelNet dataset for transfer learning. We trained linear classifiers on fixed features for this task and measured the classification accuracy.

Full finetuning. We now show results on the same task but using finetuning, *i.e.*, all the parameters of the backbone model are updated. We use the PointNet++ backbone to extract per-point feature and apply max-pooling to get the final global feature vector. We then apply one linear layer to get the final class labels. We use SGD+momentum optimizer with an initial learning rate 0.01. We use multistep LR scheduler with $10\times$ weight decrease at 8, 16 and 24 epochs. We train for 28 epochs. We apply the data augmentation used in [67] during training. We use the same configuration for training from scratch and fine-tuning, and we only load the pretrained PointNet++ backbone during fine-tuning. For linear probing, we fix the pretrained weight and only fine-tune the last linear layer. In Table 9, we compare the fine-tuning and train from scratch results for both linear probing and full fine-tuning. The pretrained PointNet++ model provides consistent improvements.

B.2.2 Label efficiency of PointNet++

We follow the same settings as § 4.2.4 of the main paper and evaluate the label efficiency of the PointNet++ pretrained using DepthContrast. In Figure 1 (main paper) we showed that DepthContrast pretraining was label efficient on the ScanNet and SUNRGBD datasets. In Fig 8, we show the label efficiency plots for Matterport3D and S3DIS downstream detection tasks. Our results are consistently better than training from scratch across all the detection benchmarks used in the paper.

Task	Scratch	DepthContrast
ModelNet Linear (Accuracy)	50.7	85.4
ModelNet Finetune (Accuracy)	88.2	91.3

Table 9: ModelNet classification. We evaluate models by linear probing and full fine-tuning. We pretrain PointNet++ models using DepthContrast on the ScanNet-vid dataset.

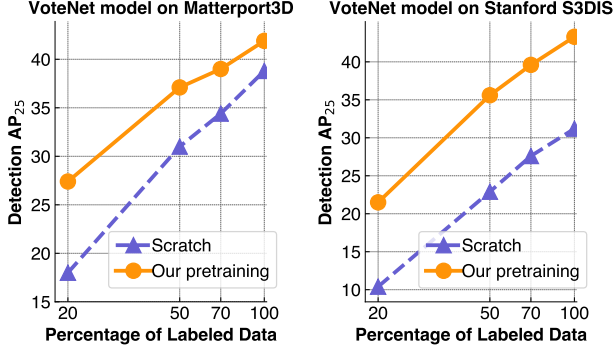


Figure 8: Label-efficiency for Matterport3D and S3DIS detection tasks. We pretrain a PointNet++ backbone model using our DepthContrast method on the ScanNet-vid dataset.

Task	Scratch	DepthContrast
ScanNet segmentation (mIOU)	70.3	71.2

Table 10: ScanNet scene segmentation. We evaluate UNet models after finetuning on the ScanNet scene segmentation task. We pretrain the DepthContrast UNet models on the ScanNet-vid dataset.

B.3. Experimental Details for UNet

For all the UNet fine-tuning evaluation, we use SGD+momentum optimizer with an initial learning rate 0.1. We use Polynomial LR scheduler with a power factor of 0.9. Weight decay is 0.0001. For voxel size, we use 0.05(5cm) for S3DIS and Synthia and 0.04(4cm) for ScanNet. We use the original data augmentation techniques in [17]. We use batch size 48 for S3DIS and 56 for ScanNet and Synthia. We train the model with 8 V100 GPUs with data parallelism for 20000 iterations for all three tasks. We use the same configuration for training from scratch and fine-tuning, and we only load the pretrained UNet backbone during fine-tuning.

B.3.1 ScanNet Segmentation

For ScanNet scene segmentation, due to memory issue, we increased the voxel size from the default 0.02(2cm) to 0.04(4cm), which leads to different scratch training results compared to [17]. Although we are pretraining and fine-tuning on the same dataset, our approach still provides improvements as shown in Table 10.

B.3.2 Label efficiency of UNet

We pretrain a UNet model using DepthContrast on the ScanNet-vid dataset. We finetune this model for scene segmentation task. In Fig 9, we show the data efficiency plot for S3DIS scene segmentation dataset. Our approach pro-

vides consistent performance boost for different percentages of training data used.

C. Experimental Details for Lidar Data

We now present the experimental details when using DepthContrast on LiDAR 3D data (Section 5.3 of the main paper).

C.1. Fine-tuning Details

We use the original configuration from PointRCNN [80] and Part A² [81] to get the scratch training results for different splits of training dataset. For PointRCNN, we use

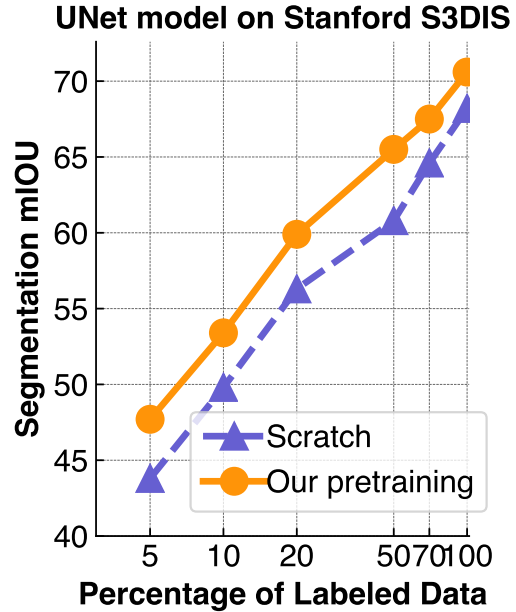


Figure 9: Label-efficiency for S3DIS scene segmentation task using the UNet model on voxel input. UNet model was pretrained on ScanNet-vid using our DepthContrast method.

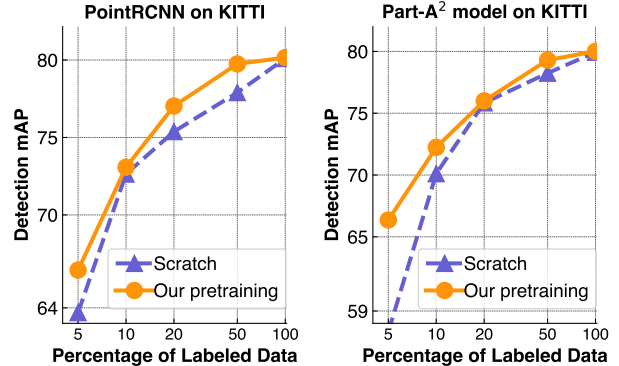


Figure 10: Label-efficiency for KITTI car detection at moderate difficulty level. We evaluate on the val split of the KITTI dataset.

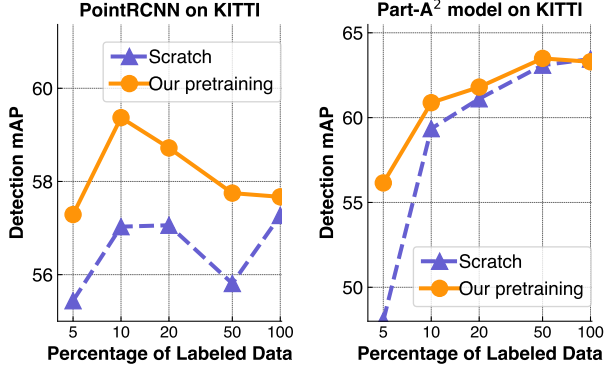


Figure 11: Label-efficiency evaluation for KITTI pedestrian detection at moderate difficulty level. We use the val split of the KITTI dataset.

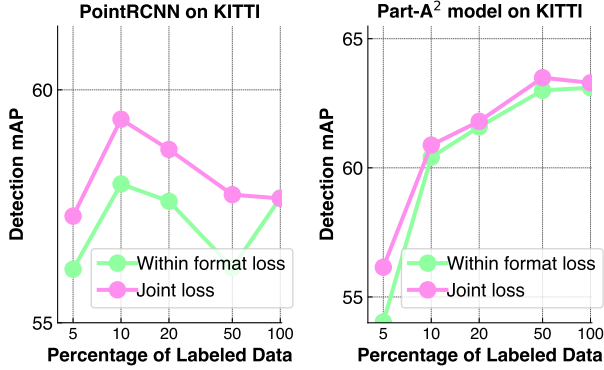


Figure 12: Comparison between within format and joint training loss. Label-efficiency evaluation for pedestrian detection at moderate difficulty level of the KITTI val split

AdamW optimizer [54] with an initial learning rate 0.01, weight decay 0.01, and momentum 0.9. For Part A^2 , we also use AdamW optimizer [54] with an initial learning rate 0.003, weight decay 0.01, and momentum 0.9. We use batch size 24 for PointRCNN and batch size 16 for Part A^2 . We train both models for 80 epochs, and the learning rate will drop by $10\times$ at 35 and 45 epochs. For both methods, we apply the same data augmentation and processing pipelines in [96].

We use the same configuration for training from scratch and fine-tuning for 5%, 10%, 20% and 50% of the labeled training data. For 100% training data, we observed overfitting issue for classes with fewer training instances, such as cyclist and pedestrian. Thus, we increased the initial learning rate by $2\times$. For PointRCNN, we also decreased the number of training epochs to 60 and set the first weight drop at 30 epochs. After modifying those parameters, we observed that the performance of scratch training didn't change.

C.2. Results on KITTI

In the main paper, due to space constraints, we only showed the results on the cyclist class in the KITTI dataset. We present results on the remainder classes.

We show the label efficiency evaluation results for car detection in KITTI in Fig 10. For simplicity, we only show the results at moderate difficulty level. The results at other difficulty levels maintain similar pattern. DepthContrast provides a performance boost with fewer training instances, especially with 5% of labeled data.

In Fig 11, we show the label efficiency evaluation results for pedestrian detection in KITTI. Our pretraining provides consistent gain over scratch for PointRCNN model. For Part A^2 , our pretraining also provides significant performance boost with fewer training instances.

Advantage of joint training over within-format loss. Similar to Section 4.3 of the main paper, we analyze if training jointly with the within and across-format losses provides better transfer performance. Specifically, we compare (Equation 3 of the main paper) our full joint loss with the within-format only loss. We pretrain separate models with these losses and evaluate them on the KITTI dataset.

In Fig 12, we show the label efficiency evaluation results for pedestrian detection. We can see that with joint loss training, the model provides more performance gain with fewer training instances, which proves that the benefit of joint training holds across different pretraining data and architectures.