

# ShapeLLM: Universal 3D Object Understanding for Embodied Interaction

Zekun Qi<sup>1¶</sup> Runpei Dong<sup>1¶F</sup> Shaochen Zhang<sup>1</sup> Haoran Geng<sup>2</sup> Chunrui Han<sup>3</sup> Zheng Ge<sup>3</sup>  
 He Wang<sup>2</sup> Li Yi<sup>4†</sup> Kaisheng Ma<sup>4†</sup>

<sup>¶</sup> Internship at MEGVII <sup>F</sup> Project lead <sup>†</sup> Corresponding authors

<sup>1</sup> Xi'an Jiaotong University <sup>2</sup> Peking University <sup>3</sup> MEGVII Technology <sup>4</sup> Tsinghua University

<https://qizekun.github.io/ShapeLLM>

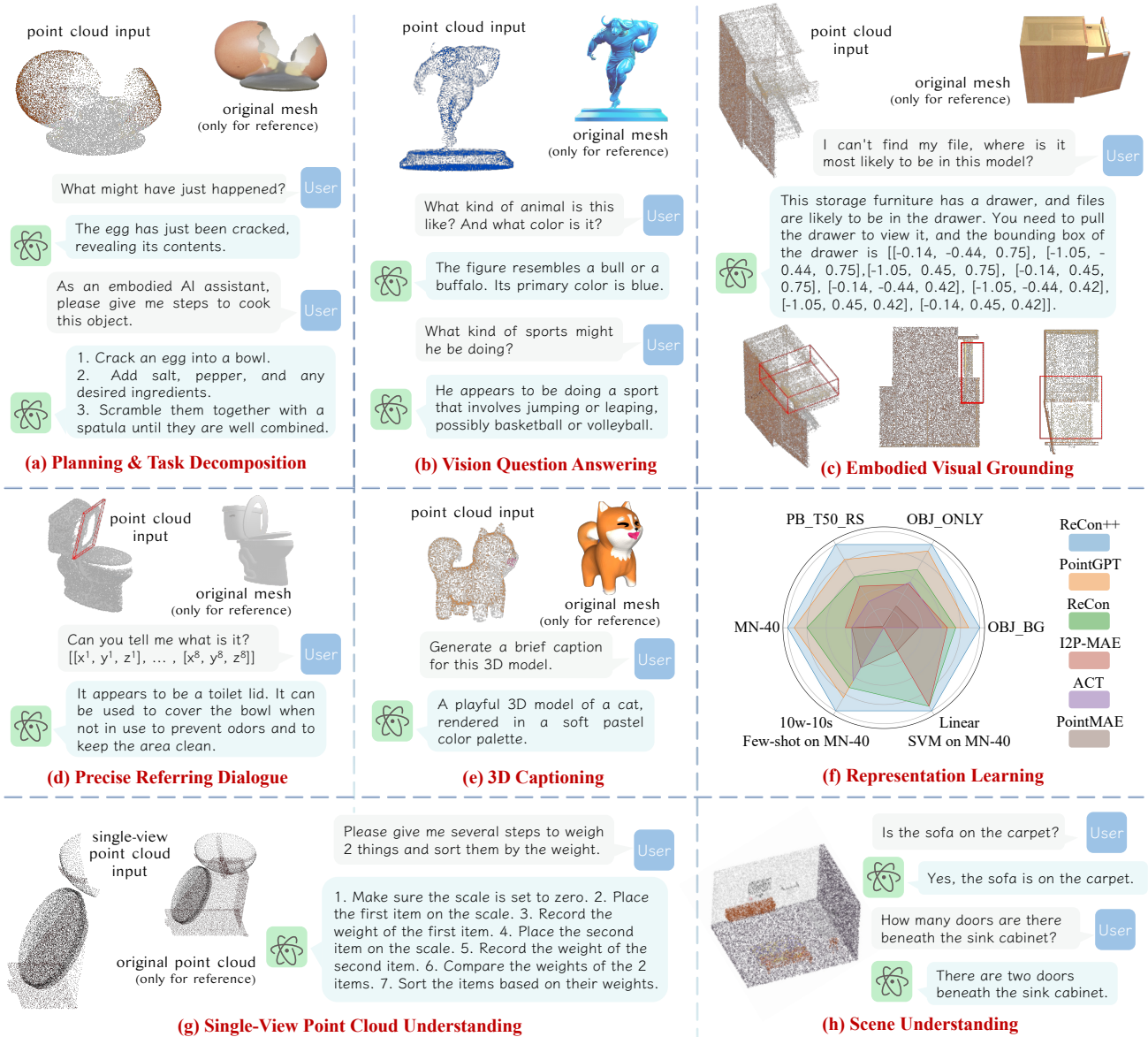


Figure 1. **Demonstrations of SHAPELLM.** We present SHAPELLM, a multi-modal large language model designed for embodied scenes. SHAPELLM can generate accurate and prompt responses to user queries by leveraging comprehensive prior knowledge of the language. This figure demonstrates the powerful capabilities of SHAPELLM in *planning & task decomposition*, *3D vision question answering* and *embodied vision grounding*, etc. The predicted bounding box is shown in the red box from three distinct perspectives.

## Abstract

This paper presents SHAPeLLM, the first 3D Multimodal Large Language Model (LLM) designed for embodied interaction, exploring a universal 3D object understanding with 3D point clouds and languages. SHAPeLLM is built upon an improved 3D encoder by extending RECON [129] to RECON++ that benefits from multi-view image distillation for enhanced geometry understanding. By utilizing RECON++ as the 3D point cloud input encoder for LLMs, SHAPeLLM is trained on constructed instruction-following data and tested on our newly human-curated evaluation benchmark, 3D MM-Vet. RECON++ and SHAPeLLM achieve state-of-the-art performance in 3D geometry understanding and language-unified 3D interaction tasks, such as embodied visual grounding.

## 1. Introduction

3D shape understanding, serving as a fundamental capability for molding intelligent systems in both digital and physical worlds, has witnessed tremendous progress in graphics, vision, augmented reality, and embodied robotics. However, to be effectively deployed by real-world agents, several critical criteria must be fulfilled: (i) Sufficient 3D geometry information needs to be captured for accurate spatial and structure processing [10, 13, 79, 126]. (ii) Models should be endowed with a foundational knowledge of the embodied interaction fashion with objects — often physically — for functional comprehension [52, 65–67, 80, 125, 192, 193]. (iii) A universal interface is required as a bridge between information encoding and decoding, which could help translate high-order instructions for agent reactions like dialogue response and embodied feedback [27, 71, 194].

Recent advancements in Large Language Models (LLMs) [11, 116, 132, 133, 150] have demonstrated unprecedented success of foundational knowledge and unified reasoning capabilities across tasks [7, 20, 28, 37, 39, 70, 73, 78, 124]. It makes it possible to utilize language as a universal interface that enables the comprehensive commonsense knowledge embedded in LLMs to enhance understanding of 3D shapes. This is particularly evident in physically-grounded tasks, where the wealth of commonsense knowledge simplifies the interpretation of an object’s functionality, mobility, and dynamics, etc. However, the aforementioned challenges remain when incorporating LLMs for 3D object understanding — especially embodied interaction that relies on precise geometry — currently under-explored.

The question is: *What makes better 3D representations that bridge language models and interaction-oriented 3D object understanding?* In this work, we introduce SHAPeLLM that meets the requirements, which is established based on the following three designing policies:

- i. **3D Point Clouds as Inputs** Some concurrent works [54] recently propose to use point cloud-rendered images [186] as multimodal LLMs’ inputs and demonstrate effectiveness. However, these works fail to achieve accurate 3D geometry understanding and often suffer from a well-known visual hallucination issue [87, 138, 196]. Compared to 2D images, 3D point clouds provide a more accurate representation of the physical environment, encapsulating sparse yet highly precise geometric data [1, 36, 127]. Moreover, 3D point clouds are crucial in facilitating embodied interactions necessitating accurate 3D structures like 6-DoF object pose estimation [84, 154, 156, 160, 167].
- ii. **Selective Multi-View Distillation** Interacting with objects typically necessitates an intricate 3D understanding that involves knowledge at various levels and granularities. For instance, a whole-part *high-level* semantic understanding is needed for interactions like opening a large cabinet, while detailed, *high-resolution* (i.e., *low-level*) semantics are crucial for smaller objects like manipulating a drawer handle [177]. However, existing works mainly distill single-view high-resolution object features from 2D foundation models [134], providing a complementary understanding [36, 129, 169]. The potential of multi-view images, which offer abundant multi-level features due to view variation and geometry consistency [9, 58, 63, 79, 101, 143], is often neglected. SHAPeLLM extends RECON [129] to RECON++ as the 3D encoder by integrating multi-view distillation. To enable the model to selectively distill views that enhance optimization and generalization, inspired by Carion et al., RECON++ is optimized through adaptive selective matching using the Hungarian algorithm [82].
- iii. **3D Visual Instruction Tuning** Instruction tuning has been proven effective in improving LLMs’ alignment capability [117, 121]. To realize various 3D understanding tasks with a universal language interface, SHAPeLLM is trained through instruction-following tuning on constructed language-output data. However, similar to 2D visual instruction tuning [4, 93], the data-dessert issue [36] is even worse since no object-level VQA data is available, unlike 2D [92]. To validate the efficacy of SHAPeLLM, we first construct ~45K instruction-following data using the advanced GPT-4V(ision) [115] on the processed Objaverse dataset [29] and 30K embodied part understanding data from GPartNet [47] for supervised fine-tuning. Following MM-Vet [179], we further develop a novel evaluation benchmark named 3D MM-Vet. This benchmark is designed to assess the core vision-language capabilities, including embodied interaction in a 3D context, thereby stimulating future research. The 3D MM-Vet benchmark comprises 59 diverse Internet<sup>1</sup> 3D objects and 232 human-written question-answer pairs.

<sup>1</sup>URL & License.

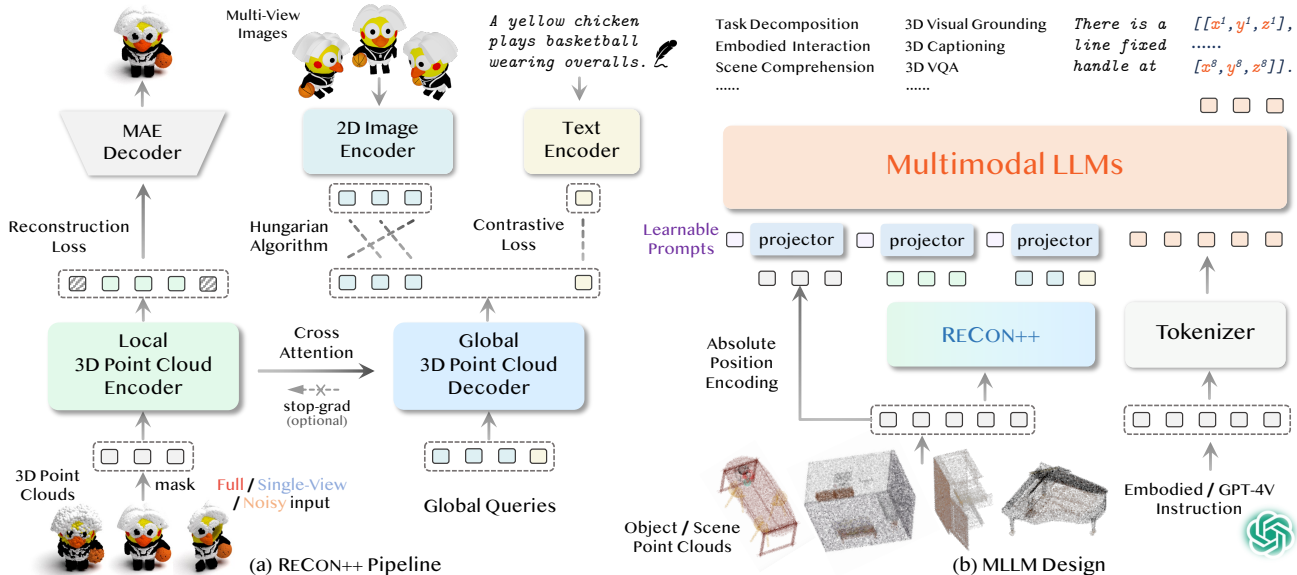


Figure 2. **Overview of our SHAPeLLM framework.** (a) The introduced RECON++ pipeline incorporates the required 3D encoder. (b) The comprehensive design of the MLLM, featuring an instruction-mode tokenizer and the integration of an aligned multi-modal representation, equips the MLLM with the capability to effectively handle 3D vision language tasks.

Through extensive experimentation, we first demonstrate that our improved 3D encoder RECON++ sets a new state-of-the-art representation transferring on both downstream fine-tuned and zero-shot 3D object recognition. Specifically, RECON++ has obtained **95.25%** and **95.0%** fine-tuned accuracy on ScanObjectNN and ModelNet40, surpassing previous best records by **+1.85%** on the most challenging ScanObjectNN. Besides, RECON++ achieved **53.7%** and **65.4%** zero-shot accuracy on Objaverse-LVIS and ScanObjectNN, which is **+0.6%** and **+1.6%** higher than previous best. By utilizing our RECON++ as SHAPeLLM’s 3D encoder, SHAPeLLM successfully unifies various downstream tasks, including *3D captioning*, *3D VQA*, *embodied task planning & decomposition*, *3D embodied visual grounding*, and *3D precise referring dialogue* (See Fig. 1). On our newly constructed 3D MM-Vet benchmark, **42.7%** and **49.3%** Total accuracy have been achieved by SHAPeLLM-7B and SHAPeLLM-13B, surpassing previous best records [166] that also uses 3D point clouds by **+2.1%** and **+5.1%**, respectively. This work initiates a first step towards leveraging LLMs for embodied object interaction, and we hope our SHAPeLLM and proposed 3D MM-Vet benchmark could spur more related future research.

## 2. SHAPeLLM

In this section, we first introduce the overall architecture of SHAPeLLM. Then, we delve into two critical challenges faced in interactive 3D understanding: data dessert [36] and representation of 3D point clouds. We present the detailed design of our method to tackle these challenges, respectively.

### 2.1. Overall Architecture

The main objective of this work is interactive 3D understanding by using the LLM as a universal interface. Drawing inspiration from recent work in visual understanding [93], the proposed SHAPeLLM consists a pre-trained 3D encoder and an LLM for effective 3D representation learning and understanding, respectively. Specifically, we adopt LLaMA [150] as our LLM, building upon the success of previous work [24, 37, 93]. As for the 3D encoder, we propose a novel 3D model named RECON++ based on the recent work RECON [129] with multiple improvements as the 3D understanding generally demands more information, such as accurate spatial and multi-view details, etc. To ensure compatibility with the LLM inputs, the representation of a 3D object obtained from RECON++ undergoes a linear projection before being fed into the LLM. To further improve low-level geometry understanding, which benefits tasks like 6-DoF pose estimation, we append the absolute position encoding (APE) obtained by linear projection of 3D coordinates. Besides, we use prefix-tuning with learnable prompts [36, 37, 76, 86] to adaptively modulate the different semantics of APE and RECON++ representations.

### 2.2. How to alleviate interactive 3D understanding Data Dessert?

Most published 3D data is typically presented as 3D object-caption pairs, lacking an interactive style. Although a few concurrent works [62, 166] have attempted to construct interactive 3D understanding datasets, the questions-and-answers (Q&As) are primarily based on annotated captions, often pro-

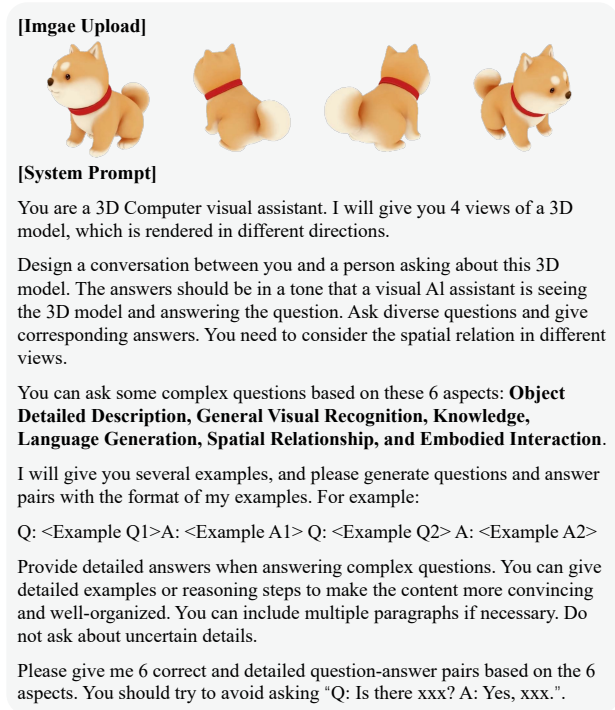


Figure 3. **Construction illustration of instruct-following data using GPT-4V [115].** Four perspective views are input into GPT-4V. In-context prompts focusing on different topics are explicitly incorporated to ensure data diversity.

viding a limited perspective without sufficient details. Additionally, those works have generally been limited to semantic understanding without considering embodied interaction. To address these limitations, our work constructs question-and-answer pairs based on multi-view images of a 3D object using GPT-4V(ision) [115]. For data diversity, we explicitly introduce six aspects as prompts, as illustrated Fig. 3. In the following, we provide the details about data collection and construction regarding *general semantic understanding* and *embodied object understanding*, respectively.

**Data** Objaverse-LVIS [29, 107] and GPartNet [47] are data sources. Objaverse-LVIS covers 1,156 LVIS [55] categories, and we sample Top-10 “likes”<sup>2</sup> 3D objects per category and generate Q&A pairs per sample. After filtering out noisy Q&As, we obtain ~45K instruction-following samples. We use 12 categories from GPartNet by removing “Remote” to avoid too many tiny boxes, which leads to filtered ~30K Q&A samples constructed from ~8K parts of the ~4K objects states covering ~1.1K different objects.

**General Semantic Understanding** This aims to enhance the model’s generalization abilities in visual recognition, knowledge integration, spatial understanding, and other aspects. We prompt GPT4-V to generate Q&As in six different aspects based on images captured from four different views of a 3D subject, as illustrated in Fig. 3.

<sup>2</sup>“Likes” statistics can be found at [Sketchfab](#).

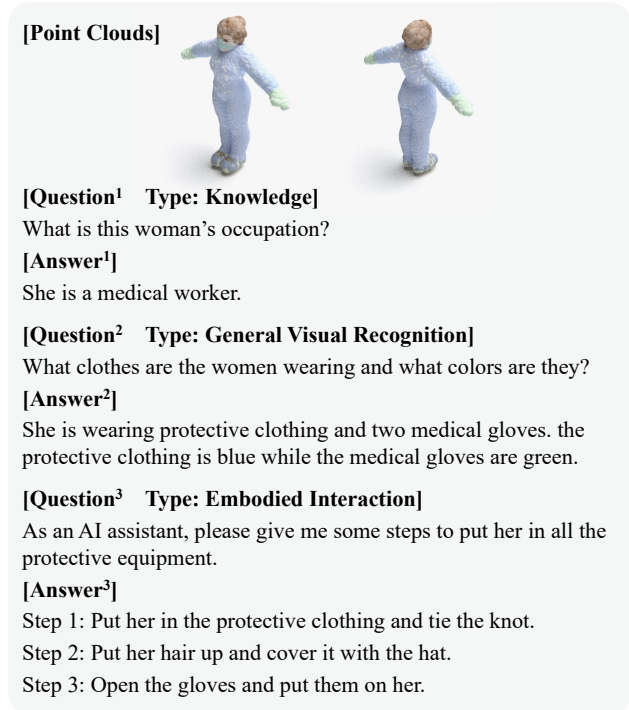


Figure 4. **3D MM-Vet dataset sample.** A wealth of precise evaluation metrics enable a comprehensive assessment.

**Embodied Object Understanding** A comprehensive understanding of the spatial positions and semantics at the part level is crucial to facilitate effective object grasping and interaction in embodied scenarios. Fortunately, the GPartNet [47] provides rich part annotations, including semantics and poses, which are instrumental in constructing instruction-tuning data for embodied interactive parts of a subject. Specifically, given a 3D object, questions are formulated based on the semantics of its different parts, and answers are constructed in both the semantics and 3D positions. The positions are represented as 6-DoF 3D bounding boxes in a straightened Python multidimensional list format, denoted as  $[[x_1, y_1, z_1], [x_2, y_2, z_2], \dots, [x_8, y_8, z_8]]$ , to meet characteristics of the textual dialogues response in LLMs. The canonical space of the object determines the sequence of coordinates. Using bounding box coordinates leverages the inherent spatial relationship, allowing LLMs to readily learn these patterns and generate accurate output coordinates. This approach can offer specific position information for embodied manipulation.

### 2.3. RECON++: Scaling Up 3D Representation Learning

Interaction with objects such as object grasping [95, 154, 167] typically requires accurate perception of 3D shape information at multi-level and multi-granularity. This imposes heightened requirements on 3D representations, calling for a higher standard of a holistic understanding of 3D geometry.

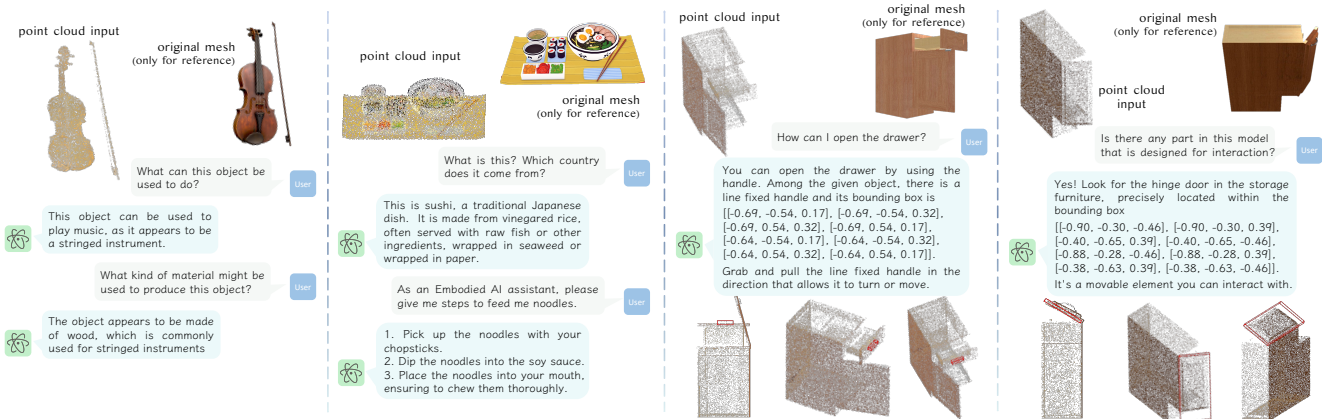


Figure 5. **Selected multimodal dialogue examples.** SHAPPELLM possesses robust capabilities in knowledge representation, reasoning, and instruction-following dialogue. With its powerful point cloud encoder RECON++, SHAPPELLM can even make accurate predictions about minute interactive components, *e.g.*, handle. The rendered mesh images are solely for visual reference here and do not constitute input data.

However, existing 3D cross-modal representation learning methods [94, 170] mainly distill high-resolution object features from single-view 2D foundation models, resulting in a unilateral shape understanding. Besides, they generally employ multi-view images as a data augmentation strategy, imposing the learned representation to the average representation of all views. Thus, the accurate 3D shape information is missing. Recently, RECON [129] utilizes contrast guided by reconstruction to address the pattern disparities between local masked data modeling and global cross-modal alignment. This results in remarkable performance in various tasks, including transfer learning, zero-shot classification, and part segmentation. However, its potential is hindered by the scarcity of pretraining data [13].

To address the above limitations, this paper proposes RECON++ with multiple improvements. First, multi-view image query tokens collaboratively comprehend the semantic information of 3D objects across different views, encompassing both RGB images and depth maps. Considering the disorderliness of pretraining data in terms of pose, we propose a cross-modal alignment method based on *bipartite matching*, which implicitly learns the pose estimation of 3D objects. Second, we *scale up* the parameters of RECON and broaden the scale of the pretraining dataset [18, 29, 107] for robust 3D representations.

Denote  $N$  as the number of multi-view images,  $I_i$  is the image feature from  $i$ -th view, and  $Q_i$  represents the global query of  $i$ -th view. Following Carion et al., we search for an optimal permutation  $\sigma$  of  $N$  elements with the lowest cost:

$$\hat{\sigma} = \arg \min_{\sigma} \sum_i^N \mathcal{L}_{\text{match}}(I_i, Q_{\sigma(i)}), \quad (1)$$

where  $\mathcal{L}_{\text{match}}(I_i, Q_{\sigma(i)})$  is a pair-wise matching cost between  $i$ -th view image features  $I_i$  and matched query  $Q_{\sigma(i)}$  with

the permutation  $\sigma$ . In practice, we employ cosine similarity as the matching cost. In this fashion, the query of each view is learned to gather accurate 3D shape information from the 3D point clouds. Concatenating the features from the local 3D point cloud encoder and global 3D point cloud decoder together provides comprehensive information for 3D understanding of multimodal LLMs.

### 3. 3D MM-Vet: 3D Multimodal Comprehension Evaluation Benchmark

A wide range of diverse visual-language capabilities is essential to develop a multimodal large language model tailored for embodied scenarios, particularly addressing task and action planning.

The model’s proficiency in processing point clouds enables it to perform general recognition tasks effortlessly, demonstrating a broad understanding of colored point clouds. This capability serves as the groundwork for more intricate tasks. Beyond 3D recognition, the LLM should exhibit competence in addressing tasks in real-world embodied scenarios. This entails unifying the aforementioned abilities to generate decomposed task actions step-by-step in an instruction-following fashion, addressing specific problems.

Hence, to formulate an evaluation system aligned with the aforementioned task description, we establish a multi-level evaluation task system encompassing four-level tasks: **General Recognition, Knowledge and Language Generation, Spatial Awareness, and Embodied Interaction**. This framework systematically and comprehensively assesses the model’s proficiency in information comprehension and language generation when processing interactive objects. The detailed descriptions of the tasks are listed as follows:

- i. **General Recognition:** Following MM-Vet [179], we assess the fundamental comprehension abilities of LLMs

Table 1. **Fine-tuned 3D recognition** on ScanObjectNN and ModelNet40. BG, ON, and RS are short for OBJ\_BG, OBJ\_ONLY, and PB\_T50\_RS, respectively. Overall accuracy (%) with voting [98] is reported. †: results with a post-pretraining stage [18].

Method	ScanObjectNN			ModelNet40	
	BG	ON	RS	1k P	8k P
<i>Supervised Learning Only</i>					
PointNet [126]	73.3	79.2	68.0	89.2	90.8
PointNet++ [127]	82.3	84.3	77.9	90.7	91.9
DGCNN [157]	82.8	86.2	78.1	92.9	-
PointMLP [108]	-	-	85.4	94.5	-
PointNetXt [131]	-	-	87.7	94.0	-
Transformer [153]	83.04	84.06	79.11	91.4	91.8
<i>with Self-Supervised Representation Learning</i>					
Point-BERT [180]	87.43	88.12	83.07	93.2	93.8
Point-MAE [119]	90.02	88.29	85.18	93.8	94.0
Point-M2AE [185]	91.22	88.81	86.43	94.0	-
ACT [36]	93.29	91.91	88.21	93.7	94.0
TAP [158]	-	-	88.5	94.0	-
VPP [130]	93.11	91.91	89.28	94.1	94.3
I2P-MAE [188]	94.15	91.57	90.11	94.1	-
ULIP-2 [170]	-	-	91.5	-	-
RECON [129]	95.35	93.80	91.26	94.5	94.7
PointGPT-B† [18]	95.8	95.2	91.9	94.4	94.6
PointGPT-L† [18]	97.2	96.6	93.4	94.7	94.9
<b>RECON++-B†</b>	<b>98.62</b>	<b>96.21</b>	<b>93.34</b>	<b>94.6</b>	<b>94.8</b>
<b>RECON++-L†</b>	<b>98.80</b>	<b>97.59</b>	<b>95.25</b>	<b>94.8</b>	<b>95.0</b>

involving both coarse- and fine-grained aspects. Coarse-grained recognition focuses on basic object attributes such as color, shape, action, *etc.* While fine-grained recognition delves into details like subparts and counting, *etc.*

- ii. **Knowledge Capability & Language Generation:** To examine the models’ capacity to understand and utilize knowledge, drawing inspiration from MMBench [100], we integrate its reasoning components. This includes knowledge spanning natural and social reasoning, physical properties, sequential prediction, math, *etc.*, evaluating gauges whether multimodal LLMs possess the requisite expertise and capacity to solve intricate tasks. We utilize customized prompts to stimulate models and extract detailed responses to evaluate language generation.
- iii. **Spatial Awareness:** In 3D, spatial awareness holds heightened significance compared to 2D due to the provided geometry information. The point clouds contain location information crucial for discerning spatial relationships between different parts. In 2D, achieving the same information intensity level would necessitate multi-view images. Therefore, our evaluation includes questions probing the ability of LLMs to understand spatial relations.
- iv. **Embodied Interaction:** The utilization scope of multimodal LLMs extends into the field of embodied interaction, facilitated by the utilization of instruction-following data. Our evaluation system tests their capacity by formally

Table 2. **Zero-shot 3D recognition** on Objaverse-LVIS [29], ModelNet40 [164] and ScanObjectNN [151]. Ensembled: pretraining with four datasets, Objaverse [29], ShapeNet [13], ABO [22] and 3D-FUTURE [42], following OpenShape [94]. †: Uni3D employs an EVA-CLIP-E [146] teacher that is larger and improves performance, while other methods employ OpenCLIP-bigG [74].

Method	Objaverse-LVIS			ModelNet40			ScanObjectNN		
	Top1	Top3	Top5	Top1	Top3	Top5	Top1	Top3	Top5
<i>2D Inference</i>									
PointCLIP [186]	1.9	4.1	5.8	19.3	28.6	34.8	10.5	20.8	30.6
PointCLIP2 [198]	4.7	9.5	12.9	63.6	77.9	85.0	42.2	63.3	74.5
<i>Trained on ShapeNet</i>									
RECON [129]	1.1	2.7	3.7	61.2	73.9	78.1	42.3	62.5	75.6
CLIP2Point [69]	2.7	5.8	7.9	49.5	71.3	81.2	25.5	44.6	59.4
ULIP [169]	6.2	13.6	17.9	60.4	79.0	84.4	51.5	71.1	80.2
OpenShape [94]	10.8	20.2	25.0	70.3	86.9	91.3	47.2	72.4	84.7
<i>Trained on Ensembled</i>									
ULIP-2 [170]	26.8	44.8	52.6	75.1	88.1	93.2	51.6	72.5	82.3
OpenShape [94]	46.8	69.1	77.0	84.4	96.5	98.0	52.2	79.7	88.7
Uni3D-B† [195]	51.7	74.1	80.8	86.3	<b>96.5</b>	<b>97.9</b>	63.8	<b>82.7</b>	90.2
Uni3D-L† [195]	53.1	75.0	81.5	86.3	<b>96.8</b>	<b>98.3</b>	58.2	81.8	89.4
<b>RECON++-B</b>	<b>53.2</b>	<b>75.3</b>	<b>81.5</b>	<b>86.5</b>	94.7	95.8	<b>63.6</b>	80.2	<b>90.6</b>
<b>RECON++-L</b>	<b>53.7</b>	<b>75.8</b>	<b>82.0</b>	<b>87.3</b>	95.4	96.1	<b>65.4</b>	<b>84.1</b>	<b>89.7</b>

requesting LLMs to provide execution steps toward an instruction. This approach aims to establish connections for handling Embodied Interaction tasks [39, 70].

To prevent any overlap with training data, our collection of 3D models is sourced exclusively from Turbosquid [142], a platform not included in the acquisition lists of Objaverse [29] and ShapeNet [13]. We meticulously curated a dataset of 59 3D models, generating 232 Q&As for evaluation purposes. In our pursuit of a precise assessment of single-task capabilities, each question is designed to test only one specific capacity outlined earlier. Every question is paired with a corresponding answer tailored to the particular 3D model, serving as the ground truth. More details and analysis can be found in Appendix B.

## 4. Experiments

### 4.1. 3D Representation Transferring with RECON++

**Fine-tuned 3D Object Recognition** In Tab. 1, we first evaluate the representation transfer learning capabilities of self-supervised RECON++ by fine-tuning on ScanObjectNN [151] and ModelNet [164], which are currently the two most challenging 3D object datasets. ScanObjectNN is a collection of  $\sim 15K$  3D object point clouds from the real-world scene dataset ScanNet [23], which involves 15 categories. ModelNet is one of the most classical 3D object datasets collected from clean 3D CAD models, which includes  $\sim 12K$  meshed 3D CAD models covering 40 categories. Following PointGPT [18], we adopt the intermediate fine-tuning strategy and use the post-pretraining

Table 3. **Zero-shot 3D multimodal comprehension evaluation of core VL capabilities in 3D context** on 3D MM-Vet. **Rec**: General Visual Recognition, **Know**: Knowledge, **Gen**: Language Generation, **Spat**: Spatial Awareness, **Emb**: Embodied Interaction.

Method	Input	Rec	Know	Gen	Spat	Emb	Total
LLaVA-13B [93]	1-View Img.	40.0	55.3	51.3	43.2	51.1	47.9
DreamLLM-7B [37]	4-View Img.	42.2	54.4	50.8	48.9	54.5	50.3
GPT-4V [115]	1-View Img.	53.7	59.5	61.1	54.7	59.0	57.4
GPT-4V [115]	4-View Img.	65.1	69.1	61.4	52.9	65.5	63.4
PointBind&LLM [54]	Point Cloud	16.9	13.0	18.5	32.9	40.4	23.5
PointLLM-7B [166]	Point Cloud	40.6	49.5	34.3	29.1	48.7	41.2
PointLLM-13B [166]	Point Cloud	46.6	48.3	38.8	45.2	50.9	46.6
<b>SHAPELLM-7B</b>	Point Cloud	<b>45.7</b>	<b>42.7</b>	<b>43.4</b>	<b>39.9</b>	<b>64.5</b>	<b>47.4</b>
<b>SHAPELLM-13B</b>	Point Cloud	<b>46.8</b>	<b>53.0</b>	<b>53.9</b>	<b>45.3</b>	<b>68.4</b>	<b>53.1</b>

stage to transfer the general semantics learned through self-supervised pretraining on ShapeNetCore [13]. For a fair comparison, our Base and Large models adopt the same architecture as PointGPT regarding layers, hidden size, and attention heads. Tab. 1 shows that: (i) RECON++ exhibits representation performance significantly surpassing that of other baselines, achieving state-of-the-art results. (ii) Particularly, RECON++ achieves a remarkable accuracy of 95.25% on the most challenging ScanObjectNN PB\_T50\_RS benchmark, boosting the Transformer baseline by +16.14%.

**Zero-Shot 3D Open-World Recognition** Similar to CLIP [134], our model aligns the feature space of languages and other modalities, which results in a zero-shot open-world recognition capability. In Tab. 2, we compare the zero-shot 3D open-world object recognition models to evaluate the generalizable recognition capability. Following OpenShape [94], we evaluate on ModelNet [164], ScanObjectNN [151], and Objaverse-LVIS [29]. Objaverse-LVIS is a benchmark involving  $\sim 47K$  clean 3D models of 1,156 LVIS categories [55]. We compare RECON++ with 2D inference methods, ShapeNet pretrained methods, and “Ensembled” datasets-pretrained methods. It can be concluded from Tab. 2: i) Compared to 2D inference and ShapeNet-pretrained methods, RECON++ demonstrates significantly superior performance, showing the necessity of *3D point clouds as inputs* and *scaling up*. ii) Compared to state-of-the-art methods trained on “Ensembled” datasets, RECON++ demonstrates superior or on-par performance across all benchmarks. Notably, RECON++-L achieves a remarkable Top-1 accuracy, which is +0.6% and +7.2% higher than Uni3D-L on the most challenging Objaverse-LVIS and ScanObjectNN benchmarks, respectively.

## 4.2. Multimodal Comprehension with SHAPELLM

**Quantitative Analysis** To assess the comprehensive capabilities of SHAPELLM, we first quantitatively compare various baselines and our model on the proposed 3D MM-Vet benchmark using GPT-4. Following ModelNet-C [136] and ModelNet40-C [144], we construct 3D MM-Vet-C to benchmark the robustness against 3D corruptions.

Table 4. **Zero-shot 3D multimodal comprehension evaluation of robustness** on 3D MM-Vet-C. **Clean**: no corruptions. **Single-View**: randomly select a camera viewpoint within the unit sphere and generate a **single viewpoint** within the FoV on polar coordinates. **Jitter**: Gaussian jittering with noise  $\epsilon \sim \mathcal{N}(0, \sigma^2)$  and  $\sigma = 0.01$ . **Rotate**: random SO(3) rotation sampling over X-Y-Z Euler angle  $(\alpha, \beta, \gamma) \sim \mathcal{U}(-\theta, \theta)$  and  $\theta = \pi/6$ .

Method	3D MM-Vet-C Variants			
	Clean	Single-View	Jitter	Rotate
PointBind&LLM [54]	23.5	20.4	19.7	19.5
PointLLM-7B [166]	41.2	33.6	38.8	40.6
PointLLM-13B [166]	46.6	41.3	42.3	44.2
<b>SHAPELLM-7B</b>	<b>47.4</b>	<b>38.3</b>	<b>45.8</b>	<b>42.7</b>
<b>SHAPELLM-13B</b>	<b>53.1</b>	<b>43.6</b>	<b>47.8</b>	<b>49.3</b>

- **3D MM-Vet.** Tab. 3 shows the detailed results of SHAPELLM on different tasks of 3D MM-Vet. It can be observed that SHAPELLM significantly outperforms PointLLM [166] across various metrics, particularly in Embodied Tasks. This substantiates our model’s versatile capability in addressing real-world scenario tasks.
- **3D MM-Vet-C.** Tab. 4 shows the comparison of model robustness against “single-view”, “jitter” and “rotate” corruptions, which are the most common corruptions in real-world scenarios. The results demonstrate significantly superior robustness of SHAPELLM against corruption, indicating stronger potential in real-world applicability.

**Qualitative Analysis** Fig. 5 illustrates qualitative examples of SHAPELLM in multimodal dialogue. SHAPELLM is capable of supporting general VQA, embodied task and action planning, as well as 6-DoF pose estimation. Notably, due to the strict spatial relationship inherent in 6-DoF bounding box coordinates, we observe that LLMs easily grasp such patterns and consistently produce valid coordinates.

## 5. Discussions

### 5.1. What is learned from multi-view alignment?

Fig. 6 illustrates the visualization of the attention map in the last cross-attention layer, documenting the image query to which each local patch in the attention map primarily attends. It provides evidence that multi-view alignment achieves geometrically informed spatial understanding, which may im-

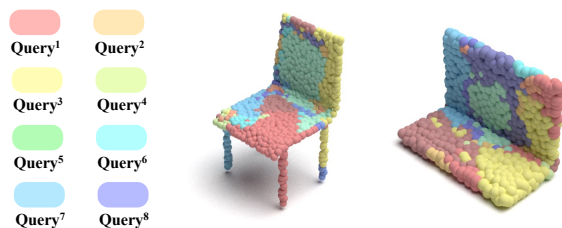








Figure 6. **Visualization of multi-view query results.** The distinct colors serve to denote distinct image queries.

Table 5. **3D referring expression grounding** on GPartNet. Accuracy with an IoU threshold of 0.25 is reported. †: Fine-tuned on GPartNet images. ‡: Inference with 3 in-context demonstrations.

Method	Input							Avg
LLaVA-13B [93]	1-View Img.	0.0	0.0	0.0	0.0	0.0	0.0	0.0
LLaVA-13B [93]	4-View Img.	0.0	0.0	0.0	0.0	0.0	0.0	0.0
LLaVA-13B† [93]	1-View Img.	1.8	9.3	3.8	0.0	2.1	11.1	4.4
LLaVA-13B† [93]	4-View Img.	2.5	13.7	7.7	0.0	4.3	11.1	6.2
GPT-4V [115]	4-View Img.	0.0	0.0	0.0	0.0	0.0	0.0	0.0
GPT-4V‡ [115]	4-View Img.	0.1	1.6	0.0	0.0	0.0	0.0	0.3
<b>SHAPELLM-7B</b>	Point Cloud	<b>5.9</b>	<b>25.8</b>	<b>11.5</b>	<b>3.4</b>	<b>5.1</b>	<b>11.1</b>	<b>10.5</b>
<b>SHAPELLM-13B</b>	Point Cloud	<b>7.6</b>	<b>26.7</b>	<b>11.5</b>	<b>6.7</b>	<b>6.8</b>	<b>11.1</b>	<b>11.7</b>

explicitly encompass the estimation of the object pose and a more profound knowledge of 3D spatial relationships.

## 5.2. Is SHAPELLM grounded in physical worlds?

Tab. 5 compares SHAPELLM with image-only methods on 3D referring expression grounding (REG) of 6-DoF poses on GPartNet. The results show that: i) Image-only methods cannot perform zero-shot geometry-necessary 6-DoF pose estimation. ii) Compared to image-only methods with 2D to 6-DoF pose estimation fine-tuning or in-context prompting, SHAPELLM still performs significantly better. It demonstrates the necessity of geometry and the difficulty of the ill-posed 2D to 6-DoF pose estimation problem.

## 5.3. Can SHAPELLM generalize to unseen objects?

Fig. 7 shows the part understanding examples of unseen objects. While SHAPELLM’s 6-DoF pose estimation is trained on GPartNet, which primarily consists of *indoor articulated furniture*. It has demonstrated promising generalization potential of spatial understanding on the *open-world objects*, paving ways for scaling up spatial-awareness training.

## 6. Related Works

**Interaction-oriented 3D Understanding** Interaction with 3D objects typically involves concept-only interaction and physical-grounded interaction [15]. The former works focus on 3D perception and semantic parsing, such as 3D object recognition and scene perception [99, 126, 127, 157, 159]. By utilizing language for open-ended interaction in 3D, a number of works demonstrate successful 3D scene QA [109, 174], grounding [16], and captioning [17]. Recently, some works propose to utilize foundation models like LLMs or CLIP for open-ended 3D object recognition [94, 186, 198] and scene segmentation [122, 181]. Guo & Zhang *et al.* [54] utilizes ImageBind [49] and LLaMA-Adapter [187] to realize point cloud-based interactive QA. Following LLaVA, PointLLM [166] conducts supervised fine-tuning by constructing a visual instruction-following dataset. Other works focus on scene-level tasks utilizing comprehensive 2D features [68, 199] or 3D features dis-

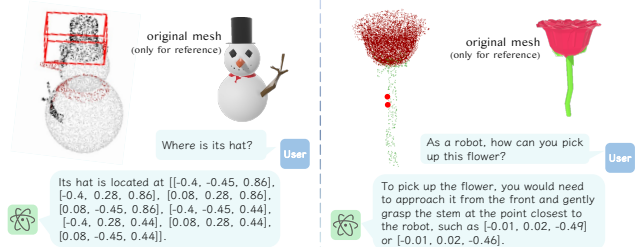


Figure 7. **Part understanding examples** of unseen objects.

tilled from 2D images into LLMs [62, 68, 199]. The second kind of interaction typically requires physical understanding in 3D, such as part understanding [47, 96, 105, 113], 6-DoF pose estimation [84, 97, 156, 160, 177], particularly useful for human-object interaction (HOI) and robotic manipulation [45–48, 50, 85, 95, 128, 139, 154, 167, 178] and complex robotic planning [14, 34, 39, 72, 90, 141]. In this work, we focus on both physical and conceptual interactions with 3D shapes for embodied understanding.

**Multimodal Large Language Models** Multimodal comprehension, which allows human interaction with textual and visual elements, has witnessed significant advancements, particularly in extending LLMs like LLaMA [21, 149, 150]. The early efforts predominantly revolved around integrating LLMs with various downstream systems by employing it as an agent [6, 57, 88, 140, 148, 155, 161, 171, 172]. Significant success has been demonstrated within this plugin-style framework. Due to the remarkable capabilities of LLMs, aligning the visual semantic space with language through parameter-efficient tuning [2, 64, 83, 173, 187, 197] and instruction tuning [24, 37, 93, 168] has emerged as the prevailing approach in current research. To further enhance interactive capabilities, some approaches have been developed towards visual-interactive multimodal comprehension by precisely referring to instruction tuning [19, 123, 189, 191]. Another family advances the developments of LLMs endowed with content creation beyond comprehension, notable efforts include DreamLLM [37], GILL [81], Emu [145, 147], SEED [44], NeXt-GPT [162], and Kosmos-G [118].

## 7. Conclusions

This paper presents SHAPELLM, a 3D multimodal LLM for embodied interaction, capable of generalizable recognition and embodied interaction comprehension. We first propose a novel 3D point cloud encoder, RECON++, by utilizing multi-view distillation and scaling up 3D representation learning, which serves as the foundation 3D representation encoder for SHAPELLM. Then, we perform 3D visual instruction tuning on constructed instruction-following data for general and embodied comprehension. We also established a 3D evaluation benchmark, *3D MM-Vet*, severing as assessing the 4-level capacity in embodied interaction scenarios, varying from basic perception to control statements generation.



## References

- [1] Panos Achlioptas, Olga Diamanti, Ioannis Mitliagkas, and Leonidas J. Guibas. Learning representations and generative models for 3d point clouds. In *Int. Conf. Mach. Learn. (ICML)*, 2018. 2
- [2] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katie Millican, Malcolm Reynolds, Roman Ring, Eliza Rutherford, Serkan Cabi, Tengda Han, Zhitao Gong, Sina Samangooei, Marianne Monteiro, Jacob Menick, Sebastian Borgeaud, Andrew Brock, Aida Nematzadeh, Sahand Sharifzadeh, Mikolaj Binkowski, Ricardo Barreira, Oriol Vinyals, Andrew Zisserman, and Karen Simonyan. Flamingo: a visual language model for few-shot learning. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. 8
- [3] 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 *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2016. 22
- [4] Yutong Bai, Xinyang Geng, Kartikeya Mangalam, Amir Bar, Alan L. Yuille, Trevor Darrell, Jitendra Malik, and Alexei A. Efros. Sequential modeling enables scalable learning for large vision models. *CoRR*, abs/2312.00785, 2023. 2
- [5] Satanjeev Banerjee and Alon Lavie. METEOR: an automatic metric for MT evaluation with improved correlation with human judgments. In *Proceedings of the Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization@ACL 2005, Ann Arbor, Michigan, USA, June 29, 2005*, 2005. 18
- [6] James Betker, Goh Gabriel, Li Jing, Tim Brooks, Jianfeng Wang, Linjie Li, Long Ouyang, Juntang Zhuang, Joyce Lee, Yufei Guo, Wesam Manassra, Prafulla Dhariwal, Casey Chu, Yunxin Jiao, and Aditya Ramesh. Improving image generation with better captions. 2023. 8
- [7] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri S. Chatterji, Annie S. Chen, Kathleen Creel, Jared Quincy Davis, Dorotya Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah D. Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshthe Khani, Omar Khattab, Pang Wei Koh, Mark S. Krass, Ranjay Krishna, Rohith Kudipudi, and et al. On the opportunities and risks of foundation models. *CoRR*, abs/2108.07258, 2021. 2
- [8] Bernhard E. Boser, Isabelle Guyon, and Vladimir Vapnik. A training algorithm for optimal margin classifiers. In *ACM Conf. Comput. Learn. Theory (COLT)*, pages 144–152. ACM, 1992. 20
- [9] Gary Bradski and Stephen Grossberg. Recognition of 3-d objects from multiple 2-d views by a self-organizing neural architecture. In *From Statistics to Neural Networks: Theory and Pattern Recognition Applications*, pages 349–375. Springer, 1994. 2
- [10] Alexander M. Bronstein, Michael M. Bronstein, Leonidas J. Guibas, and Maks Ovsjanikov. Shape google: Geometric words and expressions for invariant shape retrieval. *ACM Trans. Graph.*, 30(1):1:1–1:20, 2011. 2
- [11] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2020. 2
- [12] Nicolas Carion, Francisco Massa, Gabriel Synnaeve, Nicolas Usunier, Alexander Kirillov, and Sergey Zagoruyko. End-to-end object detection with transformers. In *Eur. Conf. Comput. Vis. (ECCV)*, 2020. 2, 5
- [13] Angel X. Chang, Thomas A. Funkhouser, Leonidas J. Guibas, Pat Hanrahan, Qi-Xing Huang, Zimo Li, Silvio Savarese, Manolis Savva, Shuran Song, Hao Su, Jianxiong Xiao, Li Yi, and Fisher Yu. Shapenet: An information-rich 3d model repository. *CoRR*, abs/1512.03012, 2015. 2, 5, 6, 7, 21
- [14] Matthew Chang, Théophile Gervet, Mukul Khanna, Sriram Yenamandra, Dhruv Shah, So Yeon Min, Kavita Shah, Chris Paxton, Saurabh Gupta, Dhruv Batra, Roozbeh Mottaghi, Jitendra Malik, and Devendra Singh Chaplot. GOAT: GO to any thing. *CoRR*, abs/2311.06430, 2023. 8
- [15] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brian Ichter, Danny Driess, Pete Florence, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. SpatialVLM: Endowing vision-language models with spatial reasoning capabilities. *CoRR*, abs/2401.12168, 2023. 8
- [16] Dave Zhenyu Chen, Angel X. Chang, and Matthias Nießner. Scanrefer: 3d object localization in RGB-D scans using natural language. In *Eur. Conf. Comput. Vis. (ECCV)*, 2020. 8
- [17] Dave Zhenyu Chen, Ali Gholami, Matthias Nießner, and Angel X. Chang. Scan2cap: Context-aware dense captioning in RGB-D scans. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2021. 8
- [18] Guangyan Chen, Meiling Wang, Yi Yang, Kai Yu, Li Yuan, and Yufeng Yue. Pointgpt: Auto-regressively generative pre-training from point clouds. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 5, 6, 20
- [19] Keqin Chen, Zhao Zhang, Weili Zeng, Richong Zhang, Feng Zhu, and Rui Zhao. Shikra: Unleashing multimodal LLM’s referential dialogue magic. *CoRR*, abs/2306.15195, 2023. 8
- [20] Xi Chen, Josip Djolonga, Piotr Padlewski, Basil Mustafa, Soravit Changpinyo, Jialin Wu, Carlos Riquelme Ruiz, Sebastian Goodman, Xiao Wang, Yi Tay, Siamak Shakeri,

- Mostafa Dehghani, Daniel Salz, Mario Lucic, Michael Tschannen, Arsha Nagrani, Hexiang Hu, Mandar Joshi, Bo Pang, Ceslee Montgomery, Paulina Pietrzyk, Marvin Ritter, A. J. Piergiovanni, Matthias Minderer, Filip Pavetic, Austin Waters, Gang Li, Ibrahim Alabdulmohsin, Lucas Beyer, Julien Amelot, Kenton Lee, Andreas Peter Steiner, Yang Li, Daniel Keysers, Anurag Arnab, Yuanzhong Xu, Keran Rong, Alexander Kolesnikov, Mojtaba Seyedhosseini, Anelia Angelova, Xiaohua Zhai, Neil Houlsby, and Radu Soricut. Pali-x: On scaling up a multilingual vision and language model. In *Int. Conf. Learn. Represent. (ICLR)*, 2023. [2](#)
- [21] Wei-Lin Chiang, Zhuohan Li, Zi Lin, Ying Sheng, Zhanghao Wu, Hao Zhang, Lianmin Zheng, Siyuan Zhuang, Yonghao Zhuang, Joseph E. Gonzalez, Ion Stoica, and Eric P. Xing. Vicuna: An open-source chatbot impressing gpt-4 with 90%\* chatgpt quality, 2023. [8](#), [21](#)
- [22] Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu, Xi Zhang, Tomas F Yago Vicente, Thomas Dideriksen, Himanshu Arora, et al. Abo: Dataset and benchmarks for real-world 3d object understanding. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. [6](#), [21](#)
- [23] Angela Dai, Angel X Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017. [6](#), [22](#)
- [24] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven C. H. Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. [3](#), [8](#), [19](#)
- [25] Wenliang Dai, Zihan Liu, Ziwei Ji, Dan Su, and Pascale Fung. Plausible may not be faithful: Probing object hallucination in vision-language pre-training. In *Proceedings of the 17th Conference of the European Chapter of the Association for Computational Linguistics, EACL 2023, Dubrovnik, Croatia, May 2-6, 2023*, 2023. [19](#)
- [26] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. [22](#)
- [27] Abhishek Das, Satwik Kottur, Khushi Gupta, Avi Singh, Deshraj Yadav, Stefan Lee, José M. F. Moura, Devi Parikh, and Dhruv Batra. Visual dialog. *IEEE Trans. Pattern Anal. Mach. Intell. (TPAMI)*, 41(5):1242–1256, 2019. [2](#)
- [28] Joe Davison, Joshua Feldman, and Alexander M. Rush. Commonsense knowledge mining from pretrained models. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 2019. [2](#)
- [29] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs, Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse: A universe of annotated 3d objects. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [2](#), [4](#), [5](#), [6](#), [7](#), [21](#)
- [30] Jiajun Deng, Shaoshuai Shi, Peiwei Li, Wengang Zhou, Yanyong Zhang, and Houqiang Li. Voxel r-cnn: Towards high performance voxel-based 3d object detection. In *AAAI Conf. Artif. Intell. (AAAI)*, 2021. [22](#)
- [31] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: pre-training of deep bidirectional transformers for language understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT 2019, Minneapolis, MN, USA, June 2-7, 2019, Volume 1 (Long and Short Papers)*, pages 4171–4186. Association for Computational Linguistics, 2019. [22](#)
- [32] Runyu Ding, Jihan Yang, Chuhui Xue, Wenqing Zhang, Song Bai, and Xiaojuan Qi. Lowis3d: Language-driven open-world instance-level 3d scene understanding. *CoRR*, abs/2308.00353, 2023. [22](#)
- [33] Runyu Ding, Jihan Yang, Chuhui Xue, Wenqing Zhang, Song Bai, and Xiaojuan Qi. PLA: language-driven open-vocabulary 3d scene understanding. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [22](#)
- [34] Yan Ding, Xiaohan Zhang, Chris Paxton, and Shiqi Zhang. Task and motion planning with large language models for object rearrangement. *CoRR*, abs/2303.06247, 2023. [8](#)
- [35] Runpei Dong, Zhanhong Tan, Mengdi Wu, Linfeng Zhang, and Kaisheng Ma. Finding the task-optimal low-bit sub-distribution in deep neural networks. In *Int. Conf. Mach. Learn. (ICML)*, 2022. [23](#)
- [36] Runpei Dong, Zekun Qi, Linfeng Zhang, Junbo Zhang, Jianjian Sun, Zheng Ge, Li Yi, and Kaisheng Ma. Autoencoders as cross-modal teachers: Can pretrained 2d image transformers help 3d representation learning? In *Int. Conf. Learn. Represent. (ICLR)*, 2023. [2](#), [3](#), [6](#), [17](#), [20](#), [22](#)
- [37] Runpei Dong, Chunrui Han, Yuang Peng, Zekun Qi, Zheng Ge, Jinrong Yang, Liang Zhao, Jianjian Sun, Hongyu Zhou, Haoran Wei, Xiangwen Kong, Xiangyu Zhang, Kaisheng Ma, and Li Yi. DreamLLM: Synergistic multimodal comprehension and creation. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. [2](#), [3](#), [7](#), [8](#), [17](#), [21](#)
- [38] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In *Int. Conf. Learn. Represent. (ICLR)*, 2021. [21](#)
- [39] Danny Driess, Fei Xia, Mehdi S. M. Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, Wenlong Huang, Yevgen Chebotar, Pierre Sermanet, Daniel Duckworth, Sergey Levine, Vincent Vanhoucke, Karol Hausman, Marc Toussaint, Klaus Greff, Andy Zeng, Igor Mordatch, and Pete Florence. Palm-e: An embodied multimodal language model. In *Int. Conf. Mach. Learn. (ICML)*, 2023. [2](#), [6](#), [8](#)
- [40] Nico Engel, Vasileios Belagiannis, and Klaus Dietmayer. Point transformer. *IEEE Access*, 9:134826–134840, 2021. [22](#)

- [41] Haoqiang Fan, Hao Su, and Leonidas J. Guibas. A point set generation network for 3d object reconstruction from a single image. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017. [22](#)
- [42] Huan Fu, Rongfei Jia, Lin Gao, Mingming Gong, Binqiang Zhao, Steve Maybank, and Dacheng Tao. 3d-future: 3d furniture shape with texture. *International Journal of Computer Vision*, 129:3313–3337, 2021. [6](#), [21](#)
- [43] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing, EMNLP 2021, Virtual Event / Punta Cana, Dominican Republic, 7-11 November, 2021*, 2021. [18](#)
- [44] Yuying Ge, Yixiao Ge, Ziyun Zeng, Xintao Wang, and Ying Shan. Planting a SEED of vision in large language model. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. [8](#)
- [45] Haoran Geng, Ziming Li, Yiran Geng, Jiayi Chen, Hao Dong, and He Wang. Partmanip: Learning cross-category generalizable part manipulation policy from point cloud observations. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [8](#)
- [46] Haoran Geng, Songlin Wei, Congyue Deng, Bokui Shen, He Wang, and Leonidas Guibas. Sage: Bridging semantic and actionable parts for generalizable articulated-object manipulation under language instructions. *CoRR*, abs/2312.01307, 2023.
- [47] Haoran Geng, Helin Xu, Chengyang Zhao, Chao Xu, Li Yi, Siyuan Huang, and He Wang. Gapartnet: Cross-category domain-generalizable object perception and manipulation via generalizable and actionable parts. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [2](#), [4](#), [8](#), [23](#)
- [48] Yiran Geng, Boshi An, Haoran Geng, Yuanpei Chen, Yaodong Yang, and Hao Dong. Rlafford: End-to-end affordance learning for robotic manipulation. In *IEEE Int. Conf. Robot. Autom. (ICRA)*, 2023. [8](#)
- [49] Rohit Girdhar, Alaaeldin El-Nouby, Zhuang Liu, Mannat Singh, Kalyan Vasudev Alwala, Armand Joulin, and Ishan Misra. Imagebind: One embedding space to bind them all. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15180–15190, 2023. [8](#)
- [50] Ran Gong, Jiangyong Huang, Yizhou Zhao, Haoran Geng, Xiaofeng Gao, Qingyang Wu, Wensi Ai, Ziheng Zhou, Demetri Terzopoulos, Song-Chun Zhu, Baoxiong Jia, and Siyuan Huang. ARNOLD: A benchmark for language-grounded task learning with continuous states in realistic 3d scenes. *CoRR*, abs/2304.04321, 2023. [8](#)
- [51] Priya Goyal, Dhruv Mahajan, Abhinav Gupta, and Ishan Misra. Scaling and benchmarking self-supervised visual representation learning. In *Int. Conf. Comput. Vis. (ICCV)*, pages 6390–6399. IEEE, 2019. [20](#)
- [52] Helmut Grabner, Juergen Gall, and Luc Van Gool. What makes a chair a chair? In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2011. [2](#)
- [53] Anisha Gunjal, Jihan Yin, and Erhan Bas. Detecting and preventing hallucinations in large vision language models. *CoRR*, abs/2308.06394, 2023. [19](#)
- [54] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, and Pheng-Ann Heng. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *CoRR*, abs/2309.00615, 2023. [2](#), [7](#), [8](#), [19](#)
- [55] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2019. [4](#), [7](#)
- [56] Saurabh Gupta, Ross B. Girshick, Pablo Andrés Arbeláez, and Jitendra Malik. Learning rich features from RGB-D images for object detection and segmentation. In *Eur. Conf. Comput. Vis. (ECCV)*, 2014. [18](#)
- [57] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [8](#)
- [58] Abdullah Hamdi, Silvio Giancola, and Bernard Ghanem. MVTN: multi-view transformation network for 3d shape recognition. In *Int. Conf. Comput. Vis. (ICCV)*, pages 1–11. IEEE, 2021. [2](#), [22](#)
- [59] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified view of parameter-efficient transfer learning. In *Int. Conf. Learn. Represent. (ICLR)*, 2021. [18](#)
- [60] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked autoencoders are scalable vision learners. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. [22](#)
- [61] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *CoRR*, abs/1606.08415, 2016. [21](#)
- [62] Yining Hong, Haoyu Zhen, Peihao Chen, Shuhong Zheng, Yilun Du, Zhenfang Chen, and Chuang Gan. 3d-llm: Injecting the 3d world into large language models. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. [3](#), [8](#), [19](#)
- [63] Ji Hou, Saining Xie, Benjamin Graham, Angela Dai, and Matthias Nießner. Pri3d: Can 3d priors help 2d representation learning? In *Int. Conf. Comput. Vis. (ICCV)*, pages 5673–5682. IEEE, 2021. [2](#)
- [64] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen. Lora: Low-rank adaptation of large language models. In *Int. Conf. Learn. Represent. (ICLR)*, 2022. [8](#)
- [65] Ruizhen Hu, Chenyang Zhu, Oliver van Kaick, Ligang Liu, Ariel Shamir, and Hao Zhang. Interaction context (ICON): towards a geometric functionality descriptor. *ACM Trans. Graph.*, 34(4):83:1–83:12, 2015. [2](#)
- [66] Ruizhen Hu, Oliver van Kaick, Bojian Wu, Hui Huang, Ariel Shamir, and Hao Zhang. Learning how objects function via co-analysis of interactions. *ACM Trans. Graph.*, 35(4):47:1–47:13, 2016.
- [67] Ruizhen Hu, Wenchao Li, Oliver van Kaick, Ariel Shamir, Hao Zhang, and Hui Huang. Learning to predict part mobility from a single static snapshot. *ACM Trans. Graph.*, 36(6):227:1–227:13, 2017. [2](#)
- [68] Jiangyong Huang, Silong Yong, Xiaojian Ma, Xiongkun Linghu, Puhao Li, Yan Wang, Qing Li, Song-Chun Zhu,

- Baoxiong Jia, and Siyuan Huang. An embodied generalist agent in 3d world. *CoRR*, abs/2311.12871, 2023. 8
- [69] Tianyu Huang, Bowen Dong, Yunhan Yang, Xiaoshui Huang, Rynson W. H. Lau, Wanli Ouyang, and Wangmeng Zuo. Clip2point: Transfer CLIP to point cloud classification with image-depth pre-training. *CoRR*, abs/2210.01055, 2022. 6
- [70] Wenlong Huang, Igor Mordatch, and Deepak Pathak. One policy to control them all: Shared modular policies for agent-agnostic control. In *Int. Conf. Mach. Learn. (ICML)*, 2020. 2, 6
- [71] Wenlong Huang, Fei Xia, Ted Xiao, Harris Chan, Jacky Liang, Pete Florence, Andy Zeng, Jonathan Tompson, Igor Mordatch, Yevgen Chebotar, Pierre Sermanet, Tomas Jackson, Noah Brown, Linda Luu, Sergey Levine, Karol Hausman, and Brian Ichter. Inner monologue: Embodied reasoning through planning with language models. In *Annu. Conf. Robot. Learn. (CoRL)*, 2022. 2
- [72] Wenlong Huang, Chen Wang, Ruohan Zhang, Yunzhu Li, Jiajun Wu, and Li Fei-Fei. Voxposer: Composable 3d value maps for robotic manipulation with language models. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 8, 23
- [73] Brian Ichter, Anthony Brohan, Yevgen Chebotar, Chelsea Finn, Karol Hausman, Alexander Herzog, Daniel Ho, Julian Ibarz, Alex Irpan, Eric Jang, Ryan Julian, Dmitry Kalashnikov, Sergey Levine, Yao Lu, Carolina Parada, Kanishka Rao, Pierre Sermanet, Alexander Toshev, Vincent Vanhoucke, Fei Xia, Ted Xiao, Peng Xu, Mengyuan Yan, Noah Brown, Michael Ahn, Omar Cortes, Nicolas Sievers, Clayton Tan, Sichun Xu, Diego Reyes, Jarek Rettinghouse, Jor-nell Quiambao, Peter Pastor, Linda Luu, Kuang-Huei Lee, Yuheng Kuang, Sally Jesmonth, Nikhil J. Joshi, Kyle Jeffrey, Rosario Jauregui Ruano, Jasmine Hsu, Keerthana Gopalakrishnan, Byron David, Andy Zeng, and Chuyuan Kelly Fu. Do as I can, not as I say: Grounding language in robotic affordances. In *Annu. Conf. Robot. Learn. (CoRL)*, 2022. 2
- [74] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Han-naneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Open-clip, 2021. 6
- [75] Ajay Jaiswal, Zhe Gan, Xianzhi Du, Bowen Zhang, Zhangyang Wang, and Yinfei Yang. Compressing llms: The truth is rarely pure and never simple. *CoRR*, abs/2310.01382, 2023. 23
- [76] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge J. Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 3, 18
- [77] Yunfan Jiang, Agrim Gupta, Zichen Zhang, Guanzhi Wang, Yongqiang Dou, Yanjun Chen, Li Fei-Fei, Anima Anandkumar, Yuke Zhu, and Linxi Fan. VIMA: general robot manipulation with multimodal prompts. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 23
- [78] Zhengbao Jiang, Frank F. Xu, Jun Araki, and Graham Neubig. How can we know what language models know. *Trans. Assoc. Comput. Linguistics*, 8:423–438, 2020. 2
- [79] Takeo Kanade and Masatoshi Okutomi. A stereo matching algorithm with an adaptive window: Theory and experiment. *IEEE Trans. Pattern Anal. Mach. Intell.*, 16(9):920–932, 1994. 2
- [80] Vladimir G. Kim, Siddhartha Chaudhuri, Leonidas J. Guibas, and Thomas A. Funkhouser. Shape2pose: human-centric shape analysis. *ACM Trans. Graph.*, 33(4):120:1–120:12, 2014. 2
- [81] Jing Yu Koh, Daniel Fried, and Ruslan Salakhutdinov. Generating images with multimodal language models. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 8
- [82] Harold W Kuhn. The hungarian method for the assignment problem. *Naval research logistics quarterly*, 2(1-2):83–97, 1955. 2
- [83] Junnan Li, Dongxu Li, Silvio Savarese, and Steven C. H. Hoi. BLIP-2: bootstrapping language-image pre-training with frozen image encoders and large language models. In *Int. Conf. Mach. Learn. (ICML)*, 2023. 8
- [84] Xiaolong Li, He Wang, Li Yi, Leonidas J. Guibas, A. Lynn Abbott, and Shuran Song. Category-level articulated object pose estimation. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2020. 2, 8
- [85] Xiaoqi Li, Mingxu Zhang, Yiran Geng, Haoran Geng, Yuxing Long, Yan Shen, Renrui Zhang, Jiaming Liu, and Hao Dong. Manipllm: Embodied multimodal large language model for object-centric robotic manipulation, 2023. 8
- [86] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, 2021. 3, 17
- [87] Yifan Li, Yifan Du, Kun Zhou, Jinpeng Wang, Wayne Xin Zhao, and Ji-Rong Wen. Evaluating object hallucination in large vision-language models. *CoRR*, abs/2305.10355, 2023. 2
- [88] Yaobo Liang, Chenfei Wu, Ting Song, Wenshan Wu, Yan Xia, Yu Liu, Yang Ou, Shuai Lu, Lei Ji, Shaoguang Mao, Yun Wang, Linjun Shou, Ming Gong, and Nan Duan. Taskmatrix.ai: Completing tasks by connecting foundation models with millions of apis. *CoRR*, abs/2303.16434, 2023. 8
- [89] Chin-Yew Lin. Rouge: A package for automatic evaluation of summaries. In *Proc. Workshop on Text Summarization Branches Out, Post-Conference Workshop of ACL 2004*, 2004. 18
- [90] Kevin Lin, Christopher Agia, Toki Migimatsu, Marco Pavone, and Jeannette Bohg. Text2motion: From natural language instructions to feasible plans. *CoRR*, abs/2303.12153, 2023. 8
- [91] Fuxiao Liu, Kevin Lin, Linjie Li, Jianfeng Wang, Yaser Yacoob, and Lijuan Wang. Aligning large multi-modal model with robust instruction tuning. *CoRR*, abs/2306.14565, 2023. 19
- [92] Haotian Liu, Chunyuan Li, Yuheng Li, and Yong Jae Lee. Improved baselines with visual instruction tuning. *CoRR*, abs/2310.03744, 2023. 2

- [93] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 2, 3, 7, 8, 19, 21
- [94] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yin hao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. Openshape: Scaling up 3d shape representation towards open-world understanding. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 5, 6, 7, 8, 21
- [95] Xueyi Liu and Li Yi. GeneOH diffusion: Towards generalizable hand-object interaction denoising via denoising diffusion. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. 4, 8
- [96] Xueyi Liu, Bin Wang, He Wang, and Li Yi. Few-shot physically-aware articulated mesh generation via hierarchical deformation. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. 8
- [97] Xueyi Liu, Ji Zhang, Ruizhen Hu, Haibin Huang, He Wang, and Li Yi. Self-supervised category-level articulated object pose estimation with part-level SE(3) equivariance. In *Int. Conf. Learn. Represent. (ICLR)*, 2023. 8
- [98] Yongcheng Liu, Bin Fan, Shiming Xiang, and Chunhong Pan. Relation-shape convolutional neural network for point cloud analysis. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2019. 6
- [99] Yunze Liu, Junyu Chen, Zekai Zhang, Jingwei Huang, and Li Yi. Leaf: Learning frames for 4d point cloud sequence understanding. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. 8
- [100] Yuan Liu, Haodong Duan, Yuanhan Zhang, Bo Li, Songyang Zhang, Wangbo Zhao, Yike Yuan, Jiaqi Wang, Conghui He, Ziwei Liu, Kai Chen, and Dahua Lin. Mmbench: Is your multi-modal model an all-around player? *CoRR*, abs/2307.06281, 2023. 6
- [101] Yuan Liu, Cheng Lin, Zijiao Zeng, Xiaoxiao Long, Lingjie Liu, Taku Komura, and Wenping Wang. Syncdreamer: Generating multiview-consistent images from a single-view image. *CoRR*, abs/2309.03453, 2023. 2
- [102] Ze Liu, Zheng Zhang, Yue Cao, Han Hu, and Xin Tong. Group-free 3d object detection via transformers. In *Int. Conf. Comput. Vis. (ICCV)*, 2021. 22
- [103] Ilya Loshchilov and Frank Hutter. SGDR: stochastic gradient descent with warm restarts. In *Int. Conf. Learn. Represent. (ICLR)*, 2017. 22
- [104] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *Int. Conf. Learn. Represent. (ICLR)*, 2019. 22
- [105] Cewu Lu, Hao Su, Yonglu Li, Yongyi Lu, Li Yi, Chi-Keung Tang, and Leonidas J. Guibas. Beyond holistic object recognition: Enriching image understanding with part states. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2018. 8
- [106] Shiyang Lu, Haonan Chang, Eric Pu Jing, Abdeslam Boularias, and Kostas E. Bekris. OVIR-3D: open-vocabulary 3d instance retrieval without training on 3d data. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 22
- [107] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 4, 5
- [108] Xu Ma, Can Qin, Haoxuan You, Haoxi Ran, and Yun Fu. Rethinking network design and local geometry in point cloud: A simple residual MLP framework. In *Int. Conf. Learn. Represent. (ICLR)*. OpenReview.net, 2022. 6
- [109] Xiaojian Ma, Silong Yong, Zilong Zheng, Qing Li, Yitao Liang, Song-Chun Zhu, and Siyuan Huang. SQA3D: situated question answering in 3d scenes. In *Int. Conf. Learn. Represent. (ICLR)*, 2023. 8
- [110] Haley MacLeod, Cynthia L. Bennett, Meredith Ringel Morris, and Edward Cutrell. Understanding blind people’s experiences with computer-generated captions of social media images. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, page 5988–5999, New York, NY, USA, 2017. Association for Computing Machinery. 19
- [111] Jiageng Mao, Yujing Xue, Minzhe Niu, Haoyue Bai, Jiashi Feng, Xiaodan Liang, Hang Xu, and Chunjing Xu. Voxel transformer for 3d object detection. In *Int. Conf. Comput. Vis. (ICCV)*, 2021. 22
- [112] Daniel Maturana and Sebastian A. Scherer. Voxnet: A 3d convolutional neural network for real-time object recognition. In *IEEE/RSJ Int. Conf. Intell. Robot. and Syst. (IROS)*, pages 922–928. IEEE, 2015. 22
- [113] Kaichun Mo, Shilin Zhu, Angel X Chang, Li Yi, Subarna Tripathi, Leonidas J Guibas, and Hao Su. Partnet: A large-scale benchmark for fine-grained and hierarchical part-level 3d object understanding. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2019. 8
- [114] OpenAI. Introducing chatgpt. 2022. 21
- [115] OpenAI. Gpt-4v(ision) system card, 2023. 2, 4, 7, 8
- [116] OpenAI. GPT-4 technical report. *CoRR*, abs/2303.08774, 2023. 2, 19, 21
- [117] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll L. Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, John Schulman, Jacob Hilton, Fraser Kelton, Luke Miller, Maddie Simens, Amanda Askell, Peter Welinder, Paul F. Christiano, Jan Leike, and Ryan Lowe. Training language models to follow instructions with human feedback. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. 2
- [118] Xichen Pan, Li Dong, Shaohan Huang, Zhiliang Peng, Wenhua Chen, and Furu Wei. Kosmos-g: Generating images in context with multimodal large language models. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. 8
- [119] Yatian Pang, Wenxiao Wang, Francis E. H. Tay, Wei Liu, Yonghong Tian, and Li Yuan. Masked autoencoders for point cloud self-supervised learning. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 6, 20, 22
- [120] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. 2002. 18
- [121] Baolin Peng, Chunyuan Li, Pengcheng He, Michel Galley, and Jianfeng Gao. Instruction tuning with GPT-4. *CoRR*, abs/2304.03277, 2023. 2
- [122] Songyou Peng, Kyle Genova, Chiyu “Max” Jiang, Andrea Tagliasacchi, Marc Pollefeys, and Thomas A. Funkhouser. Openscene: 3d scene understanding with open vocabularies.

- In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. 8, 22
- [123] Zhiliang Peng, Wenhui Wang, Li Dong, Yaru Hao, Shaohan Huang, Shuming Ma, and Furu Wei. Kosmos-2: Grounding multimodal large language models to the world. *CoRR*, abs/2306.14824, 2023. 8
- [124] Fabio Petroni, Tim Rocktäschel, Sebastian Riedel, Patrick S. H. Lewis, Anton Bakhtin, Yuxiang Wu, and Alexander H. Miller. Language models as knowledge bases? In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 2019. 2
- [125] Sören Pirk, Vojtech Krs, Kai-Mo Hu, Suren Deepak Rajasekaran, Hao Kang, Yusuke Yoshiyasu, Bedrich Benes, and Leonidas J. Guibas. Understanding and exploiting object interaction landscapes. *ACM Trans. Graph.*, 36(3):31:1–31:14, 2017. 2
- [126] Charles Ruizhongtai Qi, Hao Su, Kaichun Mo, and Leonidas J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 77–85, 2017. 2, 6, 8, 17, 22
- [127] Charles Ruizhongtai Qi, Li Yi, Hao Su, and Leonidas J. Guibas. Pointnet++: Deep hierarchical feature learning on point sets in a metric space. In *Adv. Neural Inform. Process. Syst. (NIPS)*, pages 5099–5108, 2017. 2, 6, 8, 17, 22
- [128] Haozhi Qi, Ashish Kumar, Roberto Calandra, Yi Ma, and Jitendra Malik. In-hand object rotation via rapid motor adaptation. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 8
- [129] Zekun Qi, Runpei Dong, Guofan Fan, Zheng Ge, Xiangyu Zhang, Kaisheng Ma, and Li Yi. Contrast with reconstruct: Contrastive 3d representation learning guided by generative pretraining. In *Int. Conf. Mach. Learn. (ICML)*, 2023. 2, 3, 5, 6, 17, 20, 21, 22
- [130] Zekun Qi, Muzhou Yu, Runpei Dong, and Kaisheng Ma. VPP: efficient conditional 3d generation via voxel-point progressive representation. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 6, 17, 20, 22
- [131] Guocheng Qian, Yuchen Li, Houwen Peng, Jinjie Mai, Hasan Abed Al Kader Hammoud, Mohamed Elhoseiny, and Bernard Ghanem. Pointnext: Revisiting pointnet++ with improved training and scaling strategies. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. 6, 22
- [132] Alec Radford, Karthik Narasimhan, Tim Salimans, and Ilya Sutskever. Improving language understanding by generative pre-training. 2018. 2
- [133] Alec Radford, Jeffrey Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019. 2
- [134] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning transferable visual models from natural language supervision. In *Int. Conf. Mach. Learn. (ICML)*, pages 8748–8763. PMLR, 2021. 2, 7, 22
- [135] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP 2019, Hong Kong, China, November 3-7, 2019*, 2019. 18
- [136] Jiawei Ren, Liang Pan, and Ziwei Liu. Benchmarking and analyzing point cloud classification under corruptions. In *Int. Conf. Mach. Learn. (ICML)*, 2022. 7
- [137] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, 2018. 19
- [138] Anna Rohrbach, Lisa Anne Hendricks, Kaylee Burns, Trevor Darrell, and Kate Saenko. Object hallucination in image captioning. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing, Brussels, Belgium, October 31 - November 4, 2018*, 2018. 2
- [139] William Shen, Ge Yang, Alan Yu, Jansen Wong, Leslie Pack Kaelbling, and Phillip Isola. Distilled feature fields enable few-shot language-guided manipulation. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 8
- [140] Yongliang Shen, Kaitao Song, Xu Tan, Dongsheng Li, Weiming Lu, and Yueting Zhuang. Hugginggpt: Solving AI tasks with chatgpt and its friends in huggingface. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 8
- [141] Haochen Shi, Huazhe Xu, Samuel Clarke, Yunzhu Li, and Jiajun Wu. Robocook: Long-horizon elasto-plastic object manipulation with diverse tools. In *Annu. Conf. Robot. Learn. (CoRL)*, 2023. 8
- [142] Shutterstock. Turbosquid. <https://www.turbosquid.com/>. 6
- [143] Hang Su, Subhansu Maji, Evangelos Kalogerakis, and Erik G. Learned-Miller. Multi-view convolutional neural networks for 3d shape recognition. In *Int. Conf. Comput. Vis. (ICCV)*, 2015. 2, 22
- [144] Jiachen Sun, Qingzhao Zhang, Bhavya Kailkhura, Zhiding Yu, Chaowei Xiao, and Z Morley Mao. Modelnet40-c: A robustness benchmark for 3d point cloud recognition under corruption. In *ICLR 2022 Workshop on Socially Responsible Machine Learning*. 7
- [145] Quan Sun, Yufeng Cui, Xiaosong Zhang, Fan Zhang, Qiyang Yu, Zhengxiong Luo, Yueze Wang, Yongming Rao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative multimodal models are in-context learners. *CoRR*, abs/2312.13286, 2023. 8
- [146] Quan Sun, Yuxin Fang, Ledell Wu, Xinlong Wang, and Yue Cao. EVA-CLIP: improved training techniques for CLIP at scale. *CoRR*, abs/2303.15389, 2023. 6
- [147] Quan Sun, Qiyang Yu, Yufeng Cui, Fan Zhang, Xiaosong Zhang, Yueze Wang, Hongcheng Gao, Jingjing Liu, Tiejun Huang, and Xinlong Wang. Generative pretraining in multimodality. *CoRR*, abs/2307.05222, 2023. 8
- [148] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *CoRR*, abs/2303.08128, 2023. 8
- [149] Rohan Taori, Ishaan Gulrajani, Tianyi Zhang, Yann Dubois, Xuechen Li, Carlos Guestrin, Percy Liang, and Tatsunori B.

- Hashimoto. Stanford alpaca: An instruction-following llama model. [https://github.com/tatsu-lab/stanford\\_alpaca](https://github.com/tatsu-lab/stanford_alpaca), 2023. 8
- [150] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, Aurélien Rodriguez, Armand Joulin, Edouard Grave, and Guillaume Lample. Llama: Open and efficient foundation language models. *CoRR*, abs/2302.13971, 2023. 2, 3, 8, 21
- [151] Mikaela Angelina Uy, Quang-Hieu Pham, Binh-Son Hua, Thanh Nguyen, and Sai-Kit Yeung. Revisiting point cloud classification: A new benchmark dataset and classification model on real-world data. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 1588–1597, 2019. 6, 7, 21, 22
- [152] Vladimir Vapnik. *Statistical learning theory*. Wiley, 1998. 20
- [153] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need. In *Adv. Neural Inform. Process. Syst. (NIPS)*, pages 5998–6008, 2017. 6, 20, 22
- [154] Weikang Wan, Haoran Geng, Yun Liu, Zikang Shan, Yaodong Yang, Li Yi, and He Wang. Unidexgrasp++: Improving dexterous grasping policy learning via geometry-aware curriculum and iterative generalist-specialist learning. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. 2, 4, 8
- [155] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlikar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *CoRR*, abs/2305.16291, 2023. 8
- [156] He Wang, Srinath Sridhar, Jingwei Huang, Julien Valentin, Shuran Song, and Leonidas J. Guibas. Normalized object coordinate space for category-level 6d object pose and size estimation. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2019. 2, 8
- [157] Yue Wang, Yongbin Sun, Ziwei Liu, Sanjay E. Sarma, Michael M. Bronstein, and Justin M. Solomon. Dynamic graph CNN for learning on point clouds. *ACM Trans. Graph.*, 38(5):146:1–146:12, 2019. 6, 8
- [158] Ziyi Wang, Xumin Yu, Yongming Rao, Jie Zhou, and Jiwen Lu. Take-a-photo: 3d-to-2d generative pre-training of point cloud models. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. 6
- [159] Hao Wen, Yunze Liu, Jingwei Huang, Bo Duan, and Li Yi. Point primitive transformer for long-term 4d point cloud video understanding. In *Eur. Conf. Comput. Vis. (ECCV)*, 2022. 8
- [160] Yijia Weng, He Wang, Qiang Zhou, Yuzhe Qin, Yueqi Duan, Qingnan Fan, Baoquan Chen, Hao Su, and Leonidas J. Guibas. CAPTRA: category-level pose tracking for rigid and articulated objects from point clouds. In *Int. Conf. Comput. Vis. (ICCV)*, 2021. 2, 8
- [161] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *CoRR*, abs/2303.04671, 2023. 8
- [162] Shengqiong Wu, Hao Fei, Leigang Qu, Wei Ji, and Tat-Seng Chua. Next-gpt: Any-to-any multimodal LLM. *CoRR*, abs/2309.05519, 2023. 8
- [163] Tong Wu, Guandao Yang, Zhibing Li, Kai Zhang, Ziwei Liu, Leonidas J. Guibas, Dahua Lin, and Gordon Wetzstein. Gpt-4v(ision) is a human-aligned evaluator for text-to-3d generation. *CoRR*, abs/2401.04092, 2024. 19
- [164] 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/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, pages 1912–1920, 2015. 6, 7, 20, 22
- [165] 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 *Eur. Conf. Comput. Vis. (ECCV)*, pages 574–591. Springer, 2020. 22
- [166] Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. Pointllm: Empowering large language models to understand point clouds. *CoRR*, abs/2308.16911, 2023. 3, 7, 8, 18, 19, 21
- [167] Yinchen Xu, Weikang Wan, Jialiang Zhang, Haoran Liu, Zikang Shan, Hao Shen, Ruicheng Wang, Haoran Geng, Yijia Weng, Jiayi Chen, Tengyu Liu, Li Yi, and He Wang. Unidexgrasp: Universal robotic dexterous grasping via learning diverse proposal generation and goal-conditioned policy. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. 2, 4, 8
- [168] Zhiyang Xu, Ying Shen, and Lifu Huang. Multiinstruct: Improving multi-modal zero-shot learning via instruction tuning. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (ACL) (Volume 1: Long Papers)*, 2023. 8
- [169] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. ULIP: learning unified representation of language, image and point cloud for 3d understanding. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. 2, 6
- [170] Le Xue, Ning Yu, Shu Zhang, Junnan Li, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. ULIP-2: towards scalable multimodal pre-training for 3d understanding. *CoRR*, abs/2305.08275, 2023. 5, 6, 18
- [171] Rui Yang, Lin Song, Yanwei Li, Sijie Zhao, Yixiao Ge, Xiu Li, and Ying Shan. Gpt4tools: Teaching large language model to use tools via self-instruction. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2023. 8
- [172] Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Ehsan Azarnasab, Faisal Ahmed, Zicheng Liu, Ce Liu, Michael Zeng, and Lijuan Wang. MM-REACT: prompting chatgpt for multimodal reasoning and action. *CoRR*, abs/2303.11381, 2023. 8
- [173] Qinghao Ye, Haiyang Xu, Guohai Xu, Jiabo Ye, Ming Yan, Yiyang Zhou, Junyang Wang, Anwen Hu, Pengcheng Shi, Yaya Shi, Chenliang Li, Yuanhong Xu, Hehong Chen, Junfeng Tian, Qian Qi, Ji Zhang, and Fei Huang. mplug-owl: Modularization empowers large language models with multimodality. *CoRR*, abs/2304.14178, 2023. 8
- [174] Shuquan Ye, Dongdong Chen, Songfang Han, and Jing Liao. 3d question answering. *IEEE Transactions on Visualization and Computer Graphics*, 2022. 8

- [175] Li Yi, Vladimir G Kim, Duygu Ceylan, I-Chao Shen, Mengyan Yan, Hao Su, Cewu Lu, Qixing Huang, Alla Sheffer, and Leonidas Guibas. A scalable active framework for region annotation in 3d shape collections. *ACM Trans. Graph.*, 35(6):1–12, 2016. [22](#)
- [176] Li Yi, Hao Su, Xingwen Guo, and Leonidas J. Guibas. Syncspecnn: Synchronized spectral CNN for 3d shape segmentation. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2017. [22](#)
- [177] Li Yi, Haibin Huang, Difan Liu, Evangelos Kalogerakis, Hao Su, and Leonidas J. Guibas. Deep part induction from articulated object pairs. *ACM Trans. Graph.*, 37(6):209, 2018. [2](#), [8](#)
- [178] Yang You, Bokui Shen, Congyue Deng, Haoran Geng, He Wang, and Leonidas J. Guibas. Make a donut: Language-guided hierarchical emd-space planning for zero-shot deformable object manipulation. *CoRR*, abs/2311.02787, 2023. [8](#)
- [179] Weihao Yu, Zhengyuan Yang, Linjie Li, Jianfeng Wang, Kevin Lin, Zicheng Liu, Xinchao Wang, and Lijuan Wang. Mm-vet: Evaluating large multimodal models for integrated capabilities. *CoRR*, abs/2308.02490, 2023. [2](#), [5](#)
- [180] Xumin Yu, Lulu Tang, Yongming Rao, Tiejun Huang, Jie Zhou, and Jiwen Lu. Point-bert: Pre-training 3d point cloud transformers with masked point modeling. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. [6](#), [17](#), [20](#), [22](#)
- [181] Junbo Zhang, Runpei Dong, and Kaisheng Ma. CLIP-FO3D: learning free open-world 3d scene representations from 2d dense CLIP. In *Int. Conf. Comput. Vis. Worksh. (ICCV Workshop)*, 2023. [8](#), [22](#)
- [182] Linfeng Zhang, Chenglong Bao, and Kaisheng Ma. Self-distillation: Towards efficient and compact neural networks. *IEEE Trans. Pattern Anal. Mach. Intell.*, 44(8):4388–4403, 2022. [23](#)
- [183] Linfeng Zhang, Xin Chen, Runpei Dong, and Kaisheng Ma. Region-aware knowledge distillation for efficient image-to-image translation. In *Brit. Mach. Vis. Conf. (BMVC)*, 2023.
- [184] Linfeng Zhang, Runpei Dong, Hung-Shuo Tai, and Kaisheng Ma. Pointdistiller: Structured knowledge distillation towards efficient and compact 3d detection. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [23](#)
- [185] Renrui Zhang, Ziyu Guo, Peng Gao, Rongyao Fang, Bin Zhao, Dong Wang, Yu Qiao, and Hongsheng Li. Pointm2AE: Multi-scale masked autoencoders for hierarchical point cloud pre-training. In *Adv. Neural Inform. Process. Syst. (NeurIPS)*, 2022. [6](#), [20](#)
- [186] Renrui Zhang, Ziyu Guo, Wei Zhang, Kunchang Li, Xupeng Miao, Bin Cui, Yu Qiao, Peng Gao, and Hongsheng Li. Pointclip: Point cloud understanding by CLIP. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2022. [2](#), [6](#), [8](#)
- [187] Renrui Zhang, Jiaming Han, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, Peng Gao, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *CoRR*, abs/2303.16199, 2023. [8](#)
- [188] Renrui Zhang, Liuhui Wang, Yu Qiao, Peng Gao, and Hongsheng Li. Learning 3d representations from 2d pre-trained models via image-to-point masked autoencoders. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [6](#), [20](#)
- [189] Shilong Zhang, Peize Sun, Shoufa Chen, Min Xiao, Wenqi Shao, Wenwei Zhang, Kai Chen, and Ping Luo. Gpt4roi: Instruction tuning large language model on region-of-interest. *CoRR*, abs/2307.03601, 2023. [8](#)
- [190] Yue Zhang, Yafu Li, Leyang Cui, Deng Cai, Lemao Liu, Tingchen Fu, Xinting Huang, Enbo Zhao, Yu Zhang, Yulong Chen, Longyue Wang, Anh Tuan Luu, Wei Bi, Freda Shi, and Shuming Shi. Siren’s song in the AI ocean: A survey on hallucination in large language models. *CoRR*, abs/2309.01219, 2023. [19](#)
- [191] Liang Zhao, En Yu, Zheng Ge, Jinrong Yang, Haoran Wei, Hongyu Zhou, Jianjian Sun, Yuang Peng, Runpei Dong, Chunrui Han, and Xiangyu Zhang. Chatspot: Bootstrapping multimodal llms via precise referring instruction tuning. *CoRR*, abs/2307.09474, 2023. [8](#)
- [192] Xi Zhao, He Wang, and Taku Komura. Indexing 3d scenes using the interaction bisector surface. *ACM Trans. Graph.*, 33(3):22:1–22:14, 2014. [2](#)
- [193] Juntian Zheng, Qingyuan Zheng, Lixing Fang, Yun Liu, and Li Yi. CAMS: canonicalized manipulation spaces for category-level functional hand-object manipulation synthesis. In *IEEE/CVF Conf. Comput. Vis. Pattern Recog. (CVPR)*, 2023. [2](#)
- [194] Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena. *CoRR*, abs/2306.05685, 2023. [2](#)
- [195] Junsheng Zhou, Jinsheng Wang, Baorui Ma, Yu-Shen Liu, Tiejun Huang, and Xinlong Wang. Uni3d: Exploring unified 3d representation at scale. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. [6](#), [21](#)
- [196] Yiyang Zhou, Chenhang Cui, Jaehong Yoon, Linjun Zhang, Zhun Deng, Chelsea Finn, Mohit Bansal, and Huaxiu Yao. Analyzing and mitigating object hallucination in large vision-language models. In *Int. Conf. Learn. Represent. (ICLR)*, 2024. [2](#), [19](#)
- [197] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *CoRR*, abs/2304.10592, 2023. [8](#)
- [198] Xiangyang Zhu, Renrui Zhang, Bowei He, Ziyu Guo, Ziyao Zeng, Zipeng Qin, Shanghang Zhang, and Peng Gao. Pointclip v2: Prompting clip and gpt for powerful 3d open-world learning. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. [6](#), [8](#)
- [199] Ziyu Zhu, Xiaoqian Ma, Yixin Chen, Zhidong Deng, Siyuan Huang, and Qing Li. 3d-vista: Pre-trained transformer for 3d vision and text alignment. In *Int. Conf. Comput. Vis. (ICCV)*, 2023. [8](#)



# Contents

<b>1. Introduction</b>	<b>2</b>
<b>2. SHAPELLM</b>	<b>3</b>
2.1. Overall Architecture . . . . .	3
2.2. How to alleviate interactive 3D understanding <i>Data Dessert?</i> . . . . .	3
2.3. RECON++: <i>Scaling Up</i> 3D Representation Learning . . . . .	4
<b>3. 3D MM-Vet: 3D Multimodal Comprehension   Evaluation Benchmark</b>	<b>5</b>
<b>4. Experiments</b>	<b>6</b>
4.1. 3D Representation Transferring with RE- CON++ . . . . .	6
4.2. Multimodal Comprehension with SHAPELLM	7
<b>5. Discussions</b>	<b>7</b>
5.1. What is learned from multi-view alignment?	7
5.2. Is SHAPELLM grounded in physical worlds?	8
5.3. Can SHAPELLM generalize to unseen objects?	8
<b>6. Related Works</b>	<b>8</b>
<b>7. Conclusions</b>	<b>8</b>
<b>A Additional Experiments</b>	<b>17</b>
A.1. Ablation Study . . . . .	17
A.1.1 SHAPELLM Architecture . . . . .	17
A.1.2 Baseline Improvement . . . . .	18
A.2 Multimodal Comprehension with SHAPELLM	18
A.3 Representation Transferring with RECON++	20
<b>B Additional Information about 3D MM-vet</b>	<b>20</b>
B.1. Evaluation System . . . . .	20
B.2. Analysis . . . . .	21
<b>C Implementation details</b>	<b>21</b>
<b>D Training details</b>	<b>22</b>
<b>E Additional Related Work</b>	<b>22</b>
E.1. 3D Representation Learning . . . . .	22
<b>F. Future Works</b>	<b>23</b>

## A. Additional Experiments

### A.1. Ablation Study

#### A.1.1 SHAPELLM Architecture

**Architecture** Let  $\mathcal{F}_\theta$  be the multimodal LLM parameterized by  $\theta$ , we use a RECON++ encoder  $\mathcal{H}_\phi$  as SHAPELLM’s 3D point cloud encoder, followed by three MLP projection layers  $\mathcal{M}_{\zeta^{\text{local}}}$  and  $\mathcal{M}_{\zeta^{\text{global}}}$  for 3D embedding projection of RECON++’s local and global representations, respectively. To facilitate geometry-necessary tasks like 6-DoF pose estimation, we use absolute position encoding (APE) with an MLP projection  $\mathcal{M}_{\zeta^{\text{APE}}}$  to provide additional precise low-level geometric information. Given the original 3D point cloud inputs  $\mathcal{P} = \{\mathbf{p}_i | i = 1, 2, \dots, N\} \in \mathbb{R}^{N \times 3}$  with  $N$  coordinates encoded in a  $(x, y, z)$  Cartesian space. Following previous works [36, 129, 180],  $N_s$  seed points are first sampled using farthest point sampling (FPS). The point cloud  $\mathcal{P}$  is then grouped into  $N_s$  neighborhoods  $\mathcal{N} = \{\mathcal{N}_i | i = 1, 2, \dots, N_s\} \in \mathbb{R}^{N_s \times K \times 3}$  with group centroids from the seed point set  $\mathcal{P}^s$ . The APE representation can be written as

$$\mathbf{E}_{\text{APE}} = \mathcal{M}_{\zeta^{\text{APE}}} \circ \mathcal{P}^s. \quad (2)$$

The local and transformation-invariant 3D geometric embeddings  $\mathbf{x}_i = \text{MAX}_{\mathbf{p}_{i,j} \in \mathcal{N}_i} (\Phi_\gamma(\xi_{i,j}))$  for  $\mathcal{P}_i, i = 1, 2, \dots, N_s$  is used as 3D token embeddings of RECON++, where  $\Phi_\gamma$  is a per-point MLP point feature extractor [126, 127] and  $\xi_{i,j}$  is the feature of  $j$ -th neighbour point  $\mathbf{p}_{i,j}$  in the neighbourhood  $\mathcal{N}_i$ . Let  $\{\mathbf{g}_q^{\text{image}}\}_{q=1}^G$  be  $G$  multi-view image global queries and  $\mathbf{g}^{\text{text}}$  be the global text query. RECON++ outputs the local and global 3D point cloud representations by taking 3D embeddings and global queries as inputs:

$$\left[ \mathbf{e}_{\text{local}}, \mathbf{e}_{\text{global}} \right] = \left[ \mathcal{H}_\phi \left( [\mathcal{P}^s, \{\mathbf{g}_q^{\text{image}}\}_{q=1}^G, \mathbf{g}^{\text{text}}] \right) \right], \quad (3)$$

and the representation to SHAPELLM is:

$$\left[ \mathbf{E}_{\text{local}}, \mathbf{E}_{\text{global}} \right] = \left[ \mathcal{M}_{\zeta^{\text{local}}} \circ \mathbf{e}_{\text{local}}, \mathcal{M}_{\zeta^{\text{global}}} \circ \mathbf{e}_{\text{global}} \right]. \quad (4)$$

In addition, inspired by prefix-tuning [86] and dream queries [37], we append  $Q$ -length learnable embeddings  $\{\mathbf{d}_q^{\text{APE}}\}_{q=1}^Q, \{\mathbf{d}_q^{\text{local}}\}_{q=1}^Q, \{\mathbf{d}_q^{\text{global}}\}_{q=1}^Q$  as visual prompts representation [130]  $\mathbf{E}_{\text{prompt}}$  for adaptively modulating different semantic information encoded in APE, local and global RECON++ representations, respectively.

Formally, the encoded 3D representations to SHAPELLM can be written as:

$$\left[ \{\mathbf{d}_q^{\text{APE}}\}_{q=1}^Q, \mathbf{E}_{\text{APE}}, \{\mathbf{d}_q^{\text{local}}\}_{q=1}^Q, \mathbf{E}_{\text{local}}, \{\mathbf{d}_q^{\text{global}}\}_{q=1}^Q, \mathbf{E}_{\text{global}} \right]. \quad (5)$$

Table 6. **Ablation study on the dedicated designs of SHAPELLM architecture.** The performance of multimodal comprehension on 3D MM-Vet and referring expression grounding on GPartNet with SHAPELLM-13B is reported. Note that  $E_{\text{global}}$  is calculated with both global queries and cross-attention with local 3D embeddings.

$E_{\text{APE}}$	$E_{\text{prompt}}$	$E_{\text{local}}$	$E_{\text{global}}$	3D MM-Vet	GPartNet
✓	✗	✗	✗	30.8	<b>12.3</b>
✓	✓	✗	✗	32.0	11.4
✗	✗	✓	✗	42.2	10.0
✗	✓	✗	✓	50.3	10.5
✓	✗	✓	✓	52.3	10.5
✓	✓	✗	✓	50.3	11.7
✗	✗	✗	✓	52.4	11.7
✗	✗	✓	✓	49.6	10.1
✗	✓	✓	✓	51.7	10.1
✓	✓	✓	✓	<b>53.1</b>	11.7

**Input Components** Tab. 6 shows the ablation study of each input component by supervised fine-tuning with different input representations, demonstrating that it is necessary to employ all designs for achieving decent performance on both 3D comprehension and real-world grounding.

**Visual Prompt Number** Fig. 8 shows the performance of SHAPELLM using different numbers of prompts, including 1, 8, 16, 32, and 64. This ablation study has shown that a different number of prompts leads to varied improvements, and the optimal setting is 32. This observation is similar to VPT [76] where the prompts used to modulate Transformer attention should be studied [59].

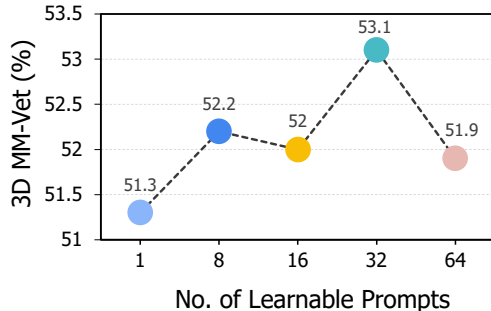


Figure 8. **Ablation study on visual prompt number.** The performance of SHAPELLM-13B on 3D MM-Vet is reported.

### A.1.2 Baseline Improvement

*Can we improve the baseline to bridge the gap between PointLLM and SHAPELLM?* In Tab. 6, we study two technical factors that are contributed by SHAPELLM: *3D point cloud encoder* and *SFT data*.

- *Improvement from encoder.* (Line 1) First, by changing PointLLM’s encoder to RECON++, a significant improvement of +4.20% is obtained. This demonstrates the significantly better 3D representation extraction of RECON++ compared to ULIP-2. It is consistent with previous findings in Tab. 1 and Tab. 2 that RECON++ outperforms

Table 7. **Ablation study on baseline improvements.** Results are tested on 3D MM-Vet with the baseline model PointLLM-13B [166] using different *point cloud encoders* and *SFT data*. Line 0 is the result with the original PointLLM configuration. Line 1 and 2 denote the results of changing the point cloud encoder to RECON++ and using our SFT data. Note that when RECON++ is used as the encoder,  $E_{\text{APE}}$  and  $E_{\text{prompt}}$  are also used by default.

	Encoder	SFT Data	Rec	Know	Gen	Spat	Emb	Total
0	ULIP-2 [170]	PointLLM	46.6	48.3	38.8	45.2	50.9	46.6
1	RECON++	PointLLM	47.5	52.8	43.6	44.9	54.5	50.8
2	RECON++	Ours	46.8	53.0	53.9	45.3	68.4	53.1

ULIP-2 by a large margin regarding 3D representation transferring learning and zero-shot learning.

- *Improvement from data.* (Line 2) As stated in Sec. 2.2, we have constructed instruction-following data for supervised fine-tuning (SFT) using GPT-4V involving comprehensive topics. By further using the SFT data curated by us, PointLLM’s performance gap to SHAPELLM has been fulfilled. This demonstrates the superiority of our SFT data, where the decent quality comes from the more advanced GPT4-V model using multi-view images and the comprehensive topics covered in the data.

## A.2. Multimodal Comprehension with SHAPELLM

**Generative 3D Object Recognition & Captioning** Following PointLLM [166], we conduct generative 3D recognition and captioning experiments. Tab. 8 shows 3D object classification overall accuracy (%) and captioning performance evaluated by GPT-4 and data-driven metrics: Sentence-BERT (S-BERT) [135] and SimCSE [43]. It can be observed that SHAPELLM consistently outperforms other methods across all metrics, demonstrating robust recognition and instruction-following capabilities.

Note that similar to PointLLM’s findings, we also notice that the 3D captioning performance evaluated by traditional metrics like BLEU-1 [120], ROUGE-L [89], and METEOR [5] are highly unreliable in accurately revealing the response quality. This is further demonstrated by human-oriented evaluation, such as the preference win rate comparison presented next.

**Singe-View Point Cloud Inputs** As stated in Sec. 4.2, we construct 3D MM-Vet-C which studies three kinds of corruptions commonly met in real-world scenarios: “single-view”, “jitter”, and “rotate”. Among these corruptions, the “single-view” issue stands out as the most critical challenge since obtaining the objects’ complete point clouds is non-trivial, similar to multi-view images. As a result, everyday real-world robots only get single-view 3D perceptions with sensors such as RGB-D [56]. Fig. 9 shows the qualitative examples of SHAPELLM-13B’s response using single-view point cloud inputs, demonstrating surprisingly outstanding robustness in processing such occluded inputs.

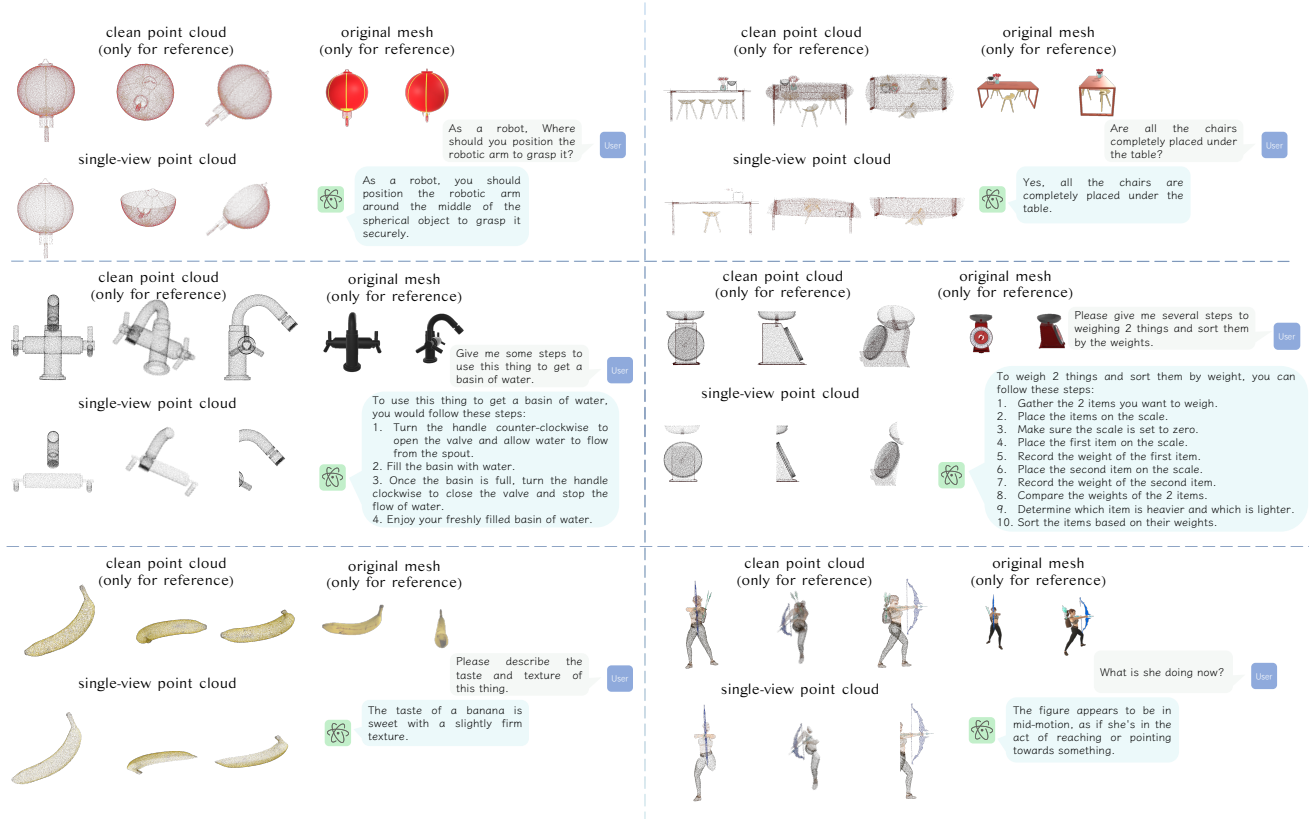


Figure 9. 3D multimodal dialogue using *single-view* point cloud inputs. All answers are generated by SHAPELLM-13B with single-view occluded inputs. SHAPELLM achieves outstanding robustness against such occlusion, which is the most commonly met in real worlds.

Table 8. **Generative 3D recognition and captioning.** The accuracy (%) averaged under the instruction-typed prompt “What is this?” and the completion-typed prompt “This is an object of” is reported.

Method	Input	Classification		Captioning		
		MN-40	Objaverse	GPT-4	S-BERT	SimCSE
InstructBLIP-7B [24]	1-View Img.	25.51	43.50	45.34	47.41	48.48
InstructBLIP-13B [24]	1-View Img.	28.69	34.25	44.97	45.90	48.86
LLaVA-7B [93]	1-View Img.	39.71	50.00	46.71	45.61	47.10
LLaVA-13B [93]	1-View Img.	36.59	51.75	38.28	46.37	45.90
3D-LLM [62]	3D Obj. + Mul.-V. Img.	-	45.25	33.42	44.48	43.68
PointLLM-7B [166]	Point Cloud	52.63	53.00	44.85	47.47	48.55
PointLLM-13B [166]	Point Cloud	52.78	54.00	48.15	47.91	49.12
<b>SHAPELLM-7B</b>	Point Cloud	<b>53.08</b>	<b>54.50</b>	<b>46.92</b>	<b>48.20</b>	<b>49.23</b>
<b>SHAPELLM-13B</b>	Point Cloud	<b>52.96</b>	<b>54.00</b>	<b>48.94</b>	<b>48.52</b>	<b>49.98</b>

**Human Win Rate Comparison** GPT-4 [116] is widely used as an evaluator in natural language and vision language processing, as seen in recent modern benchmarks like MM-Bench and MM-Vet. Recent studies [163] have demonstrated that ChatGPT-based evaluation is more closely *aligned with human preferences* compared to traditional metrics. With GPT4-turbo, the standard deviation of 3D MM-Vet is less than 0.1. To further verify the soundness of the models’ response, we also conduct human evaluation and report the win rate in Fig. 10, where SHAPELLM demonstrates superior preference by humans.

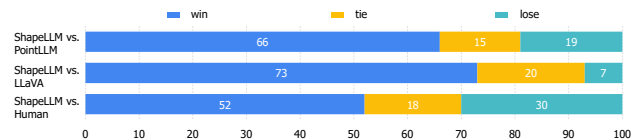


Figure 10. Win rate comparison.

**Visual Hallucination** Visual hallucination is a well-known issue in LLMs and MLLMs that generate non-existent objects or identities from the input data, significantly compromising their multimodal comprehension capabilities [25, 53, 91, 196] and may pose safety risks [110, 137]. Recent research [190] suggests that hallucination may stem from biases in training data, particularly within supervised fine-tuning data, or inappropriate generation strategies. In Fig. 11, we qualitatively demonstrate the illusion evaluation of SHAPELLM compared to other methods. We assess the model’s ability to counteract illusions by prompting it with detailed captions and misleading questions. The results in Fig. 11 demonstrate that previous methods Point-Bind&Point-LLM [54] and PointLLM [166] suffer from the problems of mis-recognition and mis-associating non-existing identities.

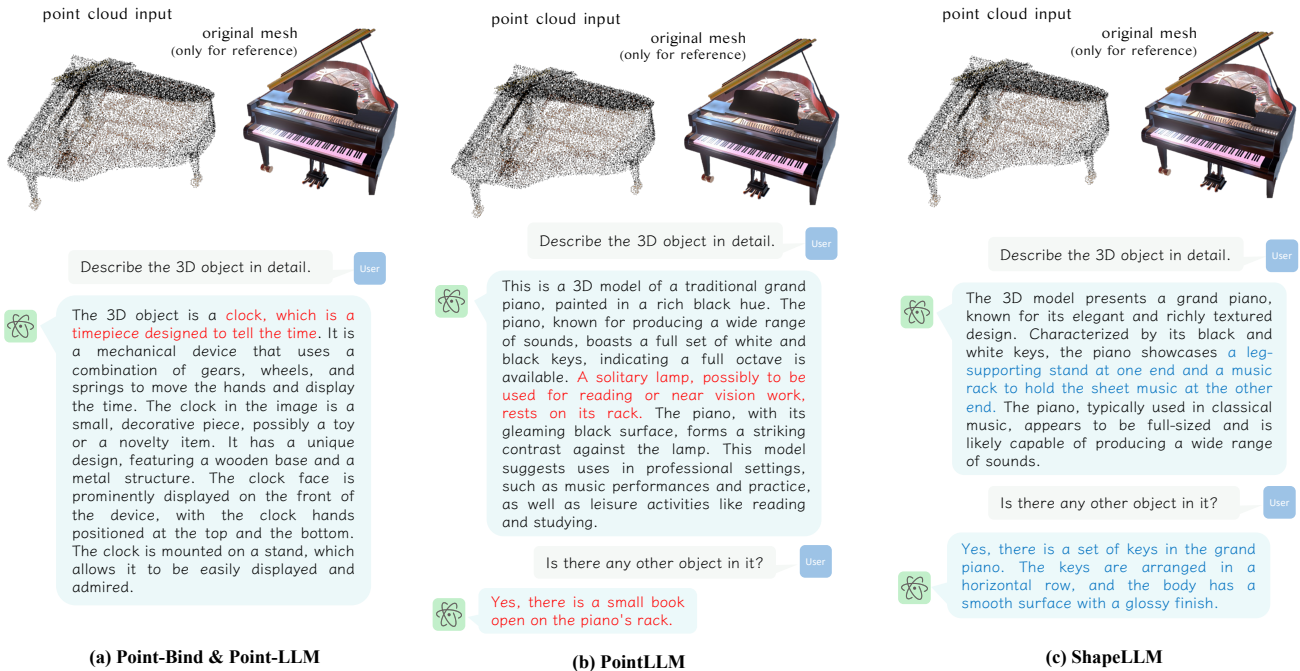


Figure 11. Qualitative comparison examples of visual hallucination.

### A.3. Representation Transferring with RECON++

**Linear SVM Evaluation** Linear SVM evaluation [8, 152] can be used to evaluate the discriminative quality of pre-trained features [51]. The results on ModelNet40 are shown in Tab. 9. It shows that our RECON++ outperforms Point-BERT, which also uses plain Transformers with contrastive objectives, by a clear margin of +6.2%. Compared to hierarchical Transformers methods, our RECON++ outperforms PointM2AE [185] by +0.7%.

Table 9. **Linear SVM classification on ModelNet40.** Overall accuracy (%) without voting is reported.

Method	Hierarchical	ModelNet40
Point-BERT [180]	✗	87.4
PointMAE [119]	✗	91.0
PointM2AE [185]	✓	92.9
ACT [36]	✗	93.1
I2P-MAE [188]	✓	93.4
RECON [129]	✗	93.4
<b>RECON++</b>	✗	<b>93.6</b>

**Few-Shot 3D Object Recognition** Few-shot learning is critical for evaluating the representation transferring capabilities in data and training efficiency. We conduct few-shot 3D object recognition experiments on the ModelNet40 [164] dataset, and the results are shown in Tab. 10. Our RE-

CON++ achieves state-of-the-art performance in all the benchmarks compared to previous works.

Table 10. **Few-shot classification results on ModelNet40.** Overall accuracy (%) without voting is reported.

Method	5-way		10-way	
	10-shot	20-shot	10-shot	20-shot
Transformer [153]	87.8 ± 5.2	93.3 ± 4.3	84.6 ± 5.5	89.4 ± 6.3
Point-BERT [180]	94.6 ± 3.1	96.3 ± 2.7	91.0 ± 5.4	92.7 ± 5.1
Point-MAE [119]	96.3 ± 2.5	97.8 ± 1.8	92.6 ± 4.1	95.0 ± 3.0
Point-M2AE [185]	96.8 ± 1.8	98.3 ± 1.4	92.3 ± 4.5	95.0 ± 3.0
ACT [36]	96.8 ± 2.3	98.0 ± 1.4	93.3 ± 4.0	95.6 ± 2.8
VPP [130]	96.9 ± 1.9	98.3 ± 1.5	93.0 ± 4.0	95.4 ± 3.1
RECON [129]	97.3 ± 1.9	98.9 ± 1.2	93.3 ± 3.9	95.8 ± 3.0
PointGPT [18]	98.0 ± 1.9	99.0 ± 1.0	94.1 ± 3.3	96.1 ± 2.8
<b>RECON++</b>	<b>98.0 ± 2.3</b>	<b>99.5 ± 0.8</b>	<b>94.5 ± 4.1</b>	<b>96.5 ± 3.0</b>

## B. Additional Information about 3D MM-vet

### B.1. Evaluation System

Unlike classification or regression tasks, language generation tasks lack a definitive ground truth that can comprehensively cover diverse real-life scenarios. Therefore, evaluating the alignment of model-generated results with the question and assessing their appropriateness becomes a challenging problem, requiring a reasonable quantitative score. Fortunately,

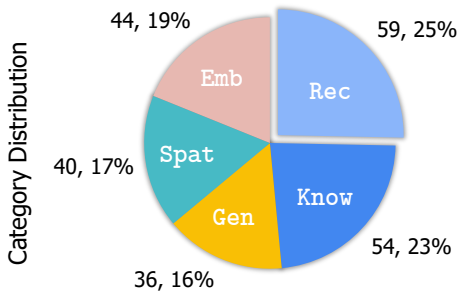


Figure 12. **The number of diverse questions** of *core VL capabilities* on 3D MM-Vet. **Rec**: General Visual Recognition, **Know**: Knowledge, **Gen**: Language Generation, **Spat**: Spatial Awareness, **Emb**: Embodied Interaction.

we have observed the recent surge in the popularity of GPT, providing us with a dependable tool for conducting open-ended evaluations.

To enhance the performance of GPT, we employ a few-shot style in-context prompt. This involves feeding GPT with prompts from evaluative examples and instructing it to generate scores. Specifically, we present prompts to obtain a score ranging from 0 to 1, indicating the degree of similarity between the model-generated answers and the ground truths we provided. When implementing this approach, we observed that results generated multiple times may vary a lot. To address it, we apply the same evaluation setting to a single answer for  $K$  iterations, obtaining the average result as the final score for a precise answer. The score of an answer  $S_a$  and the total score  $S_t$  of answer set  $A$  are calculated by:

$$S_a = \frac{\sum_{i=1}^K s_{a_i}}{K}, \quad S_t = \frac{\sum_{a \in A} S_a}{N}.$$

Here we set  $K = 5$ , and  $s_{a_i}$  is the score of the  $i_{th}$  test of answer  $a$ . The average score for a specific capability is the sum of scores in category  $C$  answer set  $A_C$ :

$$S_c = \frac{\sum_{a \in A_C} S_a}{N_c},$$

where  $N_c$  is the number of answers in each capability set.

To mitigate excessive standard deviation, we opt for GPT-4 in a series of  $K$  scoring rounds to get rounds of outputs with a standard deviation below 0.1. This choice is motivated by the enhanced stability offered by GPT-4 [116], in contrast to GPT-3.5 [114], where scores across different rounds exhibit significant variability.

## B.2. Analysis

The 3D MM-Vet evaluation benchmark consists of 5 different categories of questions. In Fig. 12 we report the distribution of problem categories. The knowledge and General

Table 11. **Details of RECON++ model variants.** This table format follows Dosovitskiy et al..

Model	Layers	Hidden size	MLP size	Heads
RECON++-S	12	384	1536	6
RECON++-B	12	768	3072	12
RECON++-L	24	1024	4096	16

Visual Recognition parts contain multiple subparts that comprehensively evaluate these capacities and thus hold higher proportions. Fig. 13 shows an example of how we prompt GPT-4 for 3D MM-Vet evaluation. Fig. 14 and Fig. 15 illustrate additional examples of 3D MM-Vet Q&As.

## C. Implementation details

**RECON++** Following the standard ViT [38] architecture, we design four different model structures consistent with prior work [94, 129, 195]. The model parameters are shown in Tab. 11. Following OpenShape [94], we employ four datasets as pretraining data, namely Objaverse [29], ShapeNet [13], ABO [22], and 3D-FUTURE [42]. Each point cloud sample has a size of  $10000 \times 6$ , where the first three dimensions represent  $xyz$  coordinates, and the latter three dimensions represent  $rgb$  values.

Regarding the masked modeling strategy, we experimented with both random masking strategies and the latest causal masking strategy. Using causal masking as initialization significantly improves transfer learning capability, as shown in the ablation experiments in Tab. 12. Specifically, the point encoder of SHAPELLM still employs the original local-guided stop-gradient strategy [129]. Additionally, to enhance global classification and retrieval capabilities, we backpropagate gradients from the global branch to the local branch in open vocabulary zero-shot experiments, as demonstrated in the ablation experiments in Tab. 12.

Table 12. **Ablation study on mask type & stop gradient.** transfer: fine-tuned 3D recognition on ScanObjectNN [151]. zero-shot: zero-shot 3D recognition on Objaverse-LVIS [29]. All experiments are conducted on RECON++-L and SHAPELLM-13B.

Mask Type	Stop Grad	Fine-Tune	Zero-Shot	3D MM-Vet
Random	✓	92.5	52.8	<b>53.1</b>
Random	✗	93.6	<b>53.7</b>	52.9
Causal	✓	<b>95.3</b>	49.8	50.7
Causal	✗	92.8	51.0	51.6

**SHAPELLM** We use the LLaMA model [150] as our LLM backbone, with the 7B and 13B Vicuna-1.1 [21] checkpoint as the default settings. We partitioned the point clouds into 512 patches using furthest point sampling and k-nearest neighbors. Similar to other MLLMs [37, 93, 166], we employ a 3-layer MLP with GELU [61] as the projector, with

Table 13. Training recipes for RECON++ and SHAPELLM.

Config	RECON++			SHAPELLM	
	HyBrid/Ensembled	ScanObjectNN	ModelNet	Cap3D	LVIS/GaPartNet
optimizer	AdamW	AdamW	AdamW	AdamW	AdamW
learning rate	5e-5	2e-5	1e-5	2e-3	2e-5
weight decay	5e-2	5e-2	5e-2	-	-
learning rate scheduler	cosine	cosine	cosine	cosine	cosine
training epochs	300	300	300	3	1
warmup epochs	10	10	10	0.03	0.03
batch size	512	32	32	256	128
drop path rate	0.1	0.2	0.2	-	-
number of points	1024/10000	2048	1024/10000	10000	10000
number of point patches	64/512	128	64/512	512	512
point patch size	32	32	32	32	32
augmentation	Rot&Scale&Trans	Rot	Scale&Trans	-	-
GPU device	8×A800	1×A800	1×A800	8×A800	8×A800

hidden layer sizes of 1,024 and 2,048, respectively. Note that different projector parameters are utilized for absolute positional encoding, local, and global features. Through training the projector, multi-scale and multi-mode features of the point cloud are mapped into the text space. After adding two special tokens, the vocabulary size becomes 32,003.

## D. Training details

**RECON++** Due to the sensitivity of the Chamfer Distance [41] loss to accuracy, all experiments were conducted at FP32 precision using 8 × 80G A800 GPUS. We still use the strategy of *contrast with reconstruct* [129]. To save parameter tuning time and improve performance, we divide the training process into two stages: the reconstruction stage based on mask modeling and the cross-modal alignment stage based on knowledge distillation. For transfer learning classification tasks, RECON++ is pretrained on 1,024 points. For zero-shot tasks and SHAPELLM tasks, RECON++ is pretrained on 10,000 points. Further details regarding the hyperparameter settings are documented in Tab. 13.

**SHAPELLM** All experiments were conducted using 8 × 80G A800 GPUs with a BF16 data type. During the multi-modal alignment stage, we train our model for one epoch with a batch size 256 and a learning rate 2e-3. During the instruction tuning stage, we train our model for one epoch with a batch size of 128 and a learning rate 2e-5. Throughout both stages, we employ flash-attention [26], the AdamW [104] optimizer, and a cosine learning rate scheduler [103]. For the entire training process, the 7B and 13B models require approximately 10 and 20 hours, respectively. Further details regarding hyper-parameters are documented in Tab. 13.

## E. Additional Related Work

### E.1. 3D Representation Learning

Research on 3D Representation Learning encompasses various methods, including point-based [126, 127], voxel-based [112], and multiview-based approaches [58, 143]. Point-based methods [40, 131] have gained prominence in object classification [151, 164] due to their sparsity yet geometry-informative representation. On the other hand, voxel-based methods [30, 130, 176] offer dense representation and translation invariance, leading to a remarkable performance in object detection [23] and segmentation [3, 175]. The evolution of attention mechanisms [153] has also contributed to the development of effective representations for downstream tasks, as exemplified by the emergence of 3D Transformers [40, 102, 111]. Notably, 3D self-supervised representation learning has garnered significant attention in recent studies. PointContrast [165] utilizes contrastive learning across different views to acquire discriminative 3D scene representations. Innovations such as Point-BERT [180] and Point-MAE [119] introduce masked modeling [31, 60] pre-training into the 3D domain. ACT [36] pioneers cross-modal geometry understanding through 2D or language foundation models such as CLIP [134] or BERT [31]. Following ACT, RECON [129] further proposes a learning paradigm that unifies generative and contrastive learning. Additionally, leveraging foundation vision-language models like CLIP [36, 134] has spurred the exploration of a new direction in open-world 3D representation learning. This line of work seeks to extend the applicability and adaptability of 3D representations in diverse and open-world/vocabulary scenarios [32, 33, 106, 122, 181].

## F. Future Works

SHAPELLM has made significant progress in advancing 3D shape understanding and embodied perception through MLLMs. Future endeavors aim to scale up embodied understanding training using datasets larger than GPartNet [47], potentially leading to open-vocabulary part-level comprehension, including 6-DoF pose estimation. Excitingly, there is a vision to establish a unified framework capable of comprehending not only 3D shapes but also entire 3D scenes. To enhance real-world applications on robots, a promising approach involves a robotics co-design that effectively connects 3D representations with downstream language-based tasks [72, 77]. Additionally, addressing efficiency for real-time deployment is crucial, emphasizing techniques like model compression [35, 75, 182–184].

**[System Prompt]**

You are a helpful AI assistant.

**[User Prompt]**

Now I will give you a question, the type of the question, an answer from model, and an answer from label. All you need to do is focus on these two answers and figure out whether they are saying the same thing about the specific type of question. Your response should only be a confidence score ranging from 0 to 100. Remember the confidence score is to evaluate how much two answers are describing the same thing. Your response confidence score should follow the scoring standard of the prompt I gave.

Firstly I will give you several question-and-answer pairs as long as their confidence score:

question 1: How many oranges will there be if 1/3 of them are removed?

question type: Knowledge

answer from model: There will be 6 left.

answer from label: As there are 9 oranges in total, there will be 6 oranges left if 1/3 of them are removed.

confidence score: 100

question 2: What is this object?

question type: General Visual Recognition

answer from model: This is a bathtub

answer from label: This is a dirty bathtub.

confidence score: 80

question 3: What is this object?

question type: General Visual Recognition

answer from model: This is a bottle of water

answer from label: This is a bottle of oil

confidence score: 50

question 4: What is holding in this boy's right hand?

question type: Spatial Recognition

answer from model: He is holding a white cup in his right hand.

answer from label: He is holding a sword in his right hand.

confidence score: 0

Next, I will give you the elements:

question: {},

question type: {},

answer from model: {},

answer from label: {}.

Please remember, while outputting the confidence score, do not include any words, just the number.

**[Example1]**

Question: When grasping it, should you apply a greater or lesser force?

question type: Embodied Interaction

answer from model: When grasping it, you should apply a lesser force to avoid damaging the doughnut.

answer from label: His cloak was attached to his shoulder armor and hung on his back.

confidence score: 90

**[Example2]**

Question: Where is its cloak?

question type: Spatial Recognition

answer from model: The cloak is draped over its shoulders.

answer from label: A lesser force.

confidence score: 100

Figure 13. GPT-4 evaluation template and examples of our 3D MM-Vet benchmark.



**[Point Clouds]**

**[Question<sup>1</sup> Type: Embodied Interaction]**

Please give me steps to roll the dice until it hits 6.

**[Answer<sup>1</sup>]**

Step 1: Roll the dice.

Step 2: If the dice hits 1, 2, 3, 4, or 5, roll it again.

Step 3: If the dice hits 6, stop rolling and show the roll of the dice.



**[Question<sup>2</sup> Type: Spatial Awareness]**

Change the 1 and 5 on the object, which number will be the opposite of the number 5?

**[Answer<sup>2</sup>]**

Change the number 1 and number 5 on this dice, number 6 will be on the opposite of number 5.

**[Question<sup>3</sup> Type: Knowledge]**

At which places will this object be used?

**[Answer<sup>3</sup>]**

This object is a dice, so it might be used at places like bars and gambling houses.

**[Point Clouds]**

**[Question<sup>1</sup> Type: General Visual Recognition]**

What subparts are there in the scene?

**[Answer<sup>1</sup>]**

There is a bag of cookies, a mug of milk and a China bowl.

**[Question<sup>2</sup> Type: Embodied Interaction]**

As an AI robot, please give me steps to mix the milk and cookies in a bowl.

**[Answer<sup>2</sup>]**

Step 1: Pour the milk into the bowl.

Step 2: Put the cookies into the bowl.

Step 3: Stir with a spoon.

**[Question<sup>3</sup> Type: Knowledge]**

Describe the physical properties of the milk.

**[Answer<sup>3</sup>]**

The milk is a kind of liquid with a white color, whose density and boiling point is higher than water while the freezing point is lower than water, has a mild, slightly sweet odor and taste.



**[Point Clouds]**

**[Question<sup>1</sup> Type: Embodied Interaction]**

I want to change the place of the spoon and the fork, please give me steps.

**[Answer<sup>1</sup>]**

Step 1: Pick up the fork and the spoon.

Step 2: Put down the spoon at the place of the fork.

Step 3: Put down the fork at the place of the spoon.

**[Question<sup>2</sup> Type: Language Generation]**

What are these objects commonly used to do in usual life? Please speak in detail.

**[Answer<sup>2</sup>]**

Forks, spoons, and table knives are common eating utensils used in everyday life for various purposes. Forks are primarily used for piercing and picking up solid food item. Spoons are used for scooping and conveying liquids, semi-liquids, and foods with a sauce or broth. Table knives, also known as dinner knives, are used for cutting and slicing food on your plate. They are typically not as sharp as kitchen knives, as their primary purpose is to assist with cutting while dining.



Figure 14. Additional Visualization example of 3D MM-Vet Q&A pairs.

**[Point Clouds]**

**[Question<sup>1</sup> Type: General Visual Recognition]**

What is this object?

**[Answer<sup>1</sup>]**

This is an ukulele.

**[Question<sup>2</sup> Type: Knowledge]**

What is the difference in shape between this object and the guitar?

**[Answer<sup>2</sup>]**

This is a ukulele that has 4 strings, while the guitar has 6 strings.

**[Question<sup>3</sup> Type: Knowledge]**

What are the physical properties when this object meets fire?

**[Answer<sup>3</sup>]**

It's basically made of wood and nylon so fire would cause great damage to it.



**[Point Clouds]**

**[Question<sup>1</sup> Type: General Visual Recognition]**

Who is this anime character?

**[Answer<sup>1</sup>]**

This is SpongerBob.

**[Question<sup>2</sup> Type: Language Generation]**

This is a classic anime character, please describe the anime it appears.

**[Answer<sup>2</sup>]**

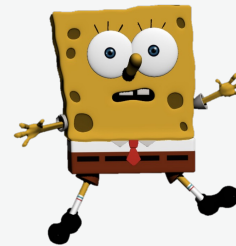
Spongebob Squarepants is an American animated comedy series set in the Pacific Ocean in a city called Bikini Bottom. The plot content of SpongeBob Squarepants basically has nothing to do with Marine knowledge, and even exaggerates to be completely inconsistent with science and common sense, such as making fire under the sea, bathing under the sea, having a lake under the sea (Cool Lake), drowning fish in the lake (Cool Lake), and burning underwater buildings.

**[Question<sup>3</sup> Type: Language Generation]**

Tell me about his best friend Patrick Star.

**[Answer<sup>3</sup>]**

Pink Starfish, Spongebob's best friend, is also one of Squidward's hate objects, wearing a pair of purple flower beach shorts, in the series always with SpongeBob to make a lot of jokes. He has no job, is a loyal customer of the Krusty Krab, and also likes to eat Haiba paste. Sometimes he works for the Krusty Krab or Sea Bully, but only as a one-day employee, and most of the work is messed up, but the artistic talent is more than Squidward and SpongeBob Squarepants.



**[Point Clouds]**

**[Question<sup>1</sup> Type: Embodied Interaction]**

Give me several steps to take the rusty barrel away from this pack.

**[Answer<sup>1</sup>]**

Step 1: Clamp the rusty barrel.

Step 2: Take it down from the height.

Step 3: Turn around and take it away from the pack.

**[Question<sup>2</sup> Type: Spatial Recognition]**

Where is the rusty barrel?

**[Answer<sup>2</sup>]**

The rusty barrel is in the top row, next to the yellow one.

**[Question<sup>3</sup> Type: Spatial Recognition]**

Please describe the spatial relation of this entirety.

**[Answer<sup>3</sup>]**

The barrels are stacked in two layers, the bottom layer is three yellow barrels, and the top layer is a yellow barrel and a rusted barrel in the gap between the bottom three buckets.



Figure 15. Additional Visualization example of 3D MM-Vet Q&A pairs.