

SeqAfford: Sequential 3D Affordance Reasoning via Multimodal Large Language Model

Supplementary Materials

Supplementary organization:

A Code and Dataset	1
B Textual Prompt Format	1
B.1. Role	1
B.2. Task	1
B.3. Examples	2
B.4. Instruction	2
B.5. Template	2
B.6. Generated Instructions Example	2
C Format of Questions and and Answers	2
D Dataset Details	2
E Detailed Methodology	2
E.1. Point Upsample	2
E.2. Mutil-Object	6
F. More Visualization Results	9
G Comparison with the Previous SOTA Method	9

A. Code and Dataset

We will release our code and dataset once the paper is accepted.

B. Textual Prompt Format

This section outlines the specific design of the Textual Prompt component included in the four generative paradigms demonstrated earlier. By designing the Textual Prompt, it helps to enhance GPT-4’s understanding and generalization ability for the relevant tasks, providing suitable introductions for subsequent tasks. We have utilized the principles of prompt engineering to design modules such as Role, Task, Example, and Instruction for the Textual prompt. Additionally, we have formatted the input in JSON format to improve GPT-4’s comprehension of the text, standardizing the input of Object Name and Affordance Type. We show the details in Fig. A1.

B.1. Role

In the role module, we have preset the textual Prompt template as follows: "Role: You are an analytical assistant specializing in robotic affordance grounding. Your expertise is in creating questions that facilitate the training of robotic affordance grounding, enabling robots to reason about task execution, such as determining the appropriate part of an object to grasp."

B.2. Task

In the task module, we have preset the textual Prompt template as follows: "Task: Description: You will be provided with the name of an object. The robot is expected to use this tool to perform a variety of everyday tasks. Along with the tool name, you will receive an affordance type that the object can afford and a list of questions that have already been generated for this tool. Guidelines: Diversity: Aim for a wide range of tasks, ensuring that there is no overlap with previous ones. Daily Tasks: questions should be common and representative of the ones encountered in daily life. Leakage Avoidance: Ensure that the generated questions tasks can only afford the given affordance."

Table A1. **Dataset Statistics.** This table showcases the distribution of our instruction-tuning benchmark for the Sequential Affordance Reasoning task, which introduces a wide range of variances and realistic simulations, including seen/unseen settings, single/sequential scenarios.

Task	Setting	Train		Test	
		Shapes	Instruction-Point Cloud Pairs	Shapes	Instruction-Point Cloud Pairs
Single	Seen	16084	157,820	2287	4566
Single	Unseen	14307	124,140	2287	4566
Sequential	Seen	8156	13393	2,946	7454

B.3. Examples

Examples illustrate the types of questions that can be generated based on the object name and affordance type. These examples guide the creator of the prompt in formulating new and diverse questions. In the task module, we have preset the textual Prompt template as follows: "Examples: When cutting a rope with these scissors, which part of it should your palm touch? If you want to lift the bag, at which point is your finger most likely to carry it? Point out the areas on the microwave ideal for opening. How would you grasp the hat to best maintain its condition? If you want to boil water, at which points on the tap would you open the water valve? And so on...".

B.4. Instruction

In the task module, we have preset the textual Prompt template as follows: "Instruction: With the provided OBJECT NAME: '+object+' and AFFORDANCE TYPE: '+affordance types+', generate fifteen new affordance grounding question tasks. Use the HISTORY of generated tasks as a reference to ensure compliance with the diversity guideline. Output should be in the JSON format."

B.5. Template

Here is the template that encapsulates all the elements of the textual prompt:

B.6. Generated Instructions Example

We illustrated detailed examples of our generated instructions in Fig. [A2](#), [A3](#), [A4](#), [A5](#)

C. Format of Questions and and Answers

We have several formats of input, we show the Question List template , Explanatory Question List template and Answer List template in Fig. [A6](#)

D. Dataset Details

Our dataset comprises 162,386 instruction-point cloud pairs in the single affordance segmentation setting and 20,847 pairs in the sequential affordance segmentation setting, making up a total of 18,371 point cloud instances across 23 object categories. We have visualized word clouds for the object categories, affordance types, and the introduction in our dataset, demonstrating the richness of our dataset. Fig. [A7](#) shows the distribution of the generated question length and Fig. [A8](#) shows the words cloud of the generated instructions.

E. Detailed Methodology

E.1. Point Upsample

Fig. [A9](#) shows the detailed method of our upsample ways. Here, we adopt Uni3D[?] as our 3D encoder, which is capable of encoding point clouds to obtain semantically rich point cloud features $F_{p_{sparse}}$. To perform dense prediction tasks, we need to obtain finer-grained point cloud features from the sparse point cloud features $F_{p_{sparse}}$. Following PointNet++[?] and DGCNN[?], we utilized their feature propagation techniques to build our feature propagation process. Specifically, we extract features from the intermediate layers of the 3D encoder and use Farthest Point Sampling (FPS) to sample different numbers of

[Role]

You are an analytical assistant specializing in robotic affordance grounding. Your expertise is in creating questions that facilitate the training of robotic affordance grounding, enabling robots to reason about task execution, such as determining the appropriate part of an object to grasp.

[Task]

Description: You will be provided with the name of an object. The robot is expected to use this tool to perform a variety of everyday tasks. Along with the tool name, you will receive an affordance type that the object can afford and a list of questions that have already been generated for this tool.

[Guidelines]

Diversity: Aim for a wide range of tasks, ensuring that there is no overlap with previous ones.

Daily Tasks: questions should be common and representative of the ones encountered in daily life.

Leakage Avoidance: Ensure that the generated questions tasks can only afford the given affordance.

[Examples]

When cutting a rope with these scissors, which part of it should your palm touch?

If you want to lift the bag, at which point is your finger most likely to carry it?

Point out the areas on the microwave ideal for opening. How would you grasp the hat to best maintain its condition?

If you want to boil water, at which points on the tap would you open the water valve?

Instruction: With the provided **OBJECT NAME** and **AFFORDANCE TYPE**, generate fifteen new affordance grounding question tasks. Use the **HISTORY** of generated tasks as a reference to ensure compliance with the diversity guideline. The output should be in the JSON format.

Figure A1. textual Prompt Format

points from the point cloud:

$$P2 = FPS(X_{point}), \quad (1)$$

$$P3 = FPS(X_{point}). \quad (2)$$

Feature propagation is then applied to obtain features $f1$ and $f2$:

$$f1 = FP(P1, P2, H8), \quad (3)$$

$$f2 = FP(P1, P3, H4). \quad (4)$$

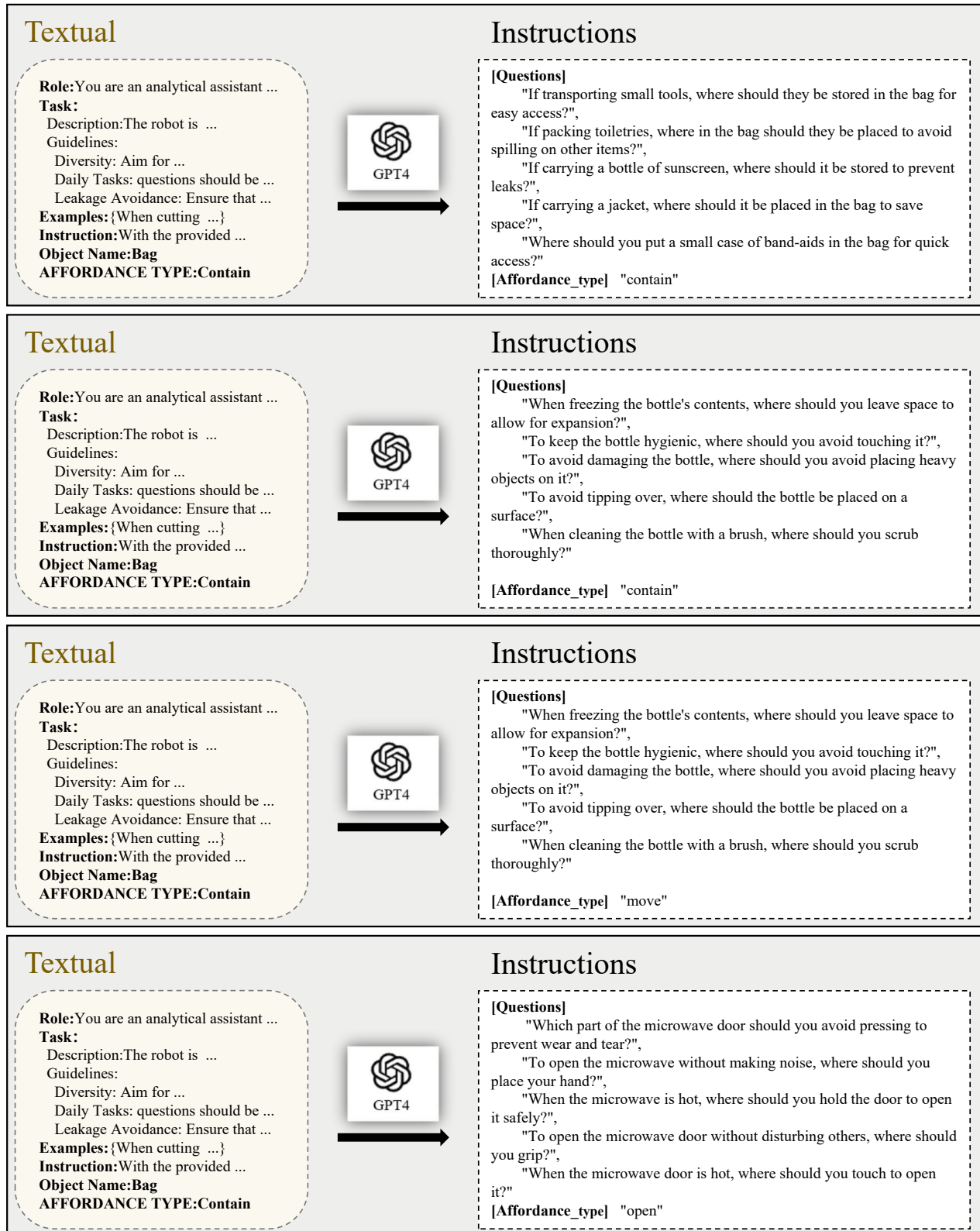


Figure A2. Generated Instructions Example:Textual Prompt

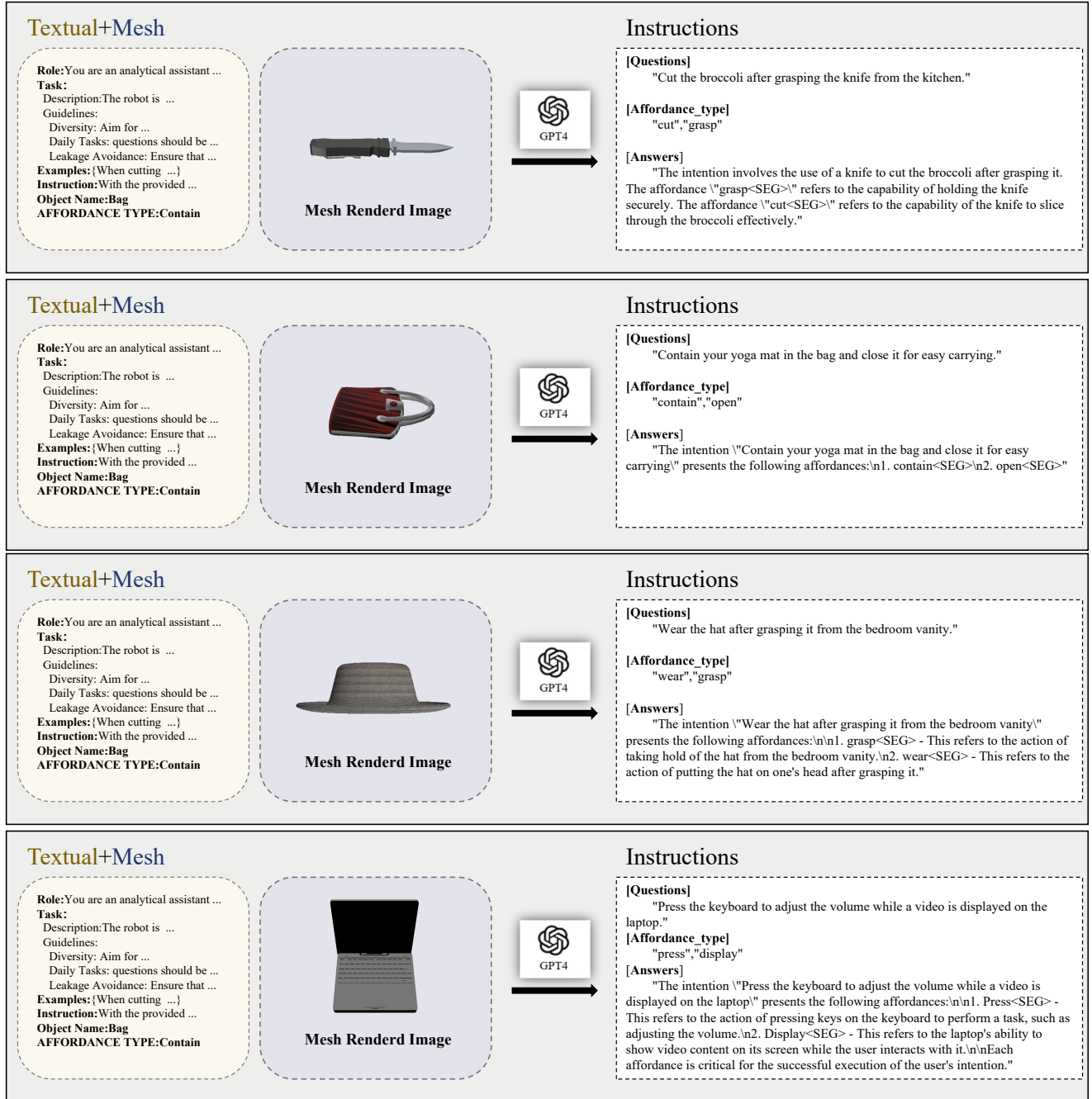


Figure A3. Generated Instructions Example:Textual Prompt with Mesh Images

Next, we follow DGCNN[?] using its upsample technique to obtain $f1'$ and $f2'$:

$$f1' = Up(P2, f1, P2.H8), \quad (5)$$

$$f2' = Up(P3, f2, P2, f1') \quad (6)$$

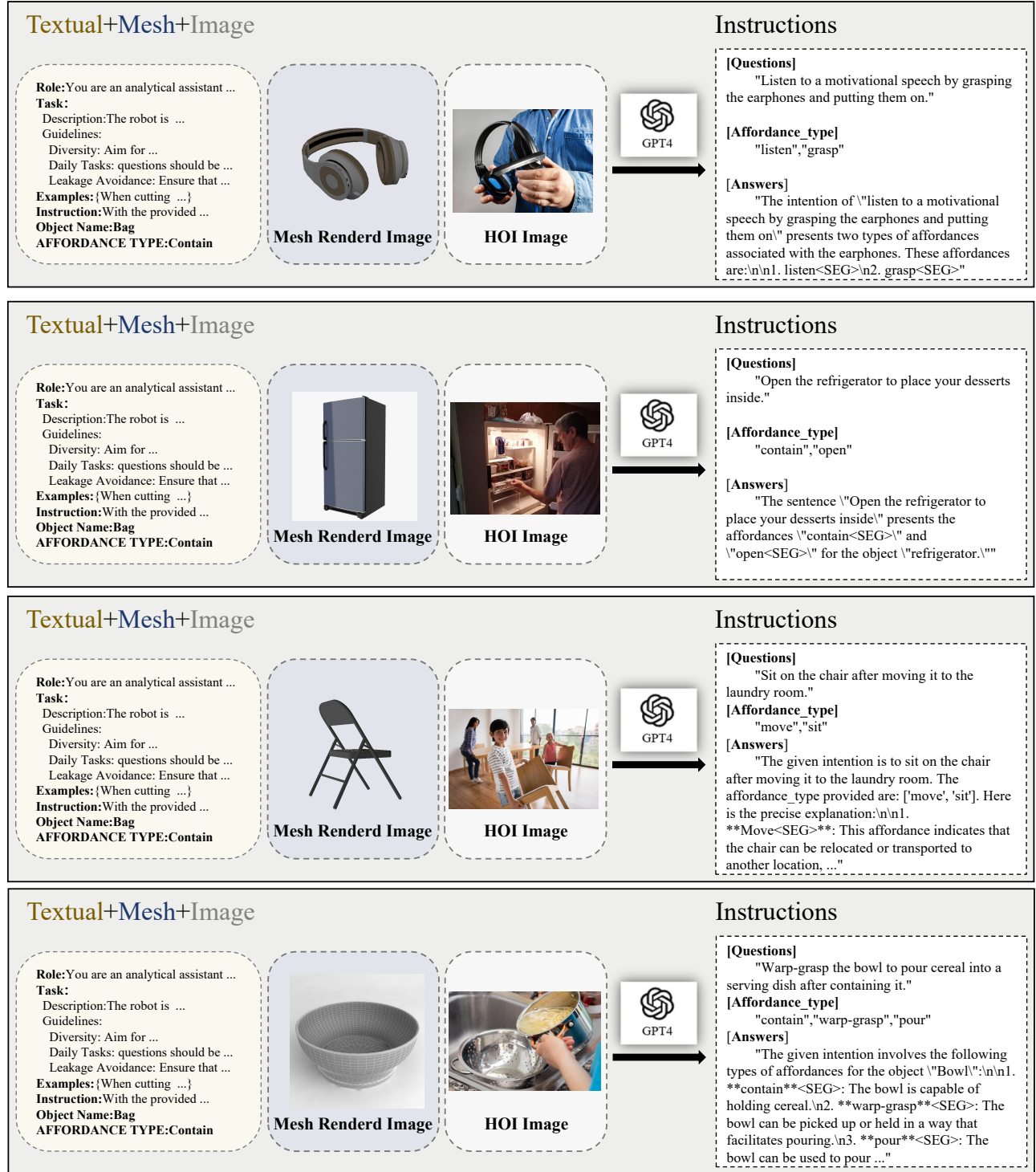


Figure A4. Genereted Instructions Example:Textual Prompt with Mesh and HOI Images

E.2. Mutil-Object

Here, we describe the detailed implementation of the Mutil-Object affordance reasoning. During training, we add $\langle SEG \rangle$ token after the object name, and our MLLM will learn the paragram and when it needs to generate text output, it will also

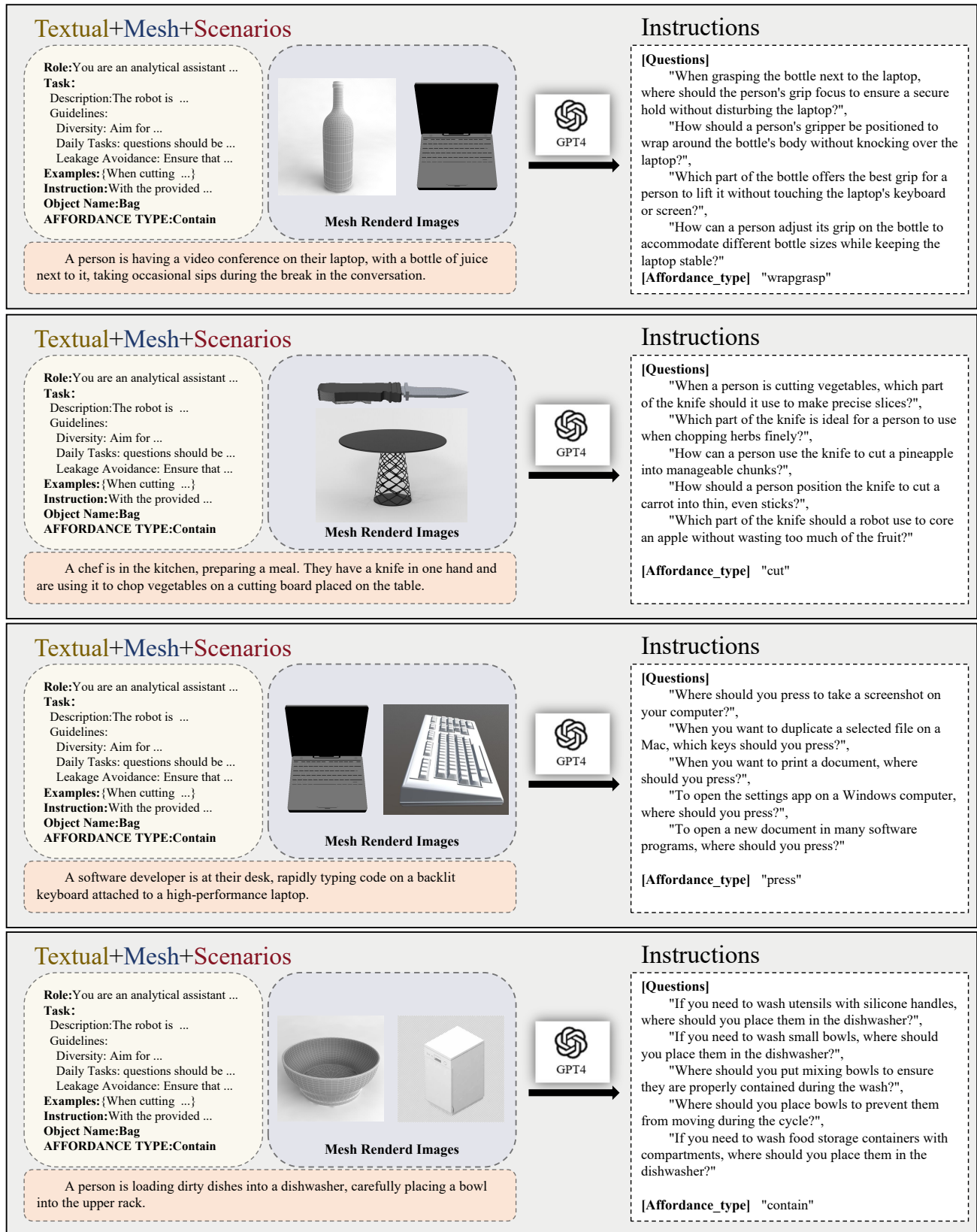


Figure A5. Genereated Instructions Example:Textual Prompt with Mesh Images and Scenarios Descriptions

[Answer List Template]

```
ANSWER_LIST = [ "{Affordance_type} are {seg}, separately.",
  "{Affordance_type} are {seg}.", "Sure, {Affordance_type} are {seg},
  separately.", "Sure, {Affordance_type} are {seg}.", "the segmentation result of
  {Affordance_type} are {seg}.", "the segmentation result of {Affordance_type}
  are {seg}, separately.", "Sure, the segmentation result of {Affordance_type} are
  {seg}.", "Sure, the segmentation result of {Affordance_type} are {seg},
  separately." ]
```

[Mutil Answer List Template]

```
MR_MULTI_ANSWER_LIST = [
  "{class_name} are {seg}, separately.",
  "{class_name} are {seg}.",
  "Sure, {class_name} are {seg}, separately.",
  "Sure, {class_name} are {seg}.",
  "the segmentation result of {class_name} are {seg}.",
  "the segmentation result of {class_name} are {seg}, separately.",
  "Sure, the segmentation result of {class_name} are {seg}.",
  "Sure, the segmentation result of {class_name} are {seg}, separately." ]
```

[Explanatory Question List Template]

```
EXPLANATORY_QUESTION_LIST = [
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Please respond with segmentation
  mask." Please output segmentation mask and explain why.
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Please output segmentation mask",
  Please output segmentation mask and explain the reason." ]
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Can you give me the segmentation
  mask?" Please output segmentation mask and give some explanation.]
```

[Question List Template]

```
LONG_QUESTION_LIST = [
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Please respond with segmentation
  mask.",
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Please output segmentation mask." ],
  DEFAULT_POINT_TOKEN + "\n" + "{sent}. Can you give me the segmentation
  mask?" ]
```

Figure A6. The details of our templates

follow this paradigm. We extract the nouns from the output of the language model, we then use the order of these nouns as the order of the model affordance reasoning. This is justified because the large language model is able to output the correct text sequence.

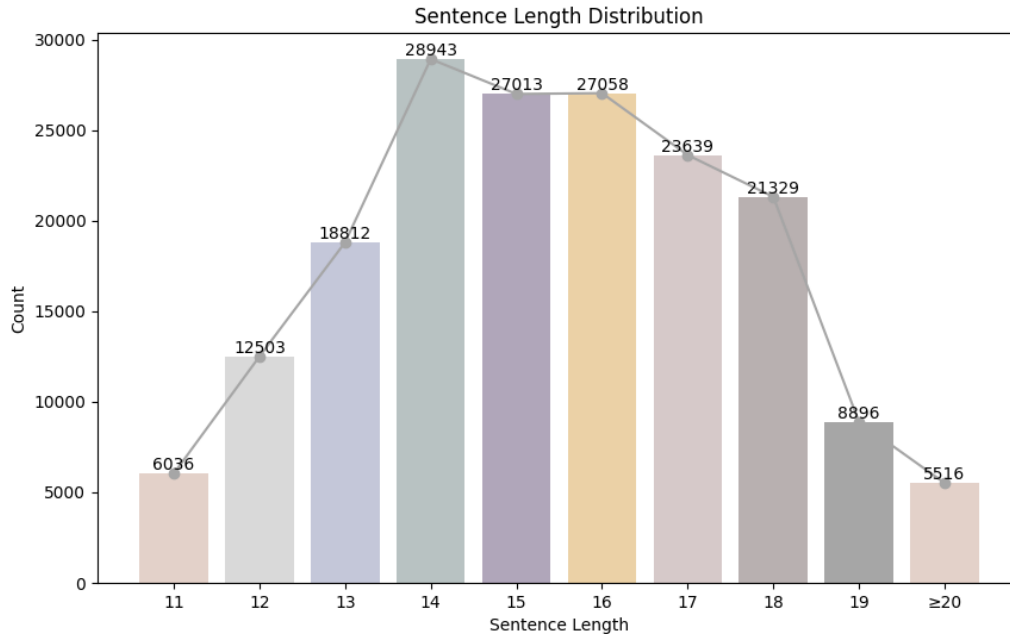


Figure A7. Question Length Distribution

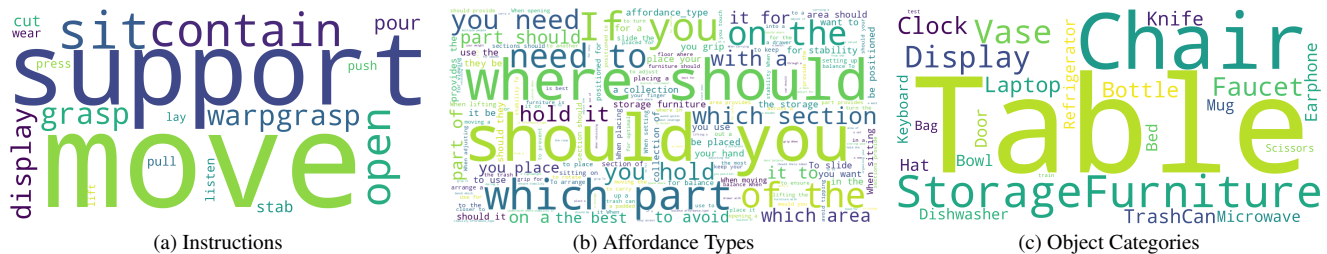


Figure A8. Word clouds of (a) instructions, (b) affordance types, and (c) object categories.

F. More Visualization Results

Here, we show more qualitative results of our model. Fig. A10 shows the additional visualization results of Single Affordance Reasoning, and Fig. A11 shows the additional visualization results of

G. Comparison with the Previous SOTA Method

We show the qualitative comparison with LASO[?]: obviously, when the question is complex, LASO[?] failed to understand the intention and do not perform well in connecting the affordance with the intention.

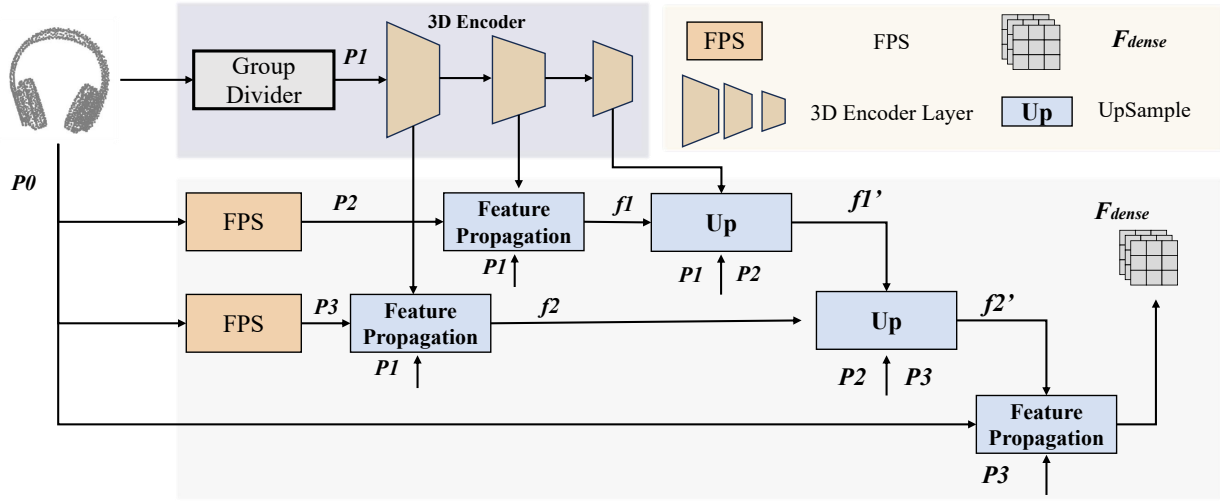


Figure A9. **Upsample.** We show the details of the Upsample method. Specifically, we use FPS and feature propagation to get the hierarchical point cloud features and upsample the sparse features into dense features

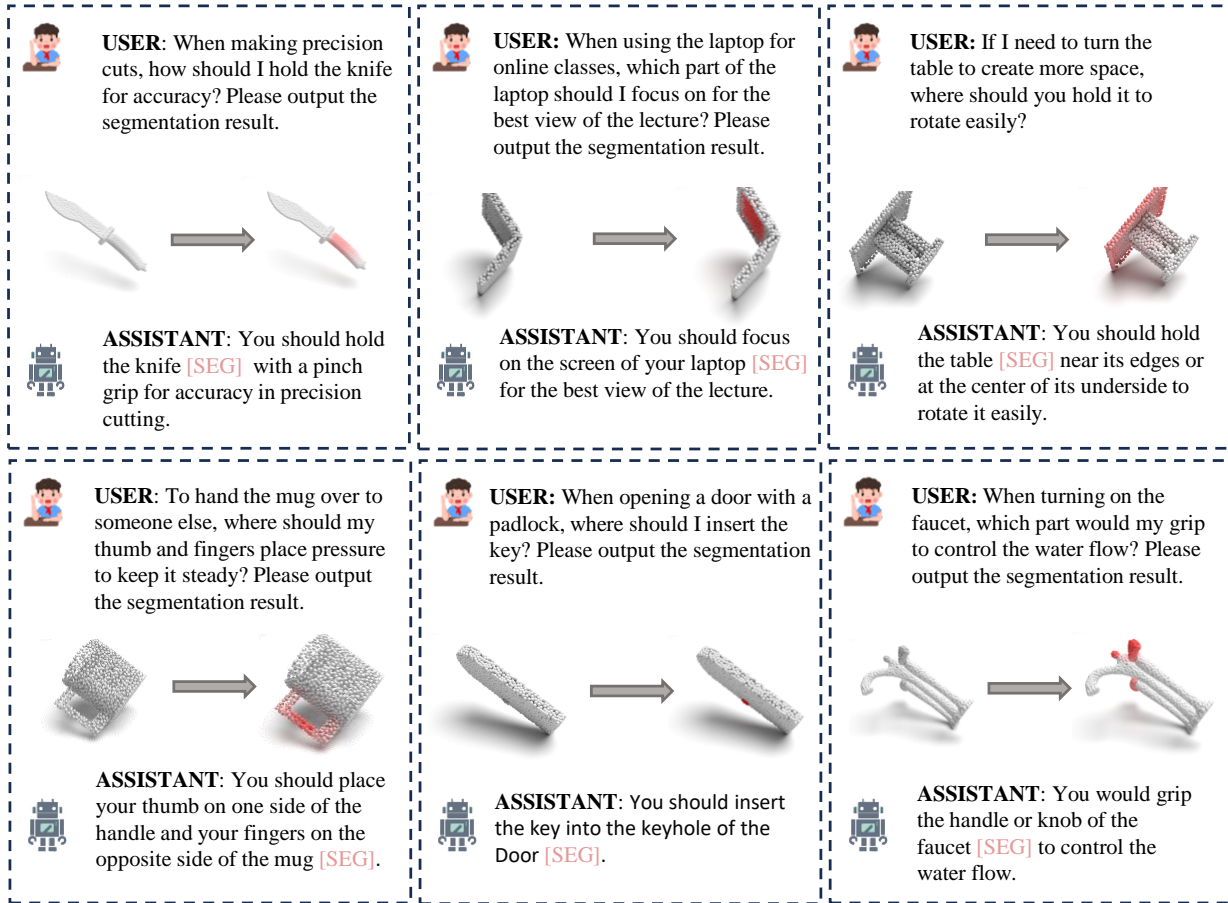


Figure A10. More Qualitative Visualization Results of Sequential Affordance Reasoning



Figure A11. More Qualitative Visualization Results of Sequential Affordance Reasoning

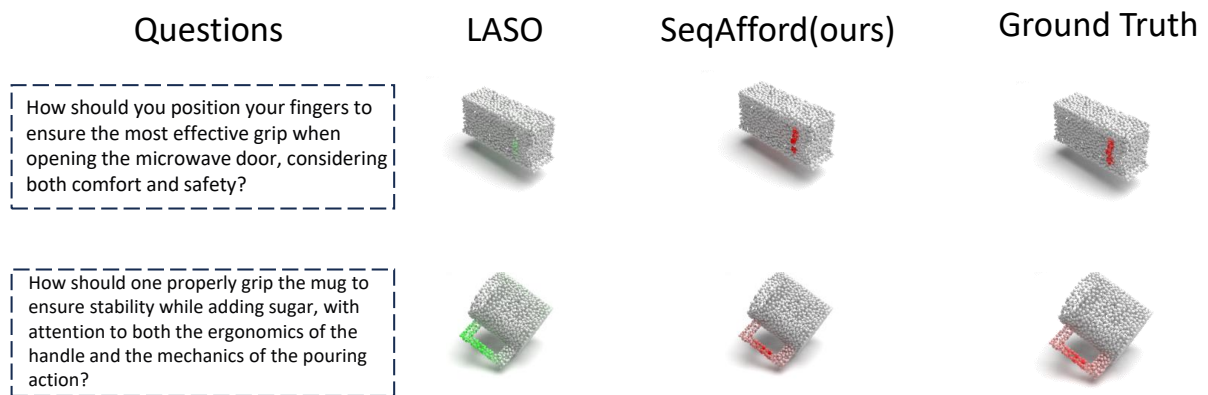


Figure A12. Comparison with the Previous SOTA Method