# PointLLM: Empowering Large Language Models to Understand Point Clouds

Runsen Xu[1,2]    Xiaolong Wang[3]    Tai Wang[2]    Yilun Chen[2]    Jiangmiao Pang[2✉]    Dahua Lin[1,2]

[1]The Chinese University of Hong Kong    [2]Shanghai AI Laboratory    [3]Zhejiang University

{runsenxu,dhlin}@ie.cuhk.edu.hk, xlking@zju.edu.cn

{wangtai,chenyilun,pangjiangmiao}@pjlab.org.cn

## Abstract

*The unprecedented advancements in Large Language Models (LLMs) have shown a profound impact on natural language processing but are yet to fully embrace the realm of 3D understanding. This paper introduces PointLLM, a preliminary effort to fill this gap, enabling LLMs to understand point clouds and offering a new avenue beyond 2D visual data. PointLLM understands colored object point clouds with human instructions and generates contextually appropriate responses, illustrating its grasp of point clouds and common sense. Specifically, it leverages a point cloud encoder with a powerful LLM to effectively fuse geometric, appearance, and linguistic information. We collect a novel dataset comprising 660K simple and 70K complex point-text instruction pairs to enable a two-stage training strategy: aligning latent spaces and subsequently instruction-tuning the unified model. To rigorously evaluate the perceptual and generalization capabilities of PointLLM, we establish two benchmarks: Generative 3D Object Classification and 3D Object Captioning, assessed through three different methods, including human evaluation, GPT-4/ChatGPT evaluation, and traditional metrics. Experimental results reveal PointLLM's superior performance over existing 2D and 3D baselines, with a notable achievement in human-evaluated object captioning tasks where it surpasses human annotators in over 50% of the samples. Codes, datasets, and benchmarks are available at https://github.com/OpenRobotLab/PointLLM.*

## 1. Introduction

Recent years have witnessed the emergence of large language models (LLMs) [6, 8, 42–44, 50, 56, 57], demonstrating awe-inspiring abilities in natural language processing. These models have become versatile tools, acting as generalized interfaces [22] to perform an array of complex
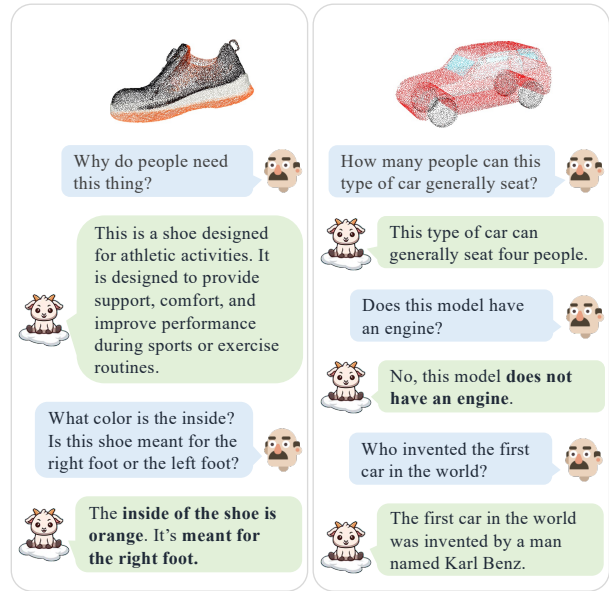
---

✉Corresponding author.



Figure 1. We introduce PointLLM, a multi-modal large language model capable of understanding colored point clouds of objects. It perceives object types, geometric, and appearance without concerns for ambiguous depth, occlusion, or viewpoint dependency.

tasks[6, 50]. However, the mastery over text-based tasks is just one aspect of what LLMs can achieve. A new horizon emerges as researchers begin to explore multi-modal LLMs, capable of processing various forms of data such as audio[26] and images[1, 27, 35, 37, 43, 70, 72].

The next step in this evolution lies in understanding 3D structures. Imagine a scenario where one can interactively create and edit 3D content through verbal commands[31, 41], or instruct a robot to manipulate objects using natural language[12]. These applications require LLMs with a nuanced and accurate understanding of 3D structures.

While existing efforts to integrate LLMs with 2D images[9, 12, 37, 72] provide a pathway to understanding 3D structures, they face difficulties such as depth ambiguity,

1

occlusion, and viewpoint dependency. Solutions like selecting optimal views or using multi-view images exist but they may be elusive due to objects' arbitrary orientations and can increase model complexity. In contrast, point clouds offer an efficient and universal 3D representation. They provide direct geometric and appearance data, enabling a more comprehensive understanding of 3D shapes, effective occlusion management, and viewpoint-independent analysis. Despite these benefits, the integration of point clouds with LLMs is still a relatively uncharted area.

In this work, we pave the way to empower large language models to understand point clouds, with a preliminary focus on 3D objects. Specifically, we present PointLLM, which accepts colored object point clouds with human instructions and generates accurate responses, reflecting its understanding of point clouds and common sense, as illustrated in Fig. 1. Enhancing LLMs' understanding of 3D object point clouds presents three problems: the absence of training data, the necessity of building a suitable model architecture, and the lack of comprehensive benchmarks and evaluation methods, each of which is addressed as follows.

**Data collection.** We collect a large-scale point-text instruction following dataset, containing 660K brief-description instructions and 70K complex instructions. The training data that guides the model in extracting meaningful representations from point clouds and responding to user instructions are especially rare in the context of object point clouds, and manual collection can be both time-consuming and expensive. To circumvent this issue, we utilize the recently introduced Cap3D [40], a large-scale 3D object captioning dataset built upon Objaverse [11]. Employing the reasoning abilities and world model of GPT-4[43], we prompt GPT-4 to generate varied instruction following data based on the contexts provided by the captions.

**Model and training.** We introduce PointLLM, which employs a pre-trained point cloud encoder for encoding point clouds into tokens and utilizes a powerful pre-trained large language model for reasoning and generating responses. Our training features a two-stage strategy[37]: alignment of the latent spaces between the encoder and the LLM, followed by instruction-tuning the unified model. This methodology ensures an effective fusion of both geometric and appearance information from point clouds with the linguistic capabilities of the language model. We also provide some empirical studies of the model's design choices.

**Benchmarks and evaluation.** We establish two distinct benchmarks: Generative 3D Object Classification and 3D Object Captioning, accompanied by a comprehensive evaluation framework, to assess the model's understanding of point clouds. Due to the generative nature, our models are directly prompted to engage in object classification on ModelNet40[64] and Objaverse[11], along with Objaverse-based captioning. As defining a single evaluation met-

ric for generative tasks is difficult, we employ three types of evaluation methods, including human evaluation, GPT-4/ChatGPT[42] evaluation, and traditional metric[4, 15, 36, 45, 51] to rigorously assess our model's perceptual and generalization capabilities.

Experimental results show that our PointLLM demonstrates substantially better performance over 2D and 3D baselines, and in over 50% of samples of object captioning, it gains higher scores than human annotators in human evaluation. To supplement these quantitative evaluations, we present a range of qualitative examples, offering a broader perspective on the real-world performance of PointLLM.

## 2. Related Work

**Multi-modal large language models.** Multi-modal Large Language Models (MLLMs) are designed to comprehend and interpret a wide range of information that extends beyond mere text-based data[67], including but not limited to images[17, 27, 37, 60, 72], audio[26], motion[30], etc. Broadly, the models can be classified into two categories. The first category includes models that employ a large language model to interface with individual, modality-specific models or APIs [20, 26, 46, 55, 63]. This approach circumvents the need for additional model training but is heavily dependent on the availability and capabilities of pre-existing models or APIs. The second category pertains to models that employ an end-to-end training strategy. There are two prominent paradigms within this category. The first involves training the model from scratch, similar to text-only LLMs, using large-scale multi-modal corpora and datasets [27, 47]. The second paradigm builds on pre-trained LLMs and unimodal encoders, thereby avoiding training from scratch [1, 3, 9, 12, 14, 17, 33, 35, 37, 53, 70–72]. This strategy typically involves a two-stage process: alignment of the unimodal encoder with the LLM's feature space, followed by instruction-based fine-tuning. In our work, we adhere to the alignment and tuning strategy, intending to construct an MLLM capable of understanding 3D object point clouds.

**Object point cloud understanding with language.** Inspired by models like CLIP [49], which bridges visual and textual modalities, similar advancements have emerged in the 3D object domain[23, 28, 38, 59, 65, 66, 69, 73]. PointCLIP[69], PointCLIPv2[73], and CLIP2Point[28] utilize depth image projections of point clouds for 3D recognition with 2D CLIP models. Others, such as ULIP[65], JM3D[59], OpenShape[38], and CG3D[23], train point cloud encoders to align with CLIP representations using triplets of point clouds, images, and texts. ULIP-2[66] and OpenShape[38] have expanded this by employing image-captioning models for automatic data generation, enhancing training triplet scalability. Cap3D [40] and UniG3D [54] adopt similar approaches for point-text dataset generation. In our work, we leverage Cap3D's captions on Ob-

javerse for automatic instruction-data generation in training PointLLM. The recently introduced 3D-LLM[25] also seeks to enable LLMs to comprehend 3D, by rendering objects into multi-view images, using 2D foundational models like CLIP[49] and SAM[32] for feature extraction, and 2D MLLMs such as BLIP[35] for output generation. Concurrently, Point-Bind LLM[18] aligns point cloud features with ImageBind[16] and uses 2D MLLMs like Imagebind-LLM[21] for generation. Though simple, it faces challenges like hallucination due to its retrieval nature. Distinctively, PointLLM provides direct and comprehensive understanding of object point clouds by end-to-end training, enabling accurate, open-ended, and free-form interactions.

## 3. Methodology

This section elucidates our strategy for the automatic generation of point-text instruction-following data. We then delve into the architecture of our model, PointLLM, which takes as input an object point cloud and user instruction and outputs corresponding responses. Lastly, we detail our loss function and two-stage training strategy.

### 3.1. Point-Text Instruction Following Data

The daunting challenge in the development of an end-to-end multi-modal LLM is procuring large-scale multi-modal instruction-following data, vital for representation learning, aligning latent spaces, and orienting the model to adhere to human intentions[1, 9, 34, 37, 72]. However, manual labeling of such data is cost-prohibitive and labor-intensive. To overcome this, we propose an automated data generation technique utilizing the large-scale point cloud captioning dataset, Cap3D[40], with the assistance of GPT-4[43]. The generated dataset adheres to a uniform instruction following template, shown in Tab. 1, and consists of brief-description instructions and complex instructions, which aid in latent space alignment and instruction tuning, respectively.

**Brief-description instructions.** The Cap3D[40] dataset provides two variations of captions for the 3D objects in Objaverse[11]: those generated by image-captioning models and those annotated by humans. While there are 660K objects accompanied by generated captions, only 40K samples have human-annotated captions. For brief-description instruction, we utilize the model-generated split due to the need for a larger data volume for aligning the latent spaces of point cloud and text modalities [37]. We created a list of 30 instructions to instruct the model to provide a succinct description of a given 3D object point cloud. A random instruction from this list is chosen as the user instruction, and the caption from Cap3D is used directly as the model response, forming a single-round instruction following sample. This results in 660K brief-description instruction data, each corresponding to a unique object point cloud.

Table 1. **Instruction following template.** {System Prompt} is the system prompt used by the pre-trained LLM, {p_tokens} are point tokens, and {Instruction} and {Response} denote user instructions and model responses. Losses are computed only for model responses and the end-of-sentence token </s>.

| {System Prompt} | |
| --- | --- |
| USER: | \<p_start>{p_tokens}\<p_end>{Instruction 1} |
| ASSISTANT: | {Response 1}</s> |
| USER: | {Instruction 2} |
| ASSISTANT: | {Response 2}</s> |
| USER: | {Instruction 3} |
| ASSISTANT: | {Response 3}</s> |

**Complex instructions.** Beyond brief descriptions, it's crucial that the model learns to understand objects from a variety of angles, responding accurately to diverse human instructions. To facilitate this, we employ GPT-4 to produce complex instruction-following data. Specifically, a caption from Cap3D is used to stimulate GPT-4 to craft a more comprehensive description that identifies the object's type, appearance, functionalities, and any other inferable information. Similar to the process for generating brief-description instructions, we also curate a set of 30 distinct prompts, each pushing the model to describe the 3D object in depth. One of these prompts is randomly coupled with the newly crafted description, forming a training sample. GPT-4 is further used to generate conversations (*i.e.*, Q&A pairs) that delve into diverse aspects of the object based on the captions, such as the object's functionality or materials, and the corresponding answers should be informative and comprehensive. For each object, GPT-4 generates 3 single-round conversations and 1 multi-round conversation with 3 Q&A pairs, all ensuring logical relevance.

With a focus on data quality, we selected 15K captions from the Cap3D human-annotated split for data generation, each comprising more than five words. After filtering incorrect GPT-4 outputs, we collected 70K complex instructions, including 15K detailed descriptions, 40K single-round conversations, and 15K multi-round conversations. The instruction lists, GPT-4 prompts, data examples, and distribution analysis can be found in App. A.

### 3.2. Model Architecture

As shown in Fig. 2, our PointLLM is a generative model that aims to complete multi-modal sentences containing both point clouds and texts. The model consists of three main components: a pre-trained point cloud encoder $f_{pe}$, a projector $f_{proj}$, and a pre-trained large language model (LLM) backbone $f_{llm}$.

The point cloud encoder $f_{pe}$ takes as input a point cloud $P \in \mathbb{R}^{n \times d}$, where $n$ is the number of points and $d$ is the feature dimension of each point. The output of the encoder
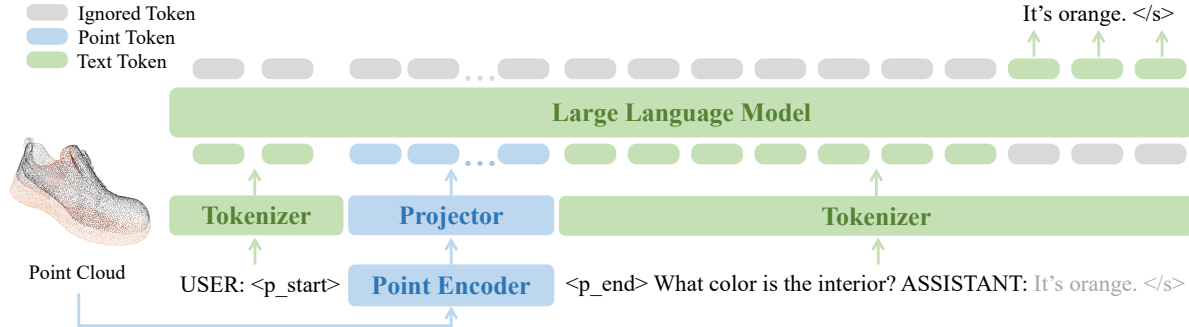
Figure 2. **An overview of PointLLM.** The point encoder extracts features from the input point cloud and the projector projects them to the latent space of the LLM backbone. The LLM backbone processes sequences of point and text tokens, and generates the predicted tokens as the output. The model is trained with a cross-entropy loss on the tokens corresponding to the model responses.

is a sequence of point features $X = (x_1, x_2, \ldots, x_m) \in \mathbb{R}^{m \times c}$, where $m$ is the number of point features and $c$ is the feature dimension. The projector $f_{proj}$ is a MLP that maps the point features $X$ to point tokens $Y = (y_1, y_2, ..., y_m) \in \mathbb{R}^{m \times c'}$, where $c'$ is the dimension of the point tokens, which is the same as the text tokens.

The LLM backbone $f_{llm}$ is a decoder-only Transformers [58], which accepts a sequence of tokens, composed of both text and point tokens. This mixed sequence of tokens is denoted as $Z = (z_1, z_2, ..., z_k) \in \mathbb{R}^{k \times c'}$, where $k$ is the total number of tokens. Utilizing a self-attention mechanism, the LLM backbone is capable of understanding the contextual relationships between different types of tokens, enabling it to generate responses based on both text and point cloud inputs. Formally, the output of the LLM backbone $f_{llm}$ is a sequence of predicted tokens $\hat{Z} = (\hat{z}_1, \hat{z}_2, ..., \hat{z}_k) \in \mathbb{R}^{k \times c'}$. The prediction of the $i$-th token, $\hat{z}_i$, is conditioned on all previous tokens, $Z_{<i} = (z_1, ..., z_{i-1})$. This can be expressed mathematically as

$$\hat{z}_i = f_{llm}(Z_{<i}). \tag{1}$$

Each $\hat{z}_i$ is passed through a final linear layer followed by a softmax operation, mapping the hidden states into a probability distribution over the vocabulary. This additional layer is denoted as $f_{vocab} : \mathbb{R}^{c'} \to \mathbb{R}^V$, where $V$ is the size of the vocabulary. The final prediction $\tilde{z}_i$ for the $i$-th token is the word in the vocabulary with the highest probability:

$$\tilde{z}_i = \arg \max_{w \in \text{vocab}} f_{vocab}(\hat{z}_i)[w]. \tag{2}$$

### 3.3. Training

**Loss function.** We train PointLLM by minimizing the negative log-likelihood of the text token at each position. Our loss function is only computed on text tokens that constitute the model's responses, including the end-of-sentence token </s>. We exclude the tokens from human instructions, ensuring that the model focuses on learning to generate accurate and coherent responses. The end-to-end nature of this training approach enables PointLLM to effectively integrate point cloud and text modalities.

**Two-stage training.** Our training procedure comprises two stages, each focusing on different aspects of the model.

During the first stage, termed the **feature alignment stage**, we freeze the parameters of the point cloud encoder and the LLM, and train only the MLP projector. At this stage, the training process uses brief-description instructions, aiming to align point features with the text token space effectively. This stage also includes the adjustment of token embeddings for the two newly added special tokens <p_start> and <p_end>.

In the second stage, referred to as the **instruction tuning stage**, we freeze the point cloud encoder while jointly training the projector and the LLM. This second stage uses complex instructions and helps the model build its ability to understand and respond to complex instructions including point cloud data.

## 4. Benchmarks and Evaluation

Evaluating the performance of a multi-modal LLM is challenging, as it's difficult to define a single metric that can capture the quality and diversity of the generated outputs. Moreover, existing benchmarks for 3D point cloud understanding are mostly based on discriminative tasks such as close-set classification or retrieval, which do not fully reflect the generative nature and open-vocabulary setting of our model. Therefore, we propose two novel benchmarks to assess our model's perceptual abilities and generalization power: Generative 3D Object Classification and 3D Object Captioning. We adopt various evaluation methods for assessing performances including human evaluation, GPT-4/ChatGPT evaluation, and traditional metric. We use GPT-4/ChatGPT as evaluators, as they demonstrate abilities to align with human judgment accurately. Please refer to App. B.2 and App. B.3 for the prompts we use and for the human verification of the GPT evaluators' correctness.

4

Table 2. **Generative 3D object classification results on the ModelNet40 test split and Objaverse.** The results show the classification accuracy under the **I**nstruction-typed (I) prompt "What is this?" and the **C**ompletion-typed (C) prompt "This is an object of ".

| Model | Input | ModelNet40 (I) | ModelNet40 (C) | Objaverse (I) | Objaverse (C) | Average |
|---|---|---|---|---|---|---|
| InstructBLIP-7B[9] | Single-V. Img. | 19.53 | 31.48 | 45.00 | 42.00 | 34.50 |
| InstructBLIP-13B[9] | Single-V. Img. | 25.97 | 31.40 | 37.00 | 31.50 | 31.47 |
| LLaVA-7B[37] | Single-V. Img. | 39.75 | 39.67 | 49.50 | 50.50 | 44.86 |
| LLaVA-13B[37] | Single-V. Img. | 37.12 | 36.06 | 53.00 | 50.50 | 44.17 |
| 3D-LLM[25] | 3D Obj. + Mul.-V. Img. | - | - | 49.00 | 41.50 | 45.25 |
| Point-Bind LLM[18] | 3D Point Cloud | 51.90 | 39.71 | 6.00 | 4.50 | 25.53 |
| **PointLLM-7B** | 3D Point Cloud | **53.44** | 51.82 | 55.00 | 51.00 | 52.82 |
| **PointLLM-13B** | 3D Point Cloud | 53.00 | **52.55** | **56.50** | **51.50** | **53.39** |

## 4.1. Generative 3D Object Classification

The task of generative 3D object classification involves prompting the model to generate the object type from its point cloud, distinguishing it from discriminative models that directly classify objects based on probability comparisons. We consider two scenarios for this task: close-set zero-shot classification and open-vocabulary classification.

**Close-set zero-shot classification.** In this scenario, the object type belongs to a fixed set of categories, and the model never sees any samples of this dataset during training. This tests the model's ability to generalize to unseen domains using its prior knowledge. We use the test split of the ModelNet40 [64] dataset as our source of data, which contains point clouds of 40 different object categories. The model is prompted to answer the object types in free form and we use ChatGPT as a post-processor to select one of the ModelNet40 categories based on the model's answer. If ChatGPT selects the correct option, then we consider the model's classification correct; otherwise, we consider it incorrect. Please refer to App. B.1 for more discussions about this task's setting.

**Open-vocabulary classification.** In this scenario, the object type is not limited to a predefined set of categories, but can be any word or phrase that identifies the object. This reflects the real-world setting where new objects can appear at any time, and the model needs to be able to recognize them without retraining. We use the Objaverse[11] dataset with human-annotated captions from Cap3D[40] as our source of data. We randomly select 200 objects from the data and use the human captions as ground truth labels. We prompt our model to classify the object with point clouds as input and collect the model's output for each object. Then we use GPT-4 as an evaluator to classify whether the model's response and the human caption are referring to the same object type. We do not require the model's response to match exactly with the human caption, as long as it conveys the same object type. For example, if the human caption is "a blue mug", then "a cup" and "a coffee mug" are all correct predictions. We opt for GPT-4 over ChatGPT in this sce-

nario, because ChatGPT tends to produce more false negatives, meaning that it considers two words or phrases are not referring to the same object type, even when they are, while GPT-4 demonstrates accurate recognition.

## 4.2. 3D Object Captioning

3D object captioning involves generating a natural language description of an object, given its point cloud representation. This is a fine-grained evaluation compared to the classification task. In our benchmark, we utilize the same 200 objects previously used for the open-vocabulary classification, and prompt our model to caption them. Human-annotated captions corresponding to these objects serve as reference ground truths for automatic evaluation.

For a comprehensive and robust evaluation, we employ three distinct methods to assess performance in this task:

1. **Human evaluation.** Human evaluators review randomly shuffled captions from various models alongside human-annotated captions for the objects. Using the official Objaverse[11] explorer, evaluators visually inspect each object and assign two scores to the captions: a correctness score and a hallucination score. The correctness score gauges the accuracy of the model in recalling object attributes (such as type, color, material, *etc.*), while the hallucination score assesses the severity of any fabricated details. Each attribute, whether correct or hallucinated, is awarded one point. Additionally, precision is calculated as the proportion of correct information within the model-generated content. For detailed scoring criteria, please refer to App. B.4.

2. **GPT-4 evaluation.** As human evaluation is both time-consuming and costly, we also employ GPT-4 as an evaluator. Given a model-generated caption and its corresponding human reference, GPT-4 identifies the aspects mentioned in the human caption and calculates the percentage of these aspects that are either correctly mentioned or partially matched in the model's caption, scoring from 0 to 100.

3. **Traditional metric evaluation.** In addition, we em-

Table 3. **3D object captioning results on Objaverse.** Models are evaluated using human evaluation, GPT-4 evaluation, and traditional metrics. A primary focus is placed on human and GPT-4 evaluation, along with data-driven metrics (Sentence-BERT and SimCSE).

| Model | Correctness | Hallucination↓ | Precision | GPT-4 | Sentence-BERT | SimCSE | BLEU-1 | ROUGE-L | METEOR |
|---|---|---|---|---|---|---|---|---|---|
| InstructBLIP-7B[9] | 2.56 | 0.77 | 76.99 | 45.34 | 47.41 | 48.48 | 4.27 | 8.28 | 12.99 |
| InstructBLIP-13B[9] | 2.58 | 1.13 | 69.56 | 44.97 | 45.90 | 48.86 | 4.65 | 8.85 | 13.23 |
| LLaVA-7B[37] | 2.76 | 0.86 | 76.30 | 46.71 | 45.61 | 47.10 | 3.64 | 7.70 | 12.14 |
| LLaVA-13B[37] | 2.43 | 0.86 | 73.97 | 38.28 | 46.37 | 45.90 | 4.02 | 8.15 | 12.58 |
| 3D-LLM[25] | 1.77 | 1.16 | 60.39 | 33.42 | 44.48 | 43.68 | **16.91** | **19.48** | **19.73** |
| **PointLLM-7B** | 3.04 | **0.66** | **82.14** | 44.85 | 47.47 | 48.55 | 3.87 | 7.30 | 11.92 |
| **PointLLM-13B** | **3.10** | 0.84 | 78.75 | **48.15** | **47.91** | **49.12** | 3.83 | 7.23 | 12.26 |
| Human | 2.67 | 0.22 | 92.46 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 | 100.00 |



Figure 3. **Win rate comparison**. PointLLM outperforms human annotations in more than half of the testing samples and exhibits a substantial advantage over other models.

ploy traditional metrics such as BLEU-1 [45], ROUGE-L [36], and METEOR [4]. Though widely used, these metrics often fall short in accurately evaluating generative tasks, as they primarily measure the overlap of n-grams or their varieties, and account less for the semantic similarity or diversity of the captions. Therefore, we incorporate two additional data-driven metrics, Sentence-BERT [51] and SimCSE[15] similarity, which compute the similarity of sentence embeddings between model-generated and human captions.

## 5. Experimental Results

### 5.1. Experimental Settings

**Implementation details.** We use the LLaMA model[57] as our LLM backbone, with the 7B and 13B Vicuna[7] checkpoint. Point-BERT[68], pre-trained with ULIP-2[66] on the Objaverse [11] dataset, serves as our point encoder. The 200 objects from Objaverse utilized for our benchmarks are not seen during any stage of the training. We utilize $n = 8192$ points and $d = 6$ dimensions for each point cloud. We assign a black color to point clouds from ModelNet40, as they lack color information. The point encoder outputs $m = 513$ point features, each with $c = 384$ dimensions. The projec-

Table 4. **Traditional metrics for different captions**. The biased scores demonstrate the limitations of these metrics.

| Caption | BLEU-1 | ROUGE-L | METEOR |
|---|---|---|---|
| **Private jet** | **100.00** | **100.00** | **100.00** |
| there is a black jet engine in a dark background | 10.00 | 18.18 | 17.86 |
| This is a 3D model of a cartoon-style commercial airplane. | 0.00 | 0.00 | 0.00 |

tor contains three linear layers with the GeLU[24] activation, which maps point features to tokens with $c' = 5120$ (7B model) or $c' = 5120$ (13B model) dimensions. As we add two additional special tokens, the vocabulary size of PointLLM is $V = 32003$. All experiments are conducted on $8 \times 80$G A100 GPUs. GPT-4 and ChatGPT in this paper all refer to OpenAI's "gpt-4-0613" and "gpt-3.5-turbo-0613" models respectively. More implementation and training details are provided in App. C.

**Baselines.** Our models are mainly compared against 3D-LLM[25] and Point-Bind LLM[18]. 3D-LLM, requiring 3D objects and multi-view images, is evaluated only on the Objaverse dataset due to its current lack of support for object point clouds. Point-Bind LLM, incompatible with colored point clouds, is included solely in object classification experiments. We also include two popular open-sourced 2D MLLMs, InstructBLIP[9] and LLaVA[37].

### 5.2. Generative 3D Object Classification

Tab. 2 shows the classification accuracy of various models on our proposed generative 3D object classification tasks. For 2D MLLMs' image inputs, we randomly sample rendered images of ModelNet40 point clouds and Objaverse objects. We prompt all the models with the same prompts of two types: the **I**nstruction-typed (I) prompt "What is this?" and the **C**ompletion-type (C) prompt "This is an object of ". 

Experimental results demonstrate PointLLM's superiority over both 2D and 3D MLLMs on ModelNet40 and

| Table 5. **Ablation on projection layers.** | | |
| --- | --- | --- |
| Hidden Dims. | 7B-Acc. | 13B-Acc. |
| N.A. | 50.63 | 52.62 |
| 1024 | 51.05 | 49.00 |
| **1024, 2048** | **52.82** | **53.39** |
| 1024, 2048, 4096 | 52.15 | 51.40 |

| Table 6. **Ablation on max pooling.** | | |
| --- | --- | --- |
| Pooling | Acc. | A100 GPU-Hours |
| 7B w/ | 48.72 | **34** |
| 7B w/o | **52.82** | 126 |
| 13B w/ | 51.10 | **56** |
| 13B w/o | **53.39** | 213 |

| Table 7. **Ablation on fine-tuning data.** | | | |
| --- | --- | --- | --- |
| Single | Multi. | Detailed | Accuracy |
| ✓ | | | 40.14 |
| ✓ | ✓ | | 45.79 |
| ✓ | ✓ | ✓ | **52.82** |

Objaverse datasets for various prompt types. Compared with 2D models, PointLLM offers direct point cloud engagement, showcasing enhanced 3D object comprehension over single-view images. This method effectively addresses challenges like occlusion and viewpoint variation, leveraging rich 3D geometry and appearance data from colored point clouds. PointLLM shows more consistent classification accuracy across different prompts than other 3D models, underlining its prompt robustness. Utilizing a pretrained point encoder and a large language model backbone, PointLLM efficiently translates point cloud data into descriptive natural language, capturing the object's identity.

The zero-shot performance on ModelNet40 further illustrates our model's aptitude for generalization. Even though ModelNet40 comprises point clouds unseen during training, PointLLM recognizes them using its pre-existing knowledge and perception abilities honed during our two-stage training. This adaptability to unseen domains and novel objects, without necessitating retraining, speaks to our model's robustness.

## 5.3. 3D Object Captioning

Tab. 3 displays the results of our 3D object captioning benchmark, averaged across objects. Each model was prompted with "Caption this 3D model in detail."

In Tab. 3 our models significantly outperform all baselines in key evaluation metrics for 3D object captioning, especially in human correctness score and GPT-4 evaluations. These scores reflect a model's ability to capture and articulate the intricate details of objects. Notably, both 7B and 13B PointLLM variants achieve the highest correctness scores, producing more accurate and detailed captions than other models, even rivaling human annotations. In addressing hallucination, a common MLLM challenge, our PointLLM-7B exhibits the lowest hallucination score and highest precision score, indicating its effectiveness in generating detailed, accurate captions with few false information. The Sentence-BERT and SimCSE results further confirm our model's capability in producing semantically rich captions closely aligned with the ground truth.

Interestingly, all 13B models, regardless of being 2D or 3D MLLMs, tend to create more hallucinated content than their 7B counterparts. This suggests that larger MLLMs may be more challenging to fine-tune for precision. The in-

vestigation of this trend in larger models and its underlying causes is an intriguing direction for future research.

We analyzed the human evaluation data to compare our models with baselines and human annotations. Win rates, calculated based on the correctness score for the 13B variants, are averaged across evaluations and presented in Fig. 3. PointLLM demonstrates notable performance, outperforming counterparts in over half of the test samples, including against human annotations (55% vs. 35%). This underscores PointLLM's ability to effectively capture and convey 3D object details, hinting at its potential for scalable, human-like captioning of 3D objects. Win rate comparison for the 7B variant and other scores are detailed in App. D.

**Limitations of traditional metrics.** Our evaluation also highlights the inadequacies of conventional metrics like BLEU-1, ROUGE-L, and METEOR, as shown in Tab. 3. While 3D-LLM scores high on these traditional metrics, often by producing shorter captions (averaging 20 words compared to others' 69+), this doesn't necessarily reflect true caption quality (witnessed by the human-evaluated metrics). Also, in the given example in Tab. 4, InstructBLIP-13B inaccurately describes a "Private jet" as a "jet engine," yet scores higher than PointLLM-13B, which correctly identifies it as an "airplane." This illustrates traditional metrics' failure to capture semantic accuracy, as they focus on n-gram overlap and may overlook the essence crucial in 3D object captioning. Consequently, we prioritize more comprehensive and reliable measures like human evaluation, GPT-4 evaluation, Sentence-BERT, and SimCSE for assessing our task.

## 5.4. Ablation Studies

In this section, we conduct ablation studies to explore various model design choices, including the number of projection layers and the use of max pooling for aggregating point cloud tokens. Additionally, we examine the impact of different data variations on the two-stage training process. The average accuracy of PointLLM on our generative classification benchmark is reported here.

**Projection layers.** While the alignment of information from different modalities to the text space using projection layers is a common practice [37, 72], the optimal number of layers required remains an open question. Our experiments,

Table 8. **Qualitative comparisons.** We show the classification and captioning results of models on ModelNet40[64] and Objaverse[11]. These samples show our PointLLM produces more accurate and detailed results than baselines and even human-annotated ground truths.

| Samples 1, 2 |  |  |
|---|---|---|
| Ground Truth | Laptop | A cartoon black monster like a dragon |
| Prompt | What is this? | Briefly caption this 3D model. |
| InstructBLIP[9] | symbol letter l | a black lizard with a sharp tooth in a dark room |
| LLaVA[37] | A small, grainy, black and white letter j. | A 3D model of a dark, menacing dragon. |
| 3D-LLM[25] | - | A black and white tiger with long legs, standing on its hind leg. |
| Point-Bind LLM[18] | This is a laptop computer. | The 3D model features a large, ornate gargoyle with a horned helmet, sitting on top of a building. |
| **PointLLM** | **The 3D model represents a notebook computer, typically a laptop.** | **The 3D model depicts a menacing black dragon, with its mouth opened wide revealing a row of sharp teeth.** |

ranging from 1 to 4 projection layers with different hidden dimensions, are detailed in Tab. 5. Results from both the 7B and 13B models indicate that 3 projection layers yield the best performance. This suggests that both an insufficient and an excessive number of layers can detrimentally affect performance. A balance in the number of layers is thus crucial for optimal model functionality.

**Max pooling.** Unlike sequential or grid-based text and image tokens, point cloud tokens are permutation invariant. Concatenating these tokens with text introduces unnecessary causal dependence and may not be optimal for feature fusion. Inspired by max pooling's symmetric properties [48], we experimented with aggregating point token information through max pooling before the projection layer. While this method didn't enhance performance as shown in Tab. 6, it greatly improved efficiency. Training time measured by 80G-A100 GPU-Hour reduced by about 75%. This underscores the challenge in developing efficient, point cloud-specific fusion mechanisms for MLLMs.

**Training data.** To determine the optimal quantity of data for feature alignment, we experimented with varying data volumes on our 7B PointLLM, maintaining constant interaction times by duplicating training epochs. Results in Fig. 4 suggest that increasing the data volume improves downstream performance, plateauing at around 600K samples. Further, as shown in Tab. 7, incorporating diverse types of instruction-following data during fine-tuning consistently yields performance improvements, underscoring the importance of our diverse instruction-following dataset.

### 5.5. Qualitative Results

Fig. 1 demonstrates PointLLM's ability to accurately perceive interior details of shoes and cars, overcoming occlusion and viewpoint challenges. This section offers a qualitative comparison of different 13B models in Tab. 8. Sam-
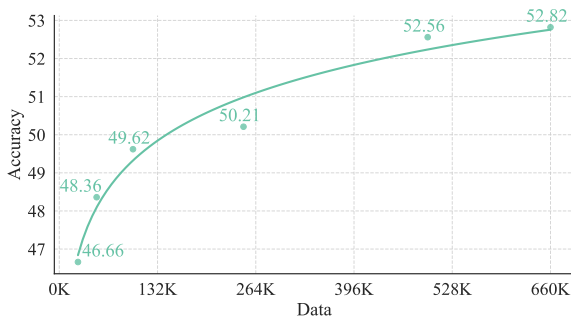


Figure 4. **Ablation on data for alignment.**

ple 1 from ModelNet40 shows a typical 2D MLLM failure: mistaking a laptop for letters due to depth perception issues inherent in single-view images. While multiple views could potentially alleviate this, they pose challenges in terms of optimal view selection and increased model complexity. Point clouds, however, directly provide object geometry, avoiding issues with depth, occlusion, or viewpoint. Sample 2 highlights PointLLM's capability to generate detailed, accurate captions, outperforming other models and even human annotations, while avoiding severe hallucinations. Additional qualitative results in App. E further illustrate the advantages of using point clouds for 3D understanding and PointLLM's superiority.

## 6. Conclusions and Future Directions

In this study, we introduce PointLLM, a novel MLLM for understanding 3D object point clouds. We also introduce a large-scale dataset and two benchmarks, complete with a comprehensive evaluation suite. All resources will be open-source for community use. Future work includes enhancing PointLLM for point cloud generation to support interactive 3D content creation, and utilizing PointLLM to automati-

cally generate high-quality 3D object captions for text-to-3D generation. We provide preliminary results related to text-to-3D generation in App. F.

**Acknowledgements.** We would like to acknowledge Xiangyu Yue for providing feedback about this paper, and thank Lihe Ding, Shaocong Dong, and Jiaming Han for their assistance with the experiments.

# Appendix

## A. Data Collection

**Instruction lists.** The 30 pre-defined instructions used to prompt the model to briefly and elaborately describe the objects are shown in Tab. 10 and Tab. 11 respectively. These prompts are generated with the assistance of GPT-4 and are coupled with captions to form our description-type data.

**Data generation with GPT-4.** In Tab. 12 we show an example of using GPT-4 for data generation as well as the system prompt of GPT-4. The input is one human-written caption provided by Cap3D[40] and the outputs are one expanded detailed caption, three single-round conversations, and one multi-round conversation. The system prompt is used for all samples, which guides the model to analyze existing captions based on the general knowledge of 3D objects and generate detailed captions, diverse Q&As, and logically connected multi-round conversations.

**Dataset distribution.** The comprehensive statistics of our newly compiled point-text instruction following dataset are detailed in Tab. 9. The dataset encompasses approximately 730K samples. Fig. 5 illustrates the length distributions of instructions and responses across various data types. In Fig. 6, we present word clouds (after removing generic words like "model", "object", etc.) and verb-noun pair distributions following [61] from our instruction following dataset, highlighting its extensive coverage of diverse topics such as color, shape, usage, material, and more.

Table 9. **Statistics of our point-text instruction following data.**

| Statistics | |
|---|---|
| Number of all samples | 731851 |
|   - brief-description type | 661577 |
|   - detailed-description type | 15055 |
|   - single-round type | 40122 |
|   - multi-round type | 15097 |
|     - multi-round responses | 45287 |
| Avg. len. of all instruction/responses (in words) | 11/17 |
|   - brief-description type | 11/15 |
|   - detailed-description type | 9/82 |
|   - single-round type | 10/15 |
|   - multi-round type | 10/21 |

## B. Benchmarks and Evaluation

### B.1. Discussions about Close-Set Classification

Initially, we consider formatting the close-set zero-shot classification task on ModelNet40[64] as a multiple-choice problem, including indexed candidate category names in the prompt, and prompting our model to select one of the 40 categories given the point cloud as input. However, since our model is not designed for multiple-choice problems but for real-world usage where it can generate any word or phrase as output, we cannot directly parse its response for evaluation. Therefore, we use ChatGPT as a post-processor to select one of the ModelNet40 categories based on the model's answer. In the meantime, we find that including category names in the prompt results in meaningless responses from InstructBLIP[9], which is the model we compare with, making meaningful comparisons challenging. Consequently, we opt for a more generalized prompt, without including the candidate lists in the prompt. This allows us to make balanced comparisons.

Including candidate lists in the prompt, we also tried to calculate the conditional probability of different options given the model's output following [62], but this method did not work well for our model. As our instruction-following training data lacks such scenarios where it's needed to choose from a fixed set of options, our model always produces very low probabilities on these options with biased results. For example, among the options "00" to "39", our model predicts very low probabilities and among these low probabilities, "00" and "39" are the highest most of the time, which leads to biased predictions. Therefore, we choose to use general prompts and utilize ChatGPT for post-processing. This approach more accurately reflects real-world scenarios where the model is expected to provide natural, free-form responses to diverse, unstructured questions.

## B.2. GPT Evaluation Prompts

**Close-set zero-shot classification.** In this task, we use ChatGPT to post-process the model output by selecting the most probable class index from the 40 ModelNet40 categories. The process is detailed in Tab. 13, where {candidate_lists} refers to the ModelNet40 category list, and {model_output} refers to the model's response. ChatGPT is required to directly output the category index, category name, and a short reason for the choice. If the description doesn't clearly refer to any one of the categories, ChatGPT must make an educated guess based on the information provided. If ChatGPT cannot infer, then "-1" is returned and a random index will be chosen as the model's classification prediction. We do not use a system prompt for ChatGPT but directly input the prompt.

**Open-vocabulary classification.** In this task, we use GPT-4 as an evaluator to classify whether the model's response and the human caption are referring to the same object type. The process is outlined in Tab. 14, where {ground_truth} and {model_output} refer to the human caption and the model's response. We do not require the model's response to match exactly with the human caption, as long as it conveys the same object type. We also directly input the prompt for GPT-4 instead of using a system prompt.

**Object captioning.** In this task, we utilize GPT-4 as an evaluator to assess model-generated captions against human-generated captions (ground truth) of 3D models. GPT-4 is tasked with identifying aspects mentioned in the human caption and calculating the percentage of these aspects that are either correctly mentioned or partially matched in the model's caption on a scale of 0 to 100, with each aspect contributing equally to the score. The evaluation process is detailed in Tab. 15, where {ground_truth} refers to the human caption, and {model_output} refers to the model's response.

## B.3. Human Verification of GPT Evaluation

To verify the effectiveness of using GPT models for evaluation, the first author manually checks the evaluation results of ChatGPT and GPT-4.

In the close-set classification task on ModelNet40, the author finds the following:

1. ChatGPT consistently outputs in the desired format, selecting the category or "-1" and providing a reason.
2. When the model output clearly refers to or hints at a category with salient information regarding one of the candidate categories, ChatGPT can accurately identify the corresponding category based on the model's output, showing a high degree of consistency with human-selected options. False negatives or false positives are rare in these cases.
3. If the model output is ambiguous, ChatGPT's selection appears random, aligning with our expectations for han-

dling such cases in classification tasks, because when the model encounters uncertainty or lacks confidence in its identification, random guessing is permissible.

For open-vocabulary classification and object captioning tasks on Objaverse, the author finds that ChatGPT underperforms in identifying the same object concept, acting as a strict judge, and producing more false negatives in classification. It often considers two words or phrases not to refer to the same object type, even when they do. In contrast, GPT-4 demonstrates accurate recognition. After reviewing 50 samples of classification results, the first author has 100% consistency with GPT-4's evaluations. As a result, we opt to use GPT-4 for the open-vocabulary and object captioning tasks on Objaverse. Examples of GPT evaluation can be found in Tab. 13, Tab. 14, and Tab. 15.

## B.4. Human Scoring Criteria

Human evaluators were employed to assess captions in the object captioning benchmark. Outputs for the same object from various models were grouped and randomly shuffled, and evaluators independently scored these captions while manually inspecting objects in the Objaverse using the official explorer at https://objaverse.allenai.org/explore.

**Scoring criteria.** The evaluation process involved assigning correctness scores and hallucination scores following these guidelines:

1. **Correctness score.**
   - Each distinct correct attribute in a model output (*e.g.*, category, color, shape, usage, material) was awarded one point. For example, a black tire correctly identified as a tire and being black would receive two points.
   - Partial correctness was graded on a scale of 0 to 1, depending on the degree of accuracy. For instance, if a model output described "a cartoon figure" but the object was specifically a cartoon horse, it would be awarded 0.5 points.
2. **Hallucination score.**
   - Hallucination points were assigned for each incorrect detail in the model output, mirroring the correctness scoring mechanism. For instance, if the model incorrectly described two yellow tires instead of four black ones, it would incur two hallucination points, one for color and one for number.
   - Repetitive inaccuracies based on one attribute were not subject to multiple penalties. As an example, erroneously mentioning a black tire when no tire existed would lead to only one hallucination point.
   - Penalties were also applied for content that was irrelevant to the object description.
3. **General considerations.**
   - Generic terms like '3D model' or 'image' were disregarded, as were references to black backgrounds or environmental colors, and viewpoints.

- Elements that were indeterminable as either correct or incorrect were not considered in the scoring.
- Within each group of evaluations, a range of scores should be established to differentiate between high and low-quality captions. After the initial scoring, a final review adjusted the scores to ensure a clear distinction between better and worse captions.

**Precision score.** We also calculate the precision score as the proportion of correct information within the model-generated content as follows:

$$\text{Precision} = \frac{\text{C. Score}}{\text{C. Score} + \text{H. Score}} \times 100\% \qquad (3)$$

where **C. Score** refers to the correctness score and **H. Score** refers to the hallucination score. The precision score in the main paper is reported after summing all **C. Score** and **H. Score** of all samples for robust evaluation.

## C. Implementation and Training

**Implementation details.** We use ULIP-2[66] to pre-train our point cloud encoder (Point-BERT[68]). ULIP-2 is a method for aligning the latent space of the point cloud encoder to that of CLIP[49] through contrastive learning, endowing the encoder with a strong zero-shot capability for 3D object recognition. As the original implementation of ULIP-2 only supports point clouds with spatial coordinates (xyz), we re-train Point-BERT from scratch with color information (xyzrgb), following the same procedure outlined in the ULIP-2 paper. For training Point-BERT, we employ ViT-L/14 trained on DataComp-1B[13] for 12.8B steps and batch size 90k (denoted as 'ViT-L/14-datacomp_xl_s13b_b90k') from OpenCLIP[29]. We use point clouds from Cap3D[40], which contains 660K objects. We filter out 3000 objects from this dataset and reserve them for future testing. These 3000 objects are not used during any stage of the entire model training and the 200 objects utilized for our benchmarks are part of these 3000 unseen objects to prevent information leakage.

**Training details.** All training are conducted on 8 × 80G A100 GPUs with BF16 data type, leveraging flash-attention [10], the AdamW [39] optimizer, and a cosine learning rate scheduler. For the feature alignment stage, we train our model for 3 epochs with a batch size of 128 and a learning rate of 2e-3. For the instruction tuning stage, we train our model for 3 epochs with a batch size of 32 and a learning rate of 2e-5. For efficiency, the 7B model completes the feature alignment and instruction tuning stages in approximately 13.3 and 2.5 hours, respectively, while the 13B model takes around 22.3 and 4.3 hours for the same stages.

## D. Win Rate Comparison

In Fig. 7, we present the win rate comparisons of the 7B and 13B model variants across different scores. Note that due to truncation errors in plotting figures, some cumulative rates (win, tie, lose) may not sum to exactly 100%.

The win rate analysis reveals that both the 7B and 13B models not only outperform baselines but also surpass human annotators in terms of correctness scores by a significant margin. This superior performance in correctly identifying object attributes underscores the models' advanced understanding and processing capabilities of 3D objects.

Regarding hallucination, our models exhibit a marked improvement over the baselines, reflecting their better ability to avoid generating incorrect or fabricated details about the objects. This is further corroborated by the precision scores, where our models demonstrate a higher ratio of correct information in their outputs compared to other models. Compared with human annotators, our models show comparable performance in about 50% of samples in terms of hallucination and precision.

However, it must be acknowledged that there is still room for improvement in reducing hallucination rates to match the levels achieved by human annotators. Striving towards the precision demonstrated by human evaluations remains a target for future enhancements.

## E. Qualitative Results

In this section, we provide the qualitative results from different datasets of the 13B models for comparison. All samples used were unseen by our models during training.

**Results on ModelNet40.** Tab. 16 illustrates the classification results from different models on the ModelNet40 dataset. These examples highlight the inherent limitations of image-based models, which depend on suitable views for accurate object identification. Notable challenges include the failure to recognize the guitar in Sample 3, the monitor in Sample 6, and the ambiguity in depth perception leading to misclassifications, such as confusing a chair with a bed in Sample 2, and a bathtub with a bowl in Sample 5. In contrast, PointLLM bypasses these challenges by using point clouds, which provide direct access to object geometry without concerns over ambiguous depth, occlusion, or viewpoint.

Moreover, there are evident hallucination issues with other models. For instance, Point-Bind LLM[18] erroneously describes a person lying on the couch in Sample 1 and someone holding a wine bottle in Sample 4. In comparison, our 13B model consistently provides accurate and realistic classifications. Its superior performance, devoid of such hallucinatory inaccuracies, underscores the model's advanced comprehension of 3D structures and its effectiveness in handling diverse object types.

**Results on Objaverse.** Due to the limited capability of Point-Bind LLM in producing meaningful outcomes, it has been excluded from our comparative analysis on Objaverse. It is noteworthy that InstructBLIP also occasionally yields

nonsensical results as in Sample 2 of Tab. 17. As depicted by the results, PointLLM consistently generates captions that are both more accurate and detailed compared to other baselines and human annotators. For example, in Sample 2, PointLLM accurately describes the golden brown eyes of an insect, a detail overlooked by human annotators who provide only a generic description, and completely missed by other models failing to identify the object type correctly. Similarly, in Samples 3 and 4, PointLLM offers elaborate descriptions encompassing shape and color without errors, in stark contrast to the simplistic captions from human annotators and erroneous information from other baselines.

**Dialogues.** Fig. 8 showcases dialogues between PointLLM and a human user, which reveal PointLLM's capacity to understand point clouds' shapes, appearances, functionalities, and more. Notably, our PointLLM is unaffected by occlusion, capable of discerning the car's internal two-seat structure and identifying a logo on the back of a shoe, tasks challenging for image inputs. Furthermore, our model engages with human instructions using common sense and avoids biases, as seen in its refusal to declare a 'best' shoe brand. Collectively, these samples validate PointLLM's proficiency in understanding point clouds and responding to human instructions both accurately and effectively.

## F. Text-to-3D Generation

The burgeoning interest in text-to-X generation tasks [5, 31, 41, 52] has led to significant advancements. Notably, [5] demonstrated that text-to-image generation models benefit greatly from training on highly descriptive, generated captions. Leveraging PointLLM's capability to generate detailed and accurate captions for 3D models, we explore its potential in enhancing text-to-3D generation models.

We employed PointLLM-13B to generate captions for Objaverse[11] objects with LVIS[19] labels. The prompt "Describe this 3D model in detail and accurately." was used for caption generation. We trained the text-to-3D generation model from [2] using captions generated by our PointLLM and the Cap3D[40] respectively for comparison.
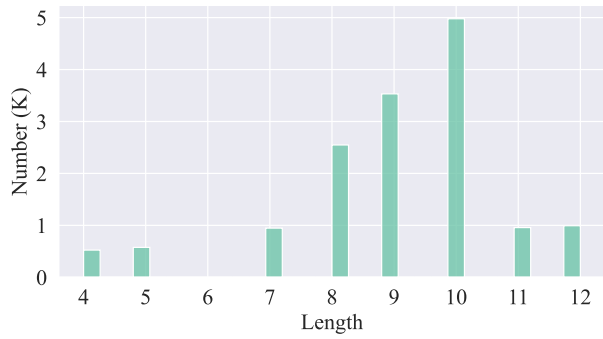
Qualitative comparisons of the generation results, as illustrated in Fig. 9, reveal that the model trained with PointLLM-generated captions generates objects more closely aligned with text prompts and exhibits more precise detailing. This highlights the advantages of utilizing detailed and accurate captions from our model for text-to-3D generation tasks, pointing towards improved fidelity and coherence in generated 3D objects.
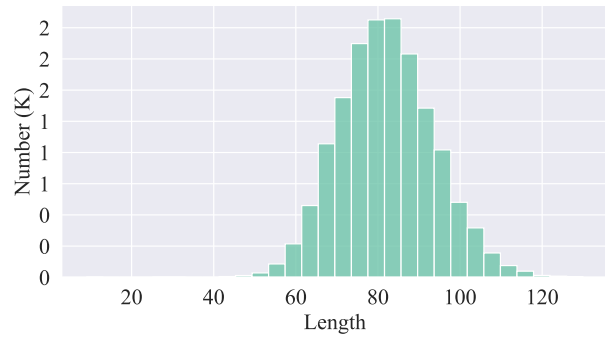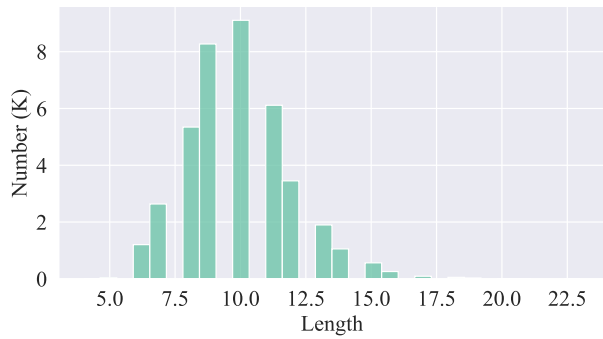
(a) **Brief description-instruction.**
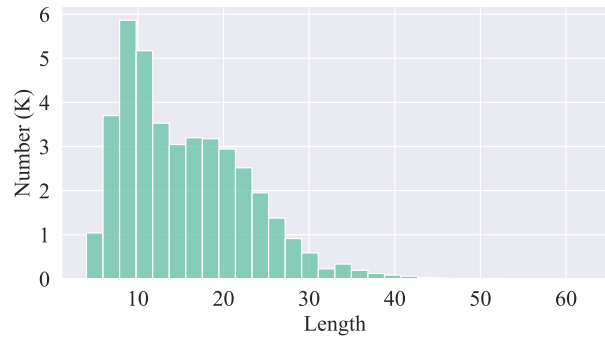
(b) **Brief description-response.**

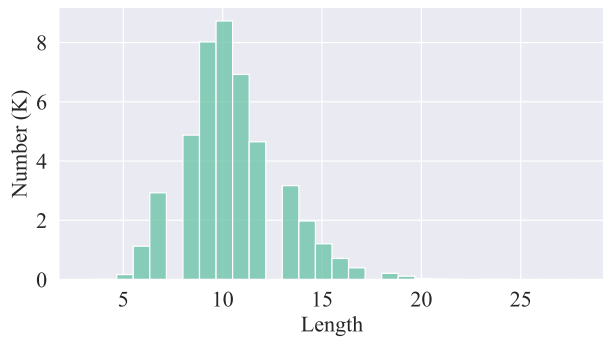(c) **Detailed description-instruction.**
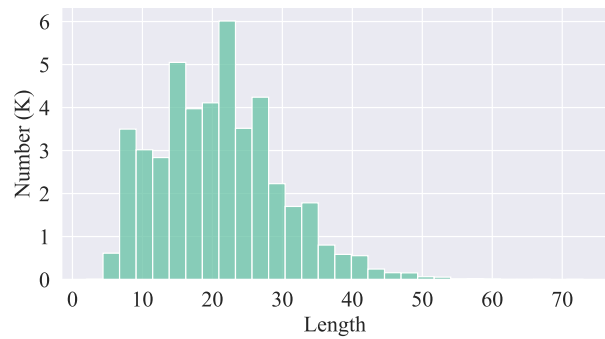
(d) **Detailed description-response.**

(e) **Single round-instruction.**

(f) **Single round-response.**

(g) **Multi round-instruction.**

(h) **Multi round-response.**

Figure 5. **Length distributions of instructions and responses of different types of point-text instruction following data.**

(a) **Brief description-instruction.**

(b) **Brief description-response.**

(c) **Detailed description-instruction.**

(d) **Detailed description-response.**

(e) **Single round-instruction.**

(f) **Single round-response.**

(g) **Multi round-instruction.**

(h) **Multi round-response.**

Figure 6. **Word distributions of instructions and responses of different types of point-text instruction following data.**

Table 10. **The instruction list for brief descriptions.** An instruction from the list is randomly selected and coupled with a human-written caption from Cap3D[40] to form a brief-description instruction following sample.

- Summarize the 3D point cloud object briefly.
- What kind of object is depicted by this point cloud?
- Provide a short explanation of this 3D structure.
- What does this collection of points represent?
- Offer a succinct summary of this 3D object.
- Can you give a brief overview of this point cloud?
- Characterize the object this point cloud is illustrating.
- Share a brief interpretation of this 3D point cloud.
- Provide an outline of this 3D shape's characteristics.
- What object is this point cloud rendering?
- Deliver a quick description of the object represented here.
- How would you describe the 3D form shown in this point cloud?
- What is the nature of the object this point cloud is representing?
- Present a compact account of this 3D object's key features.
- What can you infer about the object from this point cloud?
- Offer a clear and concise description of this point cloud object.
- How would you summarize this 3D data set?
- Give a brief explanation of the object that this cloud of points forms.
- What kind of structure does this 3D point cloud depict?
- Could you delineate the form indicated by this point cloud?
- Express in brief, what this point cloud is representing.
- Give a quick overview of the object represented by this 3D cloud.
- Convey a summary of the 3D structure represented in this point cloud.
- What kind of object is illustrated by this collection of points?
- Describe the object that this point cloud forms.
- How would you interpret this 3D point cloud?
- Can you briefly outline the shape represented by these points?
- Give a concise interpretation of the 3D data presented here.
- Explain the object this point cloud depicts succinctly.
- Offer a summary of the 3D object illustrated by this cloud.

Table 11. **The instruction list for detailed descriptions.** An instruction from the list is randomly selected and coupled with a GPT-4 generated caption to form a detailed-description instruction following sample.

- Can you tell me more about this?
- What does this represent?
- Can you describe this in more detail?
- I'm interested in this, can you explain?
- What is this object made of?
- Could you provide more info about this?
- What exactly am I looking at here?
- What is this?
- Could you describe the detailed structure of this?
- This looks interesting, can you expand on it?
- Can you explain more about this form?
- What can you tell me about the shape of this object?
- Could you delve deeper into this?
- I want to know more about this, can you help?
- Can you walk me through the details of this object?
- Can you provide a comprehensive account of this object?
- Offer a detailed interpretation of this point cloud.
- Please elucidate on the characteristics of this form.
- Could you provide an in-depth description of this structure?
- What does this cloud represent in its entirety?
- Elaborate on the details of this point cloud, please.
- Kindly furnish me with more information about this object.
- Please expand on the intricate structure of this form.
- Provide a meticulous explanation of what these points represent.
- I request a detailed breakdown of this structure.
- Give a thorough rundown of this point cloud.
- Can you offer a complete analysis of this object?
- I would like a comprehensive explanation of this form.
- Please detail the specific features of this point cloud.
- Could you elaborate extensively on what this represents?

Table 12. **An example of data generation with GPT-4.** The input consists of a human-written caption provided by Cap3D[40], guided by a system prompt to analyze the existing caption based on the general knowledge of 3D objects. The outputs include an expanded detailed caption, three single-round conversations, and one multi-round conversation.



| Uid | 4887da0aab51406dab3c5cb69ec82404 |
|---|---|
| System Prompt | Analyze the 3D object model from the given caption:<br>1. Write a new detailed caption by describing it in 50-100 words, identifying its type, appearance, and any other inferred information such as functionalities, usages, daily-life knowledge, etc. Exclude any uncertain details.<br>2. Generate 3 single-round Q&As regarding diverse aspects of the object based on the provided captions and your new captions.<br>3. Construct 1 set of 3-round Q&As, ensuring logical relevance between the questions and answers. Questions should be different from those in `"2"`.<br><br>Format your response as:<br>```json<br>{<br>  "caption": "description",<br>  "single_conversation": [{"Q": "Q", "A": "A"} x 3],<br>  "multi_conversation": [{"Q1": "Q", "A1": "A", "Q2": "Q", "A2": "A", "Q3": "Q", "A3": "A"} x 1]<br>} |
| Input | A female doll with long back cream hair and with black high boots. |
| Output | {<br>  "caption": "The 3D object model represents a detailed doll figure of a female character with long, cream-colored hair cascading down her back. It is dressed in fashionable attire with black high boots being one of its highlighted features. The doll, resembling those used in play or as collectible items, exhibits a high degree of intricate detailing, showcasing its fashionable outfit and engaging features, making it appealing for users, especially ones intrigued by doll collection or imaginative play.",<br>  "single_conversation": [<br>    {"Q": "What color is the doll's hair?", "A": "The doll's hair is cream-colored."},<br>    { "Q": "What is specific about the doll's footwear?", "A": "The specific footwear feature of the doll is its black high boots."},<br>    {"Q": "What are possible usages of this doll?", "A": "Possible usages for this doll might include being part of a toy collection, used for imaginative play, or as a detailed model for hobbyists or artists."},<br>  ]<br>  "multi_conversation": [<br>    {"Q1": "What makes this doll distinguishable?", "A1": "The doll is distinguishable by its long cream-colored hair cascading down its back and its outfit featuring black high boots.", "Q2": "Could this doll be part of a collection or series?", "A2": "Yes, based on the intricate detailing and appearance, it's possible that the doll could be part of a collection or series.", "Q3": "Who would be the likely audience or users of this doll?", "A3": "The likely users of this doll would be individuals who enjoy doll collecting, children or adults for imaginative play, or even hobb yists and artists who may use it as a model."}<br>  ]<br>} |

Table 13. **The prompt and examples of ChatGPT in close-set zero-shot classification.** ChatGPT post-processes the model output and selects the most probable class index from the available 40 categories, even if the description is vague or ambiguous. The blue placeholders {candidate_lists} and {model_output} refer to the ModelNet40 category list and the model's response, respectively.
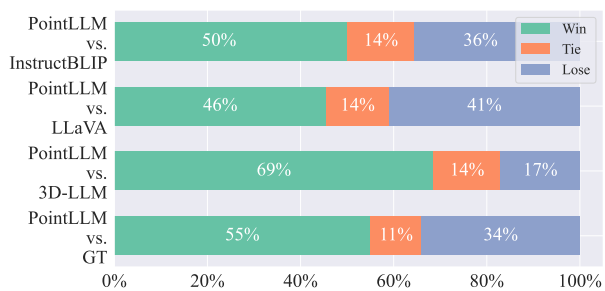
| | |
|---|---|
| Prompt | Given the following free-form description of a 3D object, please determine the most probable class index from the following 40 available categories, even if the description doesn't clearly refer to any one of them. Make your best-educated guess based on the information provided. If the description already contains a valid index, then the index should be selected. If it contains more than one valid index, then randomly select one index (specify your reason). If there is no valid index and it cannot be inferred from the information, return "-1#NA#Cannot infer". |
| | Categories: |
| | 0: airplane |
| | 1: bathtub |
| | {candidate_lists} |
| | 38: wardrobe |
| | 39: xbox |
| | Reply with the format of "index#class#short reason (no more than 10 words)". |
| | |
| | Examples: |
| | Input: This is a 3D object model of a cartoon white truck. |
| | Output: 7#car#Closest match to "car" in categories. |
| | |
| | Input: A green leaf in a flower pot. |
| | Output: 26#plant#The primary subject "leaf" directly indicates a plant. |
| | |
| | Input: It's difficult to determine the exact type of this object due to insufficient details. But it seems to be like a piece of furniture. |
| | Output: 33#table#Randomly select one kind of furniture from the list. |
| | |
| | Input: I cannot determine the specific type of the object without additional information or context. |
| | Output: -1#NA#Cannot infer. |
| | |
| | Now analyze the following: |
| | Input: {model_output} |
| | Output: |
| Example 1 | Input: This 3D model represents a war plane, typically used for combat and military purposes. It is characterized by a large, dominating hue of military gray. The detailing on the model suggests it's a fighter jet, hinting at its usage in air combats. Its posture seems ready for flight, showcasing a typical attribute of war planes - preparedness for battle. This could be used in game design or other virtual reality applications that require detailed, lifelike renderings of objects. |
| | Output: 0#airplane#The description explicitly mentions "war plane" and "fighter jet". |
| Example 2 | Input: The 3D model object is a tall, cylindrical-shaped clay container with a flat, smooth top and a slightly wider base. The cylindrical form and the structure suggest that it's possibly used for storing liquids or semi-liquids due to its wide top and flat base allowing easy pouring and stability. Its material and structure indicate it could be a common household item, like a vase or a jug. |
| | Output: 37#vase#The description includes key characteristics of a vase (tall, cylindrical shape, wide top, flat base). |
| Example 3 | Input: This is a 3D model of a sleek, black rectangular unit with its sides elegantly curved. On one of its sides, there is a distinctive circle inscribed, adding a unique character to the model. The object's form and design suggest that it can be a piece of technology, like a gadget or a home appliance, with the circle possibly indicating a button or a sensor. Its usage, however, cannot be ascertained without specific context. |
| | Output: 22#monitor#The description suggests a sleek, black rectangular unit which can be a monitor. The inscribed circle may indicate a power button or a sensor on the monitor. |

Table 14. **The prompt and examples of GPT-4 in open-vocabulary classification.** GPT-4 needs to analyze two sentences to determine if they refer to the same general object or concept, focusing on the type of object, not attributes such as color, size, or shape. The placeholders {ground_truth} and {model_output} refer to the human caption and the model's response, respectively.
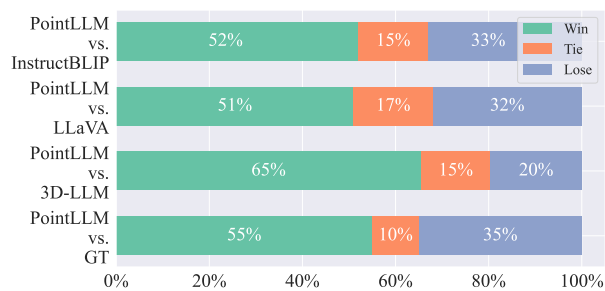
| | |
|---|---|
| Prompt | Analyze two sentences and determine if they're referring to the same general object or concept, focusing on the type of object, not attributes such as color, size, or shape. Respond with `"T"`if they refer to the same thing and `"F"`if not. Also, provide a brief rationale (no more than 20 words) for your judgment. <br><br> Example: <br> Input: 1. Spiral staircase that goes from a ground floor. 2. This is a 3D model of wooden stairs in light brown <br> Output: T#Both refer to a staircase. <br><br> Now, analyze the following: <br> Input: 1. {ground_truth} 2. {model_output} <br> Output: |
| Example 1 | Input: 1. A black and brown colored gun. 2. The 3D object is a representation of a futuristic, high-tech gun crafted from a glossy black material. Distinctive features include its metallic handrail, giving an impression of a robust mechanized design. The gun, possibly used in a sci-fi or futuristic setting, denotes advanced technology and might include functionalities such as voice recognition, aiming systems, or biometric triggers. <br> Output: T#Both refer to a gun. |
| Example 2 | Input: 1. A yellow and white fish with black stripes and fins. 2. This is a 3D model of a vibrant, polka-dotted toy fish that is predominantly orange on the body, shifting to white on the belly. The toy has dark brown spots that enhance its appearance, potentially mimicking the natural patterns found on real-life fish. It's an ideal object for educational purposes, helping to introduce children to marine life, as well as serving as a playful item in a playroom or nursery. <br> Output: T#Both refer to a fish. |
| Example 3 | Input: 1. A white cartoon scorpion with eight legs. 2. This is a 3D object model representing a cartoon version of a rare type of spider. The entire model is rendered in white, which highlights its unique and exaggerated characteristics such as multiple legs and a funnel-like body. Its cartoonish appeal makes it more appealing to a younger audience, and it could possibly be used in animations or educational materials to teach children about spiders in a less intimidating way. <br> Output: F#One is a scorpion and the other is a spider. |

Table 15. **The prompt and examples of GPT-4 in object captioning.** GPT-4 evaluates the model's response by identifying aspects mentioned in the human caption and calculating the percentage of aspects that are correctly or partially matched in the model's caption. The placeholders {ground_truth} and {model_output} refer to the human caption and the model's response, respectively.
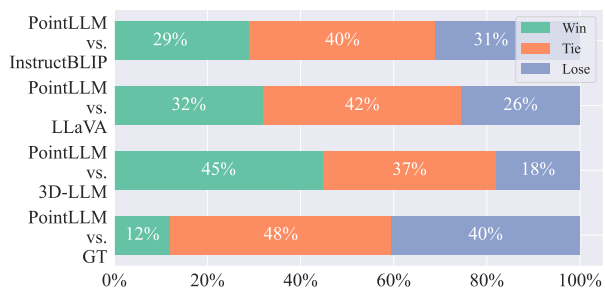
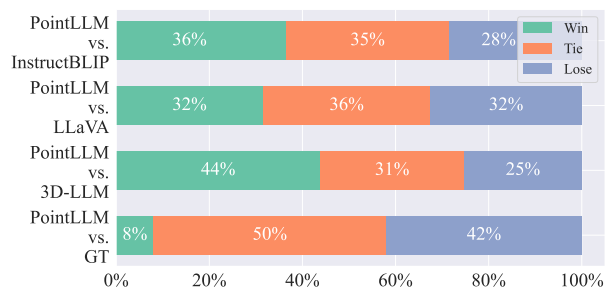| | |
|---|---|
| Prompt | Evaluate a model-generated caption against a human-generated caption (ground truth) for a 3D model. Identify the aspects mentioned in the human caption and calculate the percentage of these aspects correctly mentioned or partially matched in the model caption. Score from 0 to 100, where each aspect contributes equally to the score. Consider similar concepts for a partial score.<br><br>Provide your score (0-100) and a short justification (less than 15 words) in the format of `"score#reason"`<br><br>Example:<br>Human: A white brown skeleton<br>Model: This is a 3D model of a small, cartoon-like robot. It has a spherical body and is covered in a layer of white dust.<br>Output: 50#mention white; skeleton and robot have similar appearence.<br><br>Now score the following:<br>Human: {ground_truth}<br>Model: {model_output}<br>Output: |
| Example 1 | Human: A white presentation of a planted trees plantation region.<br>Model: The model is a representation of a table saw with a simple, four-legged design.<br>Output: 0#No aspects from the human caption are mentioned in the model's. |
| Example 2 | Human: Private jet<br>Model: This is a 3D model of a cartoon-style airplane.<br>Output: 50#mentioned airplane, but missed private jet aspect. |
| Example 3 | Human: A cartoon look like a tree in 3d<br>Model: The 3D model is an animated depiction of a tree with gray bark, characterized by its lifelike details and textured surface that mimic the realistic grain of bark.<br>Output: 100#mentioned tree, 3D and artistic/cartoonish aspect. |

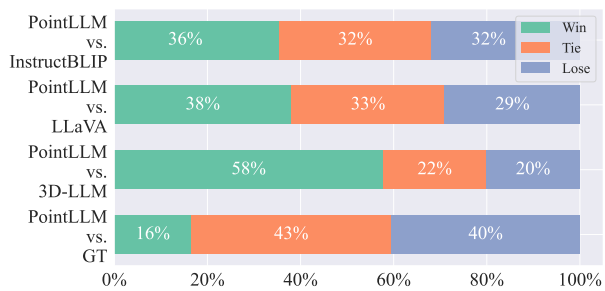Figure 7. **Win rate comparisons of 7B and 13B models across various scores.** Our models surpass baselines and human annotators in correctness scores, while exhibiting lower hallucination rates and superior precision compared to various baselines.

Table 16. **Qualitative results on ModelNet40.** The first image in each sample serves as the input for image-based models, with additional point cloud views provided for reference. PointLLM consistently and accurately identifies object types, whereas other models struggle with correct identification or produce hallucinated, incorrect content.

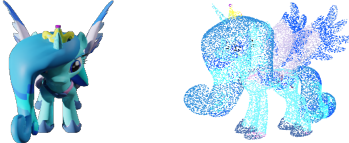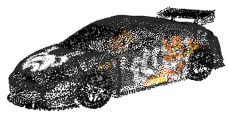| Samples 1, 2 |  |  |
|---|---|---|
| Ground Truth | Sofa | Bed |
| Prompt | What is this? | This is an object of |
| InstructBLIP[9] | sofa | person cutting meat into small pieces with a large knife |
| LLaVA[37] | The image is a black and white drawing of a couch. | The image is a black and white depiction of a chair. |
| Point-Bind LLM[18] | This is a drawing of a person laying on a couch, with a dog nearby. | This is an object of art. |
| **PointLLM** | **This 3D object is a comfortable, black, three-seater sofa.** | **This is a model of a neatly made bed, primarily composed of a solid black mattress.** |

| Samples 3, 4 |  |  |
|---|---|---|
| Ground Truth | Guitar | Bottle |
| Prompt | What is this? | This is an object of |
| InstructBLIP[9] | telescope | wine bottle |
| LLaVA[37] | The image is a close-up view of a black ink pen sitting on a white surface. | This image features a black and white depiction of a bottle. |
| Point-Bind LLM[18] | This is a bird flying in the sky. | This is an object of a person holding a bottle of wine. |
| **PointLLM** | **This 3D model represents a black electric guitar equipped with a distinctive headstock.** | **This is a 3D model of a bottle, possibly intended for liquids.** |

| Samples 5, 6 |  |  |
|---|---|---|
| Ground Truth | Bathtub | Monitor |
| Prompt | What is this? | This is an object of |
| InstructBLIP[9] | table | the foot |
| LLaVA[37] | The image features a black and white, close-up view of a grater, which is an object used for grating cheese or other food items. | The image features a small table with a black base and a white or light-colored top. |
| Point-Bind LLM[18] | This is a bowl or a vase, which is oval-shaped and has a smooth surface. | This is an object of art, a sculpture, or a decorative piece. |
| **PointLLM** | **The 3D object model represents a realistic, dark grey bathtub.** | **This is a 3D model of a desktop computer monitor.** |

Table 17. **Qualitative results on Objaverse.** The first image in each sample serves as the input for image-based models, with additional point cloud views provided for reference. PointLLM can generates more accurate and detailed captions for 3D objects than other baseline models and human annotators.



| Samples 1, 2 | | |
|---|---|---|
| Uid | 0ea33b6617174530b97d6b7a92c275fb | de8ec2a724f14fc4b54624512f80f13e |
| Prompt | What is this? | This is an object of |
| Human | A cartoon green and red like a fruit | A black insect |
| InstructBLIP[9] | an appleavatar 3d model | 26 |
| LLaVA[37] | This is a computer-generated image (CGI) or a 3D model of an apple with a leaf on it. | This is an image of a large, close-up, and dark-colored insect, possibly a big cockroach, with long antennae. |
| 3D-LLM[25] | A 3D model of a red apple. | A small, black spider with a long tail. |
| **PointLLM** | **This is a 3D model of a unique apple, distinctively adorned with a single, vibrant green leaf at the top.** | **This 3D model depicts a realistic, jet-black insect with a pair of striking, golden brown eyes.** |



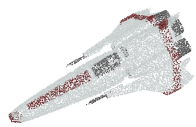| Samples 3, 4 | | |
|---|---|---|
| Uid | 059a7936ed89419ba9eae3153753ae86 | 82413f74efb84bc3b872a7a8f26cca81 |
| Prompt | Caption this 3D model in detail. | Caption this 3D model in detail. |
| Human | A blue cartoon | A white fish with orange/red spots |
| InstructBLIP[9] | There is a cute blue and white pony in a still, looking cheerful and cheerful. It's standing next to a larger blue-colored object, potentially a barrier or structure, which provides the backdrop for this scene. The overall image is designed to create a fun, whimsical setting for the pony. | There is a fish hanging on a string in front of a black background. The fish is swimming in the middle of the image and is a little bit up from the bottom of the screen. The fish has a red spots on the left and right sides, and is positioned slightly towards the top of the screen. |
| LLaVA[37] | The image features a 3D model of a cute, colorful cartoon unicorn with a blue and white body and a hint of pink. The unicorn has a magical aura and is set against a black background. Its long, flowing mane and tail add to its whimsical charm. The unicorn appears to be in motion, possibly flying or prancing, as it captures the viewer's imagination. | The image features a 3D model of a fish floating in a dark, empty space. The fish is positioned towards the left side of the image and appears to be swimming upwards. The fish is large and occupies a significant portion of the image, giving the impression of a life-like, computer-generated representation. The dark background and lack of other objects in the scene emphasize the fish's presence and create a sense of depth and dimension. |
| 3D-LLM[25] | A 3D model of a purple and pink pony with long hair, standing on its hind legs. | A small fish with a long tail is visible in the image. It has no distinguishing features or markings, and its color appears to be brown. |
| **PointLLM** | **A 3D model of a vibrant blue cartoon pony, characterized by its large white eyes, two wings, and a distinctive horn. The creature's overall design exudes a playful and friendly aura, making it appropriate for children's animations or toy production. The noticeably large white eyes and the unique wing and horn features endow the pony with an element of fantasy, while its blue body color gives it an overall appealing look.** | **The 3D model is of a long fish, illustrating a striking color combination of red, yellow, and grey. It's characterized by a slender tail. The yellow color is mostly seen on the body of the fish, fading into a subtler, grayish hue towards the fins and tail. The model depicts a typical long fish species, often found in the ocean, and can be used in a variety of digital applications, like animation, games, or virtual reality.** |

23

**What can we know from this?**

This is a 3D model of a sleek and stylish black racing car. The car sports a dark black body, complimented by black tinted windows and matching black tires. The design is optimized for high-speed performance, with features like a low and wide body to improve aerodynamics. The car likely has various functionalities geared towards professional racing, such as a powerful engine, detailed instrumentation, and high-performance brakes.

**How many seats are there in this car model?**

This car model has two seats.

**How many doors does this car have?**

This car model has two doors.



**Tell me about this.**

The 3D model represents a detailed spacecraft, which possesses a sleek aerodynamic shape for efficient space travel. The spacecraft is primarily colored white, providing a clean and professional look. Accentuating its design are red details, likely placed on key components such as the engine area, wing tips, or as lines running along the body. The spacecraft might be used for exploration, research, or transportation in outer space, based on its appearance and common uses of such vehicles.

**What physical principles are applied in the ascent of a spacecraft?**

The ascent of a spacecraft is achieved through the application of aerodynamics and the propulsion provided by the engines.



**What is this?**

The 3D object model represents a pair of stylish, yellow sneakers with white laces. These shoes feature a distinct white logo displayed prominently above the back of the heel. The sneakers overall design suggests a modern, trendy style, making them suitable for casual wear or sport activities. Their bright yellow color combined with the clean white detailing gives them a fresh, appealing look.

**List some brands that can make similar items.**

Nike, Adidas, and Converse are some of the brands that could make a similar item.

**Which brand is the best?**

The best brand is subjective and depends on personal preferences and the quality of the product.



**What is this used for?**

It is used for architectural visualization, interior design modeling, or virtual reality experiences.

**What else can we know about this model?**

The model is structured, with the stairs standing out as a distinct element, and the background being a brown, possibly wooden, floor.

**Why do we need a staircase?**

A staircase allows for vertical movement between floors within a building. In a residential setting, it is a common feature in houses.

Figure 8. **Dialogues between PointLLM and a human user.** The dialogues show PointLLM's ability to understand point clouds' shapes, appearances, functionalities, etc. PointLLM's resilience to occlusion enables it to perceive detailed features such as the interior structure of a car or a logo on a shoe, which are challenging for image inputs. Additionally, PointLLM demonstrates abilities to respond to human instructions with common sense, avoiding biases.
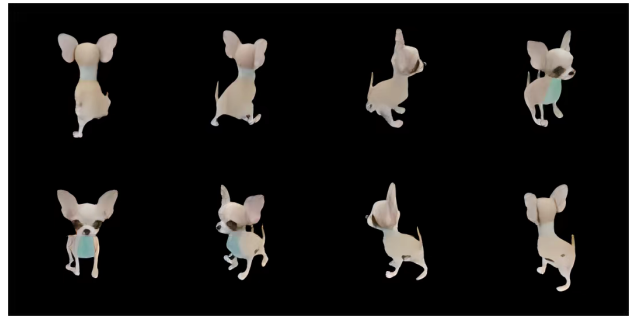
(a) **PointLLM.**

(b) **Cap3D.**

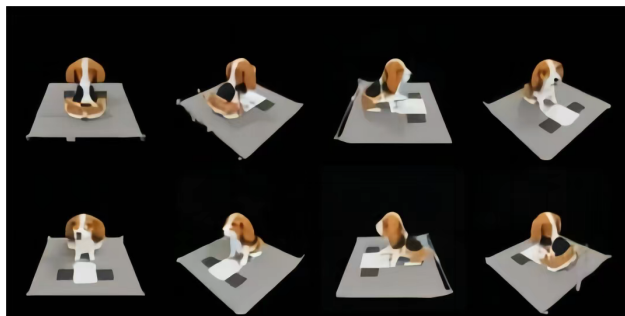Text prompt: "a chimpanzee stirring a bubbling purple potion in a cauldron."

(c) **PointLLM.**

(d) **Cap3D.**

Text prompt: "a chihuahua wearing a tutu."

(e) **PointLLM.**

(f) **Cap3D.**

Text prompt: "a confused beagle sitting at a desk working on homework."

Figure 9. **Text-to-3D generation results of models trained with different captions.** The Model trained with PointLLM-generated captions generates objects more closely aligned with text prompts and exhibits more precise detailing.

# References

[1] Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc, Antoine Miech, Iain Barr, Yana Hasson, Karel Lenc, Arthur Mensch, Katherine Millican, Malcolm Reynolds, et al. Flamingo: a visual language model for few-shot learning. In *NeurIPS*, 2022. 1, 2, 3

[2] Anonymous. Text-to-3d generation with bidirectional diffusion using both 3d and 2d priors. https://openreview.net/forum?id=V8PhVhb4pp, 2023. 12

[3] Anas Awadalla, Irena Gao, Josh Gardner, Jack Hessel, Yusuf Hanafy, Wanrong Zhu, Kalyani Marathe, Yonatan Bitton, Samir Gadre, Shiori Sagawa, et al. Openflamingo: An open-source framework for training large autoregressive vision-language models. *arXiv:2308.01390*, 2023. 2

[4] Satanjeev Banerjee and Alon Lavie. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In *ACL Workshop*, 2005. 2, 6

[5] James Betker, Gabriel Goh, 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. https://cdn.openai.com/papers/dall-e-3.pdf, 2023. 12

[6] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *NeurIPS*, 2020. 1

[7] 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. 6

[8] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv:2204.02311*, 2022. 1

[9] Wenliang Dai, Junnan Li, Dongxu Li, Anthony Meng Huat Tiong, Junqi Zhao, Weisheng Wang, Boyang Li, Pascale Fung, and Steven Hoi. Instructblip: Towards general-purpose vision-language models with instruction tuning. *arXiv:2305.06500*, 2023. 1, 2, 3, 5, 6, 8, 9, 22, 23

[10] Tri Dao, Dan Fu, Stefano Ermon, Atri Rudra, and Christopher Ré. Flashattention: Fast and memory-efficient exact attention with io-awareness. In *NeurIPS*, 2022. 11

[11] 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 *CVPR*, 2023. 2, 3, 5, 6, 8, 12

[12] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv:2303.03378*, 2023. 1, 2

[13] Samir Yitzhak Gadre, Gabriel Ilharco, Alex Fang, Jonathan Hayase, Georgios Smyrnis, Thao Nguyen, Ryan Marten, Mitchell Wortsman, Dhruba Ghosh, Jieyu Zhang, et al. Datacomp: In search of the next generation of multimodal datasets. *arXiv:2304.14108*, 2023. 11

[14] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient visual instruction model. *arXiv:2304.15010*, 2023. 2

[15] Tianyu Gao, Xingcheng Yao, and Danqi Chen. Simcse: Simple contrastive learning of sentence embeddings. *arXiv:2104.08821*, 2021. 2, 6

[16] 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 *CVPR*, 2023. 3

[17] Tao Gong, Chengqi Lyu, Shilong Zhang, Yudong Wang, Miao Zheng, Qian Zhao, Kuikun Liu, Wenwei Zhang, Ping Luo, and Kai Chen. Multimodal-gpt: A vision and language model for dialogue with humans. *arXiv:2305.04790*, 2023. 2

[18] Ziyu Guo, Renrui Zhang, Xiangyang Zhu, Yiwen Tang, Xianzheng Ma, Jiaming Han, Kexin Chen, Peng Gao, Xianzhi Li, Hongsheng Li, et al. Point-bind & point-llm: Aligning point cloud with multi-modality for 3d understanding, generation, and instruction following. *arXiv preprint arXiv:2309.00615*, 2023. 3, 5, 6, 8, 11, 22

[19] Agrim Gupta, Piotr Dollar, and Ross Girshick. Lvis: A dataset for large vocabulary instance segmentation. In *CVPR*, 2019. 12

[20] Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning without training. In *CVPR*, 2023. 2

[21] Jiaming Han, Renrui Zhang, Wenqi Shao, Peng Gao, Peng Xu, Han Xiao, Kaipeng Zhang, Chris Liu, Song Wen, Ziyu Guo, et al. Imagebind-llm: Multi-modality instruction tuning. *arXiv:2309.03905*, 2023. 3

[22] Yaru Hao, Haoyu Song, Li Dong, Shaohan Huang, Zewen Chi, Wenhui Wang, Shuming Ma, and Furu Wei. Language models are general-purpose interfaces. *arXiv:2206.06336*, 2022. 1

[23] Deepti Hegde, Jeya Maria Jose Valanarasu, and Vishal Patel. Clip goes 3d: Leveraging prompt tuning for language grounded 3d recognition. In *ICCV*, 2023. 2

[24] Dan Hendrycks and Kevin Gimpel. Gaussian error linear units (gelus). *arXiv:1606.08415*, 2016. 6

[25] 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. *NeurIPS*, 2023. 3, 5, 6, 8, 23

[26] Rongjie Huang, Mingze Li, Dongchao Yang, Jiatong Shi, Xuankai Chang, Zhenhui Ye, Yuning Wu, Zhiqing Hong, Jiawei Huang, Jinglin Liu, et al. Audiogpt: Understanding and generating speech, music, sound, and talking head. *arXiv:2304.12995*, 2023. 1, 2

[27] Shaohan Huang, Li Dong, Wenhui Wang, Yaru Hao, Saksham Singhal, Shuming Ma, Tengchao Lv, Lei Cui,

Owais Khan Mohammed, Qiang Liu, et al. Language is not all you need: Aligning perception with language models. *arXiv:2302.14045*, 2023. 1, 2

[28] Tianyu Huang, Bowen Dong, Yunhan Yang, Xiaoshui Huang, Rynson WH Lau, Wanli Ouyang, and Wangmeng Zuo. Clip2point: Transfer clip to point cloud classification with image-depth pre-training. In *ICCV*, 2023. 2

[29] Gabriel Ilharco, Mitchell Wortsman, Ross Wightman, Cade Gordon, Nicholas Carlini, Rohan Taori, Achal Dave, Vaishaal Shankar, Hongseok Namkoong, John Miller, Hannaneh Hajishirzi, Ali Farhadi, and Ludwig Schmidt. Openclip, 2021. 11

[30] Biao Jiang, Xin Chen, Wen Liu, Jingyi Yu, Gang Yu, and Tao Chen. Motiongpt: Human motion as a foreign language. *arXiv:2306.14795*, 2023. 2

[31] Heewoo Jun and Alex Nichol. Shap-e: Generating conditional 3d implicit functions. *arXiv preprint arXiv:2305.02463*, 2023. 1, 12

[32] Alexander Kirillov, Eric Mintun, Nikhila Ravi, Hanzi Mao, Chloe Rolland, Laura Gustafson, Tete Xiao, Spencer Whitehead, Alexander C Berg, Wan-Yen Lo, et al. Segment anything. *ICCV*, 2023. 3

[33] Bo Li, Yuanhan Zhang, Liangyu Chen, Jinghao Wang, Fanyi Pu, Jingkang Yang, Chunyuan Li, and Ziwei Liu. Mimic-it: Multi-modal in-context instruction tuning. *arXiv:2306.05425*, 2023. 2

[34] Junnan Li, Dongxu Li, Caiming Xiong, and Steven Hoi. Blip: Bootstrapping language-image pre-training for unified vision-language understanding and generation. 2022. 3

[35] Junnan Li, Dongxu Li, Silvio Savarese, and Steven Hoi. Blip-2: Bootstrapping language-image pre-training with frozen image encoders and large language models. *arXiv:2301.12597*, 2023. 1, 2, 3

[36] Chin-Yew Lin. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, 2004. 2, 6

[37] Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning. *arXiv:2304.08485*, 2023. 1, 2, 3, 5, 6, 7, 8, 22, 23

[38] Minghua Liu, Ruoxi Shi, Kaiming Kuang, Yinhao Zhu, Xuanlin Li, Shizhong Han, Hong Cai, Fatih Porikli, and Hao Su. Openshape: Scaling up 3d shape representation towards open-world understanding. *arXiv preprint arXiv:2305.10764*, 2023. 2

[39] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 11

[40] Tiange Luo, Chris Rockwell, Honglak Lee, and Justin Johnson. Scalable 3d captioning with pretrained models. *arXiv:2306.07279*, 2023. 2, 3, 5, 9, 11, 12, 15, 17

[41] Alex Nichol, Heewoo Jun, Prafulla Dhariwal, Pamela Mishkin, and Mark Chen. Point-e: A system for generating 3d point clouds from complex prompts. *arXiv:2212.08751*, 2022. 1, 12

[42] OpenAI. Chatgpt. https://openai.com/blog/chatgpt, 2022. 1, 2

[43] OpenAI. Gpt-4 technical report. *arXiv:2303.08774*, 2023. 1, 2, 3

[44] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. In *NeurIPS*, 2022. 1

[45] Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. Bleu: a method for automatic evaluation of machine translation. In *ACL*, 2002. 2, 6

[46] Shishir G Patil, Tianjun Zhang, Xin Wang, and Joseph E Gonzalez. Gorilla: Large language model connected with massive apis. *arXiv preprint arXiv:2305.15334*, 2023. 2

[47] 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. *arXiv:2306.14824*, 2023. 2

[48] Charles R Qi, Hao Su, Kaichun Mo, and Leonidas J Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, 2017. 8

[49] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. Learning transferable visual models from natural language supervision. 2021. 2, 3, 11

[50] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. In *JMLR*, 2020. 1

[51] Nils Reimers and Iryna Gurevych. Sentence-bert: Sentence embeddings using siamese bert-networks. *arXiv:1908.10084*, 2019. 2, 6

[52] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. High-resolution image synthesis with latent diffusion models. In *CVPR*, 2022. 12

[53] Yixuan Su, Tian Lan, Huayang Li, Jialu Xu, Yan Wang, and Deng Cai. Pandagpt: One model to instruction-follow them all. *arXiv:2305.16355*, 2023. 2

[54] Qinghong Sun, Yangguang Li, ZeXiang Liu, Xiaoshui Huang, Fenggang Liu, Xihui Liu, Wanli Ouyang, and Jing Shao. Unig3d: A unified 3d object generation dataset. *arXiv:2306.10730*, 2023. 2

[55] Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. *arXiv:2303.08128*, 2023. 2

[56] InternLM Team. Internlm: A multilingual language model with progressively enhanced capabilities. https://github.com/InternLM/InternLM, 2023. 1

[57] Hugo Touvron, Thibaut Lavril, Gautier Izacard, Xavier Martinet, Marie-Anne Lachaux, Timothée Lacroix, Baptiste Rozière, Naman Goyal, Eric Hambro, Faisal Azhar, et al. Llama: Open and efficient foundation language models. *arXiv:2302.13971*, 2023. 1, 6

[58] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017. 4

[59] Haowei Wang, Jiji Tang, Jiayi Ji, Xiaoshuai Sun, Rongsheng Zhang, Yiwei Ma, Minda Zhao, Lincheng Li, Zeng Zhao, Tangjie Lv, et al. Beyond first impressions: Integrating joint

multi-modal cues for comprehensive 3d representation. In *ACM MM*, 2023. 2

[60] Wenhai Wang, Zhe Chen, Xiaokang Chen, Jiannan Wu, Xizhou Zhu, Gang Zeng, Ping Luo, Tong Lu, Jie Zhou, Yu Qiao, et al. Visionllm: Large language model is also an open-ended decoder for vision-centric tasks. *arXiv:2305.11175*, 2023. 2

[61] Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah A Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions. *arXiv:2212.10560*, 2022. 9

[62] Jason Wei, Maarten Bosma, Vincent Zhao, Kelvin Guu, Adams Wei Yu, Brian Lester, Nan Du, Andrew M Dai, and Quoc V Le. Finetuned language models are zero-shot learners. In *ICLR*, 2021. 9

[63] Chenfei Wu, Shengming Yin, Weizhen Qi, Xiaodong Wang, Zecheng Tang, and Nan Duan. Visual chatgpt: Talking, drawing and editing with visual foundation models. *arXiv:2303.04671*, 2023. 2

[64] Zhirong Wu, Shuran Song, Aditya Khosla, Fisher Yu, Linguang Zhang, Xiaoou Tang, and Jianxiong Xiao. 3d shapenets: A deep representation for volumetric shapes. In *CVPR*, 2015. 2, 5, 8, 9

[65] Le Xue, Mingfei Gao, Chen Xing, Roberto Martín-Martín, Jiajun Wu, Caiming Xiong, Ran Xu, Juan Carlos Niebles, and Silvio Savarese. Ulip: Learning a unified representation of language, images, and point clouds for 3d understanding. In *CVPR*, 2023. 2

[66] 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. *arXiv:2305.08275*, 2023. 2, 6, 11

[67] Shukang Yin, Chaoyou Fu, Sirui Zhao, Ke Li, Xing Sun, Tong Xu, and Enhong Chen. A survey on multimodal large language models. *arXiv:2306.13549*, 2023. 2

[68] 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 *CVPR*, 2022. 6, 11

[69] 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 *CVPR*, 2022. 2

[70] Renrui Zhang, Jiaming Han, Chris Liu, Peng Gao, Aojun Zhou, Xiangfei Hu, Shilin Yan, Pan Lu, Hongsheng Li, and Yu Qiao. Llama-adapter: Efficient fine-tuning of language models with zero-init attention. *arXiv:2303.16199*, 2023. 1, 2

[71] 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. *arXiv:2307.03601*, 2023.

[72] Deyao Zhu, Jun Chen, Xiaoqian Shen, Xiang Li, and Mohamed Elhoseiny. Minigpt-4: Enhancing vision-language understanding with advanced large language models. *arXiv:2304.10592*, 2023. 1, 2, 3, 7

[73] 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 *ICCV*, 2023. 2