# FM-OV3D: Foundation Model-based Cross-modal Knowledge Blending for Open-Vocabulary 3D Detection

**Dongmei Zhang**[1*], **Chang Li**[1*], **Ray Zhang**[2*], **Shenghao Xie**[3]
**Wei Xue**[4], **Xiaodong Xie**[1], **Shanghang Zhang**[1†]

[1]National Key Laboratory for Multimedia Information Processing, School of Computer Science, Peking University
[2]Shanghai AI Lab, [3]Wuhan University, [4]Hong Kong University of Science and Technology
dmzhang@stu.pku.edu.cn, {1900012977, donxie, shanghang}@pku.edu.cn,
xieshenghao@whu.edu.cn, xuewei.x@gmail.com

## Abstract

The superior performances of pre-trained foundation models in various visual tasks underscore their potential to enhance the 2D models' open-vocabulary ability. Existing methods explore analogous applications in the 3D space. However, most of them only center around knowledge extraction from singular foundation models, which limits the open-vocabulary ability of 3D models. We hypothesize that leveraging complementary pre-trained knowledge from various foundation models can improve knowledge transfer from 2D pre-trained visual language models to the 3D space. In this work, we propose **FM-OV3D**, a method of **F**oundation **M**odel-based Cross-modal Knowledge Blending for **O**pen-**V**ocabulary **3D D**etection, which improves the open-vocabulary localization and recognition abilities of 3D model by blending knowledge from multiple pre-trained foundation models, achieving true open-vocabulary without facing constraints from original 3D datasets. Specifically, to learn the open-vocabulary 3D localization ability, we adopt the open-vocabulary localization knowledge of the Grounded-Segment-Anything model. For open-vocabulary 3D recognition ability, We leverage the knowledge of generative foundation models, including GPT-3 and Stable Diffusion models, and cross-modal discriminative models like CLIP. The experimental results on two popular benchmarks for open-vocabulary 3D object detection show that our model efficiently learns knowledge from multiple foundation models to enhance the open-vocabulary ability of the 3D model and successfully achieves state-of-the-art performance in open-vocabulary 3D object detection tasks. Code is released at https://github.com/dmzhang0425/FM-OV3D.git.

## Introduction

Open-vocabulary ability refers to the ability of models to generate or understand concepts that have not been explicitly included in training datasets. Pre-trained foundation models' high performances in various 2D open-vocabulary visual tasks demonstrate their strong open-vocabulary abilities (Liang et al. 2023; Liu et al. 2023a). However, utilizing vision-text pairs in training to enable such ability in 3D

models is challenging, since it is difficult to collect a sizable dataset of 3D point clouds paired with texts.

The knowledge embedded in pre-trained foundation models can potentially enhance 3D models. Despite the differences in modalities between 3D point-cloud and 2D images, they both share visual information about objects. There have been efforts to investigate how to transfer knowledge from 2D to 3D models, leading to various distinct methods (Zhang et al. 2022; Zhu et al. 2023; Zhang et al. 2023c).

However, many of them primarily extract knowledge from individual models. Considering the disparities in training objectives, model architectures, and training data among various models, pre-trained models' knowledge, abilities, or perception of the world may exhibit diversity. This diversity can potentially enhance the open-vocabulary ability of 3D models complementarily. For example, different from the contrastive vision-language knowledge in CLIP (Radford et al. 2021), SAM (Kirillov et al. 2023) is designed to segment all objects in an image, providing information about their positions and sizes. Moreover, when vocabulary is involved in training open-vocabulary models to correlate visual and textual features, existing methods only use predefined class lists or captions, failing to provide rich information about the classes themselves, thereby limiting recognition performance. As a text-generative model, GPT-3 (Brown et al. 2020) has a rich understanding of various classes and can serve as a source of textual knowledge. Therefore, we hypothesize that harnessing complementary pre-trained knowledge from different models can facilitate knowledge transfer from 2D pre-trained models to 3D space.

In this work, we propose FM-OV3D, a method of **F**oundation **M**odel-based Cross-modal Knowledge Blending for **O**pen-**V**ocabulary **3D D**etection, which improves the open-vocabulary localization and recognition ability of 3D models by incorporating knowledge from diverse pre-trained foundation models, without requiring any manual annotations. Specifically, to train the open-vocabulary localization ability of 3D models, we utilize the object localization knowledge within the Grounded-Segment-Anything model to generate 2D bounding boxes. To enhance the open-vocabulary recognition ability of 3D models, we associate the semantics among three different modalities of the same class: point cloud features from the 3D detec-

---

tor, CLIP extracted textual features of GPT-3 generated language prompts, and CLIP extracted visual features of Stable Diffusion-generated 2D visual prompts and point clouds' paired images. We perform open-vocabulary object detection in testing by comparing point cloud features and text features in a common feature space. Moreover, our method can be applied to any manually selected open-vocabulary training set since our GPT-3 language prompts and Stable Diffusion visual prompts can be generated regarding any selected class. The major contributions of our work include:

- We propose that leveraging complementary pre-trained knowledge from various foundation models can facilitate knowledge transfer from 2D pre-trained visual language models to the 3D space.

- We propose FM-OV3D, a method of foundation model-based cross-modal knowledge blending for open-vocabulary 3D detection, which incorporates knowledge of various foundation models to enhance the open-vocabulary localization and recognition ability of 3D models without requiring any manual annotation, which can be easily transferred to any 3D dataset.

- Experiments conducted on two public and commonly used open-vocabulary 3D point-cloud object detection datasets achieve *state-of-the-art* performances, demonstrating that our method is effective.

## Related Work

### Pre-trained Foundation Models

Pre-trained foundation models are trained on massive amounts of data on a pre-defined proxy task. Models learn statistical structures and grasp the intrinsic links within training data, acquiring extensive knowledge. Large language models (LLMs) like GPT-3 (Brown et al. 2020) are trained on a vast collection of internet text in self-supervised learning. They can generate human-like language responses and have been applied to various natural language processing downstream tasks (Wang et al. 2023; Ni and Li 2023). SAM (Kirillov et al. 2023) and Per-SAM (Zhang et al. 2023b) successfully incorporate visual content-relevant knowledge and have demonstrated high performances on various tasks (Liu et al. 2023b; Hu and Li 2023). However, these models' scope of knowledge is limited by insufficient training data across various modalities, training methods, and proxy task types used in training. As a result, existing pre-trained large models' applicability to downstream tasks is limited.

Recent research explores combining these pre-trained foundation models in various modalities. For example, CaFo (Zhang et al. 2023a) cascades a variety of pre-trained foundation models to achieve better image classification performance. Grounding DINO (Liu et al. 2023a), which blends the knowledge of DINO with textual prompts, has state-of-the-art results in zero-shot settings. Grounded-SAM (Kirillov et al. 2023; Liu et al. 2023a), which combines Grounding DINO (Liu et al. 2023a) with SAM (Kirillov et al. 2023), enhances detection and segmentation abilities simultaneously. However, constrained by fused models' limited

modalities, there are still multi-modal problems, for instance, object detection and segmentation in 3D scenarios, yet to be explored.

### Open-Vocabulary 2D/3D Object Detection

Open-vocabulary detection requires models to localize and recognize novel classes with training on only base classes (De Rijk et al. 2022; Bangalath et al. 2022; Rahman, Khan, and Barnes 2020; Rahman, Khan, and Porikli 2020; Zareian et al. 2021). Typically, knowledge of novel classes is indirectly implicated by cues from other modalities, for example, textual cues. To enhance open-vocabulary detection capabilities, some studies explore rich image-text pairs' semantics' extraction (Zareian et al. 2021). Some works (Rahman, Khan, and Barnes 2020; Rahman, Khan, and Porikli 2020; Zareian et al. 2021) replace visual detectors' classification layer with a visual-textual embedding to achieve robust performance in open-vocabulary settings.

In 3D point cloud detection tasks, directly applying visual-language pre-trained models faces the challenge of acquiring large-scale point cloud-text pairs and the gap between image and point cloud modalities. Existing works seek solutions in utilizing foundation models' knowledge on 3D tasks. PointCLIP series (Zhang et al. 2022; Zhu et al. 2023; Guo* et al. 2022) use CLIP to process multi-view images projected from 3D modality, and Point-Bind&Point-LLM (Guo et al. 2023) leverage LLMs and multi-modality semantics to achieve zero-shot 3D analysis. However, CLIP and LLMs are not trained for localization, and the localization in 2D space and 3D space differs significantly, which limits the localization capability of this method. Lu *et al.* (Lu et al. 2022) expands the 3D detector's vocabulary with ImageNet1K (Russakovsky et al. 2015). Concurrently, Lu *et al.* (Lu et al. 2023) proposes a divide-and-conquer strategy to connect textual information with point clouds. They might limit the 3D detector's generalization ability originating from the dataset applied. In this paper, we leverage multiple pre-trained models' knowledge from textual and image modalities, requiring no human-annotated data, enhancing our model's detection performance in open-vocabulary settings.

## Methodology

As shown in Figure 1, during training, we take raw point clouds, corresponding 2D images, and training vocabularies as input. During testing, the requirement for 2D images is eliminated, and the model relies only on testing vocabulary and the raw 3D point clouds.

Our 3D detector is required to predict the 3D bounding boxes transformed from the 2D results of Grounded-SAM to improve the open-vocabulary localization ability. Regarding 3D open-vocabulary recognition ability, we blend the knowledge of GPT-3, Stable Diffusion, and CLIP. We conduct our recognition loss utilizing 3D features extracted by 3D detector, 2D features, and textual features extracted by CLIP, leveraging CLIP's rich cross-modal knowledge. Details are explained in the following sections.
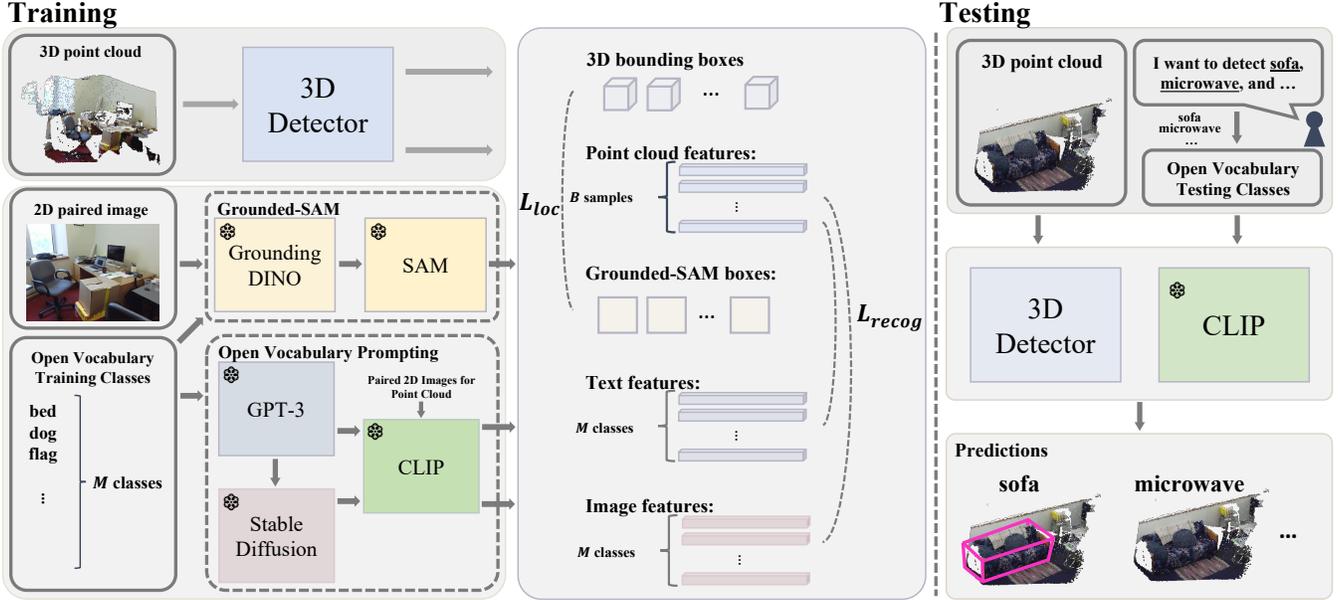
Figure 1: **The pipeline of FM-OV3D.** Given raw point clouds, corresponding 2D images, and training vocabulary, we train our model leveraging extensive knowledge from pre-trained foundation models, without requiring any annotations. $M$ and $B$ represent the size of the training vocabulary and the number of point cloud sample in a training batch, respectively. $L_{loc}$ aims at improving our 3D detector's localization ability, while $L_{recog}$ is designed to improve its recognition ability. Point cloud visualization and its paired 2D image are selected from the SUN RGB-D (Song, Lichtenberg, and Xiao 2015) dataset, while training vocabulary is from the LVIS (Gupta, Dollar, and Girshick 2019) dataset. '*' means the model is frozen.

## Open-Vocabulary 3D Localization

For the 3D detector's open-vocabulary localization ability, we employ the Grounded-Segment-Anything model, denoted as Grounded-SAM, to generate 3D bounding boxes for each point-cloud data and require our 3D detector to predict them.

After being prompted with a training vocabulary, Grounded-SAM generates 2D detection boxes on 2D images corresponding to a set of images for a given point cloud. The training vocabulary consists of the LVIS (Gupta, Dollar, and Girshick 2019) dataset that encompasses classes represented in text. From this dataset, a subset of $M$ classes $\{a_1, \ldots, a_M\}$ is selected for prompting. It's noteworthy that we are agnostic about the classes of the 2D detection boxes. The process of generating 2D bounding boxes for image $I$ is given by:

$$Box2D_I = GroundedSAM(a_1, \ldots, a_M, I) \quad (1)$$

where $Box2D_I \in \mathbf{R}^4$ represents the 2D boxes generated by Grounded-SAM. We then project these 2D bounding boxes via projection $K$ to 3D space and perform clustering to tighten the 3D bounding boxes.

$$Box3D_I = Cluster(Box2D_I \circ K^{-1}) \quad (2)$$

where $K$ is the projection matrix, which is provided in the datasets, and $Box3D_I \in \mathbf{R}^7$ represents the transformed 3D bounding boxes. $Clustering$ is a density-based clustering approach performed on points inside the projected bounding box to eliminate irrelevant outliers.

We supervise our 3D detector's predicted 3D bounding boxes $Box3D_{pc} \in \mathbf{R}^7$ by above matched 3D bounding boxes $Box3D_I$. We compute bounding box regression loss following 3DETR (Misra, Girdhar, and Joulin 2021) demonstrated by Equation 3.

$$L_{loc} = \sum_{I=1}^{P} L_{box}^{3D}(Box3D_I, Box3D_{pc}) \quad (3)$$

where $L_{box}^{3D}$ represents regression loss between $Box3D_{pc}$ and the target $Box3D_I$. $P$ represents the number of matched bounding boxes in training. By minimizing the value of $L_{loc}$, we enhance our 3D detector's open-vocabulary localization ability without requiring any annotation.

## Open-Vocabulary 3D Recognition

We improve our 3D detector's open-vocabulary recognition ability by blending knowledge of single-modal generative foundation models, including GPT-3 and Stable Diffusion models, and cross-modal discriminative models like CLIP. Specifically, we utilize GPT-3 to generate rich textual prompts and the Stable Diffusion model to generate rich 2D visual prompts, then use CLIP to extract their features.

Then we blend their knowledge by aligning object class semantics across three modalities: point cloud, images, and texts.

**Text Prompt Generation** GPT-3, with 175 billion parameters, is trained on a substantial amount of internet text in a self-supervised manner. We utilize its ability to generate

rich, detailed, human-like language descriptions on training vocabulary.

For every training class, we prompt GPT-3 to generate detailed descriptions. We adopt existing templates from (Pratt et al. 2023) and prompt GPT-3 with ten rounds each, including *"Describe what {class} look like"*, *"How can you identify {class} ?"*, *"What does {class} look like?"*, *"Describe an image from the internet of {class} "* and *"A caption of an image of {class}:".* We denote generated text prompts of $M$ training classes as $\{T_1, \ldots, T_M\}$, and the overall text prompts as $T$. The prompts generated by GPT-3 contain extensive interpretations of semantic concepts, thus providing high-quality, diverse knowledge in textual modality.

**2D Visual Prompt Generation**   We generate rich 2D images to provide our 3D model with abundant visual representations of open-vocabulary classes, broadening the vocabulary in the original 3D dataset. Stable Diffusion has an extensive textual-visual understanding and can generate synthesized 2D images according to language prompts. Therefore, utilizing GPT-generated detailed descriptions $T_1, \ldots, T_M$ for training vocabulary, we generate corresponding 2D images $SI$.

$$SI_i = SD(T_i), \quad SI = SD(T) \qquad (4)$$

where $i$ ranges from 1 to $M$, $SI$ denotes all the 2D visual prompts generated for $M$ classes in training vocabulary. Given any training vocabulary, we can expand the training data in 2D vision and language modalities utilizing our method without requiring any human annotations, tackling the problem of limited represented classes in annotated 3D datasets.

**Knowledge Blending**   The knowledge of GPT-3, Stable Diffusion, and CLIP are blended by aligning object class semantics across three modalities: point cloud features from our 3D detector, CLIP-extracted 2D image features of Stable-Diffusion generated 2D visual prompts and CLIP-extracted text features of textual prompts.

After blending, our 3D detector is trained to grasp the intrinsic links between visual objects in 3D modality and semantic concepts in 1D text modality. We first project our 3D detector's predicted bounding boxes $Box3D_{pc}$ onto paired 2D image via dataset-provided projection matrix $K$, then get image crops $I_{pc}$. The point cloud ROI features within $Box3D_{pc}$ are also extracted, denoted as $F_{pc}$. Exploiting CLIP's visual-textual knowledge, we use CLIP to extract GPT-generated prompts $T$'s textual feature $F_t$, Stable Diffusion-generated 2D image $SI$'s features $F_{2D_{SI}}$, and image crops $I_{pc}$'s features $F_{2D_{pc}}$. The combination of $F_{2D_{SI}}$ and $F_{2D_{pc}}$ is represented as $F_{2D}$. The recognition loss among point cloud features, text features, and image features is given by:

$$L_{recog} = L_{cl}(F_{pc}, F_{2D}) + L_{cl}(F_{pc}, F_t) \qquad (5)$$

Specifically, given a batch of features with size $B$ of 3D point cloud samples, $L_{cl}$ following (Oord, Li, and Vinyals 2018) can be computed as follows:

$$L_{cl}(F_1, F_2) = -\frac{1}{B} \sum_{b=1}^{B} \log(\frac{f(b, positive)}{f(b)}) \qquad (6)$$

$F_1$ can be $F_{pc}$ and $F_2$ can be $F_{2D}$ or $F_t$ in Equation 5. $f(b, positive)$ and $f(b)$ for every sample $f_b$ in the batch can be computed as follows:

$$f(b, positive) = \sum_{j=1}^{n} exp(f_b' f_j / \tau),$$
$$f(b) = \sum_{k=1}^{B} exp(f_b' f_k / \tau) \qquad (7)$$

where $\tau$ is the temperature parameter, $n$ is the number of positive samples.

The total loss function of our 3D detector can be computed as Equation 8:

$$L = L_{loc} + L_{recog} \qquad (8)$$

# Experiments

In this section, we evaluate our FM-OV3D on widely used 3D detection datasets and analyze the incorporated foundation model's effects on open-vocabulary 3D detection models. We also discuss the influence of some key parameters.

## Datasets and Evaluation Metrics

**Datasets**   We conduct experiments on public, widely used datasets **SUN RGB-D** (Song, Lichtenberg, and Xiao 2015) and **ScanNet** (Dai et al. 2017) in 3D object detection tasks. The provided point-cloud data and corresponding images, together with the matrix $K$, are used in our method.

**Evaluation Metrics**   We use mean Average Precision (AP), and Average Recall (AR) at IoU thresholds of 0.25 and 0.5, denoted as $mAP_{25}$, $mAP_{50}$, $AR_{25}$, and $AR_{50}$, as our primary metrics.

## Implementation Details

We adopt LVIS (Gupta, Dollar, and Girshick 2019) as our training vocabulary. 600 random classes sampled from our training vocabulary are used to prompt Grounded-SAM to generate 2D bounding boxes. We adopt five templates as the commands for generating GPT-3 textual prompts and later compute the mean textual feature of each class, following (Pratt et al. 2023). We apply the stable-diffusion-v1-4 model commanded by GPT-3 generated prompts. CLIP version ViT-B/32 is used for extracting features. We conduct our ablation studies on the SUN RGB-D dataset.

We train our model in 400 epochs. The base learning rate is set to 7e-4. We load 8 3D scenes onto each GPU in every batch. We adopt 3DETR (Misra, Girdhar, and Joulin 2021) as the 3D detector. Experiments are conducted on two NVIDIA GeForce RTX 3090 GPUs and A100 SXM4 80GB GPUs. In evaluation, we take our 3D detector's predicted 3D boxes as the localization result and the CLIP-predicted label of its corresponding 2D image crop as its label output.

## Performance on Open-Vocabulary 3D Object Detection

Since no prior studies have addressed open-vocabulary 3D point cloud detection problems by avoiding the need for human annotations, we compare our model's performance with

Table 1: Detection results ($AP_{25}$) on **SUN RGB-D** dataset. We report the accuracy of different classes and their mean score. * denotes our annotation-free version.

| Method | toilet | bed | chair | bathtub | sofa | dresser | scanner | fridge | lamp | desk | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GroupFree3D (Liu et al. 2021) | 0.23 | 0.04 | 1.25 | 0.03 | 0.21 | 0.21 | 0.14 | 0.10 | 0.03 | 3.02 | 0.53 |
| VoteNet (Qi et al. 2019) | 0.12 | 0.05 | 1.12 | 0.03 | 0.09 | 0.15 | 0.06 | 0.11 | 0.04 | 2.10 | 0.39 |
| H3DNet (Zhang et al. 2020) | 0.24 | 0.10 | 1.28 | 0.05 | 0.22 | 0.22 | 0.13 | 0.14 | 0.03 | 6.09 | 0.85 |
| 3DETR (Misra, Girdhar, and Joulin 2021) | 1.57 | 0.23 | 0.77 | 0.24 | 0.04 | 0.61 | 0.32 | 0.36 | 0.01 | 8.92 | 1.31 |
| OS-PointCLIP (Zhang et al. 2022) | 7.90 | 2.84 | 3.28 | 0.14 | 1.18 | 0.39 | 0.14 | 0.98 | 0.31 | 5.46 | 2.26 |
| OS-Image2Point (Xu et al. 2021) | 2.14 | 0.09 | 3.25 | 0.01 | 0.15 | 0.55 | 0.04 | 0.27 | 0.02 | 5.48 | 1.20 |
| Detic-ModelNet (Zhou et al. 2022) | 3.56 | 1.25 | 2.98 | 0.02 | 1.02 | 0.42 | 0.03 | 0.63 | 0.12 | 5.13 | 1.52 |
| Detic-ImageNet (Zhou et al. 2022) | 0.01 | 0.02 | 0.19 | 0.00 | 0.00 | 1.19 | 0.23 | 0.19 | 0.00 | 7.23 | 0.91 |
| OV-3DETIC (Lu et al. 2022) | 43.97 | 6.17 | 0.89 | **45.75** | 2.26 | **8.22** | 0.02 | 8.32 | 0.07 | **14.60** | 13.03 |
| **FM-OV3D*** | 32.40 | 18.81 | **27.82** | 15.14 | **35.40** | 7.53 | **1.95** | **9.67** | 13.57 | 7.47 | 16.98 |
| **FM-OV3D** | **55.00** | **38.80** | 19.20 | 41.91 | 23.82 | 3.52 | 0.36 | 5.95 | **17.40** | 8.77 | **21.47** |

Table 2: Detection results ($AP_{25}$) on **ScanNet** dataset. We report the accuracy of different classes and their mean score. * denotes our annotation-free version.

| Method | toilet | bed | chair | sofa | dresser | table | cabinet | bookshelf | pillow | sink | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|
| GroupFree3D (Liu et al. 2021) | 0.63 | 0.52 | 1.25 | 0.52 | 0.20 | 0.59 | 0.52 | 0.25 | 0.01 | 0.15 | 0.49 |
| VoteNet (Qi et al. 2019) | 0.04 | 0.02 | 0.12 | 0.00 | 0.02 | 0.11 | 0.07 | 0.05 | 0.00 | 0.00 | 0.04 |
| H3DNet (Zhang et al. 2020) | 0.55 | 0.29 | 1.70 | 0.31 | 0.18 | 0.76 | 0.49 | 0.40 | 0.01 | 0.10 | 0.48 |
| 3DETR (Misra, Girdhar, and Joulin 2021) | 2.60 | 0.81 | 0.90 | 1.27 | 0.36 | 1.37 | 0.99 | 2.25 | 0.00 | 0.59 | 1.11 |
| OS-PointCLIP (Zhang et al. 2022) | 6.55 | 2.29 | 6.31 | 3.88 | 0.66 | 7.17 | 0.68 | 2.05 | 0.55 | 0.79 | 3.09 |
| OS-Image2Point (Xu et al. 2021) | 0.24 | 0.77 | 0.96 | 1.39 | 0.24 | 2.82 | 0.95 | 0.91 | 0.00 | 0.08 | 0.84 |
| Detic-ModelNet (Zhou et al. 2022) | 4.25 | 0.98 | 4.56 | 1.20 | 0.21 | 3.21 | 0.56 | 1.25 | 0.00 | 0.65 | 1.69 |
| Detic-ImageNet (Zhou et al. 2022) | 0.04 | 0.01 | 0.16 | 0.01 | 0.52 | 1.79 | 0.54 | 0.28 | 0.04 | 0.70 | 0.41 |
| OV-3DETIC (Lu et al. 2022) | 48.99 | 2.63 | 7.27 | 18.64 | **2.77** | **14.34** | **2.35** | 4.54 | 3.93 | 21.08 | 12.65 |
| **FM-OV3D*** | 2.17 | 41.11 | **27.91** | 33.25 | 0.67 | 12.60 | 2.28 | **8.47** | 9.08 | 5.83 | 14.34 |
| **FM-OV3D** | **62.32** | **41.97** | 22.24 | 31.80 | 1.89 | 10.73 | 1.38 | 0.11 | **12.26** | **30.62** | **21.53** |

a 3D detection model (Lu et al. 2022) that utilizes human annotations and is exposed to open-set knowledge from other modalities. We select our open-testing classes following (Lu et al. 2022) and adopt their models discussed in an open set setting for comparison. We denote our model trained in an annotation-free setting as FM-OV3D*, and FM-OV3D represents the model trained only utilizing knowledge blending, utilizing Detic (Zhou et al. 2022) for 2D bounding box predictions. Results of our experiments on SUN RGB-D and ScanNet are shown in Table 1 and Table 2.

Our annotation-free model surpasses existing open-set 3D point cloud detector benchmarks, reaching $16.98\%$ on SUN RGB-D and $14.34\%$ on ScanNet in the $mAP_{25}$, demonstrating our model's outstanding performance on detecting 3D objects outside the training vocabulary, indicating its strong open-vocabulary ability. Furthermore, compared to OV-3DETIC (Lu et al. 2022), which achieves strong performance by leveraging the knowledge in 2D image datasets, our model blends the knowledge from both textual and 2D visual modalities. Utilizing pre-trained models' generative knowledge allows our 3D detector to grasp the intrinsic links among three modalities, without exploiting knowl-

edge from other datasets. Therefore, our method has no constraints originating from our leveraged cross-modal knowledge. Also, our method does not require human annotation in training. Our outstanding experiment results on open-vocabulary testing classes indicate that by incorporating general representations learned by various foundation models, we can bridge the gap between the limited classes in annotated 3D datasets and real-world applications, improving the 3D detector's open-vocabulary ability. Our method can be applied to any selected open-vocabulary training set without utilizing data other than raw 3D point cloud data.

Our model FM-OV3D that only leverages our knowledge blending stage demonstrates significantly enhanced performance, outperforming previous methods by $8.44\%$ on SUN RGB-D and $8.88\%$ on ScanNet in terms of $mAP_{25}$. Despite the notable improvement in overall $mAP_{25}$ performance, we observe certain classes where FM-OV3D exhibits less satisfactory results. For example, the performance on the 'scanner' of FM-OV3D falls short compared to FM-OV3D*. This suggests the potential for further enhancement of FM-OV3D by enriching its visual information related to open-vocabulary concepts.

Table 3: **Ablation study (%) of pre-trained foundation models.** 'GS' represents Grounded-SAM, and 'SD' is short for Stable Diffusion. 'Annotation' indicates whether 2D and 3D annotations are used.

| Models | Annotation | $mAP_{25}$ | $AR_{25}$ | $mAP_{50}$ | $AR_{50}$ |
|---|---|---|---|---|---|
| GPT-3 | ✓ | 18.09 | **53.87** | 1.88 | **11.58** |
| SD | ✓ | 16.34 | 47.69 | 1.10 | 8.60 |
| GPT-3 + SD | ✓ | **18.19** | 49.90 | **1.93** | 10.05 |
| GS | - | 16.47 | 55.59 | 1.84 | 11.47 |
| GS + GPT-3 | - | 15.78 | 54.23 | **1.99** | **12.99** |
| GS + GPT-3 + SD | - | **16.98** | **57.22** | 1.86 | 12.16 |

Table 4: **Ablation study (%) of the number of selected text prompts of each class.** $TP$ represents the number of GPT-3 generated text prompts selected in an individual class.

| TP | $mAP_{25}$ | $mAP_{50}$ | $AR_{25}$ | $AR_{50}$ |
|---|---|---|---|---|
| 12 | 11.66 | 1.10 | 45.37 | 7.67 |
| 25 | 11.70 | **1.38** | **48.41** | 8.26 |
| 51 | **12.50** | 1.16 | 46.73 | **8.65** |

Table 5: **Ablation study (%) of number of classes of 2D visual prompts.** $VP$ represents the number of selected Stable-Diffusion generated 2D visual prompts.

| VP | $mAP_{25}$ | $mAP_{50}$ | $AR_{25}$ | $AR_{50}$ |
|---|---|---|---|---|
| 0 | 14.09 | 2.00 | 50.21 | 12.56 |
| 3 | 12.84 | **2.24** | 51.18 | **13.11** |
| 5 | **15.49** | 2.13 | **53.48** | 11.92 |
| 7 | 12.93 | 2.20 | 47.98 | 10.20 |
| 10 | 13.76 | 1.73 | 52.14 | 12.05 |

## Ablation Study

**The effect of pre-trained foundation models.** We investigate the effects of each pre-trained foundation model in our 3D detector's training, as shown in Table 3. 'GS' denotes using the 2D bounding boxes generated by the Grounded-SAM model, 'GPT-3' represents that we utilize GPT-3 to generate text prompts according to the open-vocabulary training classes, and 'SD' denotes the 2D visual prompts generation utilizing the Stable Diffusion model. We apply bounding boxes generated by Detic (Zhou et al. 2022) as our 2D detection baseline, which are marked as annotation-needed in Table 3. The experimental results of our method are demonstrated in the last row.

Compared to the annotation-need model which is trained using only 2D visual prompts generated by Stable Diffusion or text prompts by GPT-3, the same model which is trained using both 2D visual prompts and text prompts shows better performance in $mAP_{25}$ and $mAP_{50}$. Similarly, when training the annotation-free model with both text and 2D visual prompts, its performance surpasses that of using either prompt alone or having no prompts at all. Specifically, the model incorporating both prompts achieves the best performance on the more stringent $mAP_{25}$ and $AR_{25}$ metrics. Meanwhile, the model exclusively utilizing text prompts attains the best $mAP$ and $AR$ when the threshold is set to 50. These results demonstrate the effectiveness of diverse semantic information from text-generative and cross-modal generative models even in the presence of annotations. It is reasonable since the method that only enlarges textual information lacks bridging between visual representations

and textual knowledge. Although combining rich 2D visual prompts can enrich the original 3D dataset's visual data, the model lacks direct visual-textual cross-modal understanding, and GPT-3 includes rich semantics information on target concepts and is in a more direct form of the model's object recognition predictions.

We investigate the effect of bounding boxes generated by open-vocabulary classes-prompted Grounded-SAM. In comparison to using Detic (Zhou et al. 2022) for generating 2D bounding boxes during training our detector's localization abilities, our method demonstrates outstanding performance in $AR_{25}$ and $AR_{50}$ without relying on any manual annotations, indicating its strong ability to localize objects in 3D modalities accurately.

Overall, compared with models that ditch one component in our design, the results of our method demonstrate that with the open-vocabulary 2D grounding abilities, language generative knowledge, and 2D visual generative knowledge combined, the 3D detector can fully understand the concepts' representation in three modalities, achieving superior results in open vocabulary 3D detection tasks.

**Number of selected text prompts of each class.** We generate multiple text prompts using GPT-3 in our method. We study the influence of the number of generated text prompts for an individual class on the performance of our model. The results are shown in Table 4. Considering training difficulty and time constraints, we evaluate the performance of 12, 25, and all text prompts. Experimental results indicate that uti-

Figure 2: Bounding boxes generated by the Grounded-Segment-Anything model and Detic Original pictures are selected from the **SUN RGB-D** dataset.
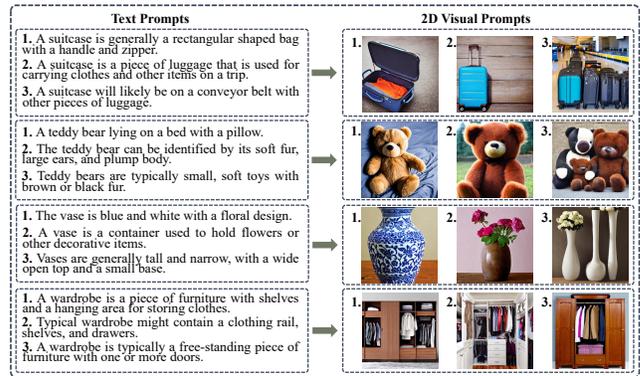


Figure 3: Visualization of GPT-generated prompts and their paired 2D visual prompts generated by Stable Diffusion on open vocabulary training dataset sampled from **LVIS**.

lizing all text prompts yields the highest $mAP_{25}$, demonstrating the effectiveness of rich textual features in enhancing the recognition capability of the detector. However, even when using only half of the text prompts (*i.e.*, 25), the model still achieves the best $mAP_{50}$ and $AR_{25}$. Although reducing the number of text prompts has a certain impact on the model's recognition performance, it remains a competitive option as it achieves comparable performance with less generation and training burden. Further reduction in the number of prompts significantly affects the model's performance, underscoring the importance of incorporating text prompts.

**Number of classes of used 2D visual prompts.** The performances of applying different numbers of classes of 2D visual prompts are shown in Table 5. We report the best $mAP_{25}$ and $AR_{25}$ when applying 5 classes of visual prompts. Compared to using a smaller number of classes or ditching 2D visual prompts, the performance achieved with 5 classes demonstrates the effectiveness of 2D visual prompts in enhancing the model's recognition capability. The model's performance declines as the number of classes increases. This can be attributed to the introduction of excessive negative samples, which introduce noise and conflicts, hindering the model's ability to learn accurate and robust feature representations and increasing the difficulty of training. Therefore, incorporating 2D visual prompts is crucial for enhancing model performance, but the number of selected classes applied should be taken into consideration to achieve optimal results.

### Qualitative Analysis

We compare the bounding boxes generated by Grounded-SAM to the 2D boxes generated by Detic (Zhou et al. 2022). We assign the class prediction of boxes generated by Grounded-SAM to a random value since we do not further use these labels for predictions. In Figure 2, objects localized by the Grounded-SAM are based on our prompted open-vocabulary classes, not facing constraints originating from the classes detected by pre-trained 2D detectors. More objects are detected when Grounded-SAM is applied. Since we use replaceable, open-vocabulary training classes to prompt the model, we achieve strong open-vocabulary ability.

We visualize GPT-3-generated text prompts and the paired 2D visual prompts generated by stable diffusion according to the GPT-3 text prompts. In Figure 3, we show that our visual prompts are diverse and variant in their direct representation of the commanded classes due to the rich GPT-3 generated descriptions, enriching the visual information and alleviating original 3D datasets' data insufficiency problem. Our visual prompts also successfully grasp the semantic concepts of the commanded classes, utilizing Stable Diffusion's vision-language knowledge.

## Conclusion

We demonstrate that leveraging complementary pre-trained knowledge from various foundation models can improve the knowledge transfer from 2D pre-trained foundation models to the 3D space, therefore enhancing the open-vocabulary ability of 3D models. We propose FM-OV3D, a foundation model-based cross-modal knowledge blending for open-vocabulary 3D object detection method that correlates multi-modal knowledge from different foundation models onto the 3D modality without requiring any human annotation. We train our 3D detector in aspects of localization and recognition. For open-vocabulary localization, we integrate the Grounded-Segment-Anything model's 2D knowledge by transforming its 2D bounding box predictions to supervise our model's localization results. For open-vocabulary recognition, we blend the knowledge from pre-trained text and image generative models and cross-modal discriminative models with knowledge in 3D modality, bridging the gap between abundant real-world classes and the insufficiency of classes in 3D datasets. We conduct experiments on SUN RGB-D and ScanNet datasets, and our experimental results demonstrate that our method is effective.

## References

Bangalath, H.; Maaz, M.; Khattak, M. U.; Khan, S. H.; and Shahbaz Khan, F. 2022. Bridging the gap between object and image-level representations for open-vocabulary detection. *Advances in Neural Information Processing Systems*, 35: 33781–33794.

Brown, T.; Mann, B.; Ryder, N.; Subbiah, M.; Kaplan, J. D.; Dhariwal, P.; Neelakantan, A.; Shyam, P.; Sastry, G.; Askell, A.; et al. 2020. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901.

Dai, A.; Chang, A. X.; Savva, M.; Halber, M.; Funkhouser, T.; and Nießner, M. 2017. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 5828–5839.

De Rijk, P.; Schneider, L.; Cordts, M.; and Gavrila, D. 2022. Structural Knowledge Distillation for Object Detection. *Advances in Neural Information Processing Systems*, 35: 3858–3870.

Guo*, Z.; Zhang*, R.; Qiu*, L.; Ma, X.; Miao, X.; He, X.; and Cui, B. 2022. CALIP: Zero-Shot Enhancement of CLIP with Parameter-free Attention. *AAAI 2023 Oral*.

Guo, Z.; Zhang, R.; Zhu, X.; Tang, Y.; Ma, X.; Han, J.; Chen, K.; Gao, P.; Li, X.; Li, H.; et al. 2023. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*.

Gupta, A.; Dollar, P.; and Girshick, R. 2019. Lvis: A dataset for large vocabulary instance segmentation. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, 5356–5364.

Hu, C.; and Li, X. 2023. When sam meets medical images: An investigation of segment anything model (sam) on multi-phase liver tumor segmentation. *arXiv preprint arXiv:2304.08506*.

Kirillov, A.; Mintun, E.; Ravi, N.; Mao, H.; Rolland, C.; Gustafson, L.; Xiao, T.; Whitehead, S.; Berg, A. C.; Lo, W.-Y.; et al. 2023. Segment anything. *arXiv preprint arXiv:2304.02643*.

Liang, F.; Wu, B.; Dai, X.; Li, K.; Zhao, Y.; Zhang, H.; Zhang, P.; Vajda, P.; and Marculescu, D. 2023. Open-vocabulary semantic segmentation with mask-adapted clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 7061–7070.

Liu, S.; Zeng, Z.; Ren, T.; Li, F.; Zhang, H.; Yang, J.; Li, C.; Yang, J.; Su, H.; Zhu, J.; et al. 2023a. Grounding dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv preprint arXiv:2303.05499*.

Liu, Y.; Zhang, J.; She, Z.; Kheradmand, A.; and Armand, M. 2023b. Samm (segment any medical model): A 3d slicer integration to sam. *arXiv preprint arXiv:2304.05622*.

Liu, Z.; Zhang, Z.; Cao, Y.; Hu, H.; and Tong, X. 2021. Group-free 3d object detection via transformers. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2949–2958.

Lu, Y.; Xu, C.; Wei, X.; Xie, X.; Tomizuka, M.; Keutzer, K.; and Zhang, S. 2022. Open-vocabulary 3d detection via image-level class and debiased cross-modal contrastive learning. *arXiv preprint arXiv:2207.01987*.

Lu, Y.; Xu, C.; Wei, X.; Xie, X.; Tomizuka, M.; Keutzer, K.; and Zhang, S. 2023. Open-vocabulary point-cloud object detection without 3d annotation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 1190–1199.

Misra, I.; Girdhar, R.; and Joulin, A. 2021. An end-to-end transformer model for 3d object detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2906–2917.

Ni, X.; and Li, P. 2023. Unified Text Structuralization with Instruction-tuned Language Models. *arXiv preprint arXiv:2303.14956*.

Oord, A. v. d.; Li, Y.; and Vinyals, O. 2018. Representation learning with contrastive predictive coding. *arXiv preprint arXiv:1807.03748*.

Pratt, S.; Covert, I.; Liu, R.; and Farhadi, A. 2023. What does a platypus look like? Generating customized prompts for zero-shot image classification. arXiv:2209.03320.

Qi, C. R.; Litany, O.; He, K.; and Guibas, L. J. 2019. Deep hough voting for 3d object detection in point clouds. In *proceedings of the IEEE/CVF International Conference on Computer Vision*, 9277–9286.

Radford, A.; Kim, J. W.; Hallacy, C.; Ramesh, A.; Goh, G.; Agarwal, S.; Sastry, G.; Askell, A.; Mishkin, P.; Clark, J.; et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*, 8748–8763. PMLR.

Rahman, S.; Khan, S.; and Barnes, N. 2020. Improved visual-semantic alignment for zero-shot object detection. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, 11932–11939.

Rahman, S.; Khan, S. H.; and Porikli, F. 2020. Zero-shot object detection: Joint recognition and localization of novel concepts. *International Journal of Computer Vision*, 128: 2979–2999.

Russakovsky, O.; Deng, J.; Su, H.; Krause, J.; Satheesh, S.; Ma, S.; Huang, Z.; Karpathy, A.; Khosla, A.; Bernstein, M.; Berg, A. C.; and Fei-Fei, L. 2015. ImageNet Large Scale Visual Recognition Challenge. *International Journal of Computer Vision (IJCV)*, 115(3): 211–252.

Song, S.; Lichtenberg, S. P.; and Xiao, J. 2015. Sun rgb-d: A rgb-d scene understanding benchmark suite. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, 567–576.

Wang, S.; Sun, X.; Li, X.; Ouyang, R.; Wu, F.; Zhang, T.; Li, J.; and Wang, G. 2023. GPT-NER: Named Entity Recognition via Large Language Models. *arXiv preprint arXiv:2304.10428*.

Xu, C.; Yang, S.; Galanti, T.; Wu, B.; Yue, X.; Zhai, B.; Zhan, W.; Vajda, P.; Keutzer, K.; and Tomizuka, M. 2021. Image2Point: 3D Point-Cloud Understanding with 2D Image Pretrained Models. *arXiv preprint arXiv:2106.04180*.

Zareian, A.; Rosa, K. D.; Hu, D. H.; and Chang, S.-F. 2021. Open-vocabulary object detection using captions. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 14393–14402.

Zhang, R.; Guo, Z.; Zhang, W.; Li, K.; Miao, X.; Cui, B.; Qiao, Y.; Gao, P.; and Li, H. 2022. Pointclip: Point cloud understanding by clip. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 8552–8562.

Zhang, R.; Hu, X.; Li, B.; Huang, S.; Deng, H.; Li, H.; Qiao, Y.; and Gao, P. 2023a. Prompt, Generate, then Cache: Cascade of Foundation Models makes Strong Few-shot Learners. *CVPR 2023*.

Zhang, R.; Jiang, Z.; Guo, Z.; Yan, S.; Pan, J.; Dong, H.; Gao, P.; and Li, H. 2023b. Personalize segment anything model with one shot. *arXiv preprint arXiv:2305.03048*.

Zhang, R.; Wang, L.; Qiao, Y.; Gao, P.; and Li, H. 2023c. Learning 3D Representations from 2D Pre-trained Models via Image-to-Point Masked Autoencoders. *CVPR 2023*.

Zhang, Z.; Sun, B.; Yang, H.; and Huang, Q. 2020. H3dnet: 3d object detection using hybrid geometric primitives. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XII 16*, 311–329. Springer.

Zhou, X.; Girdhar, R.; Joulin, A.; Krähenbühl, P.; and Misra, I. 2022. Detecting twenty-thousand classes using image-level supervision. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part IX*, 350–368. Springer.

Zhu, X.; Zhang, R.; He, B.; Guo, Z.; Zeng, Z.; Qin, Z.; Zhang, S.; and Gao, P. 2023. Pointclip v2: Prompting clip and gpt for powerful 3d open-world learning. *ICCV 2023*.