

white basket, green basket, brown basket, white bowl, red bowl, green bowl, red cup, green cup, white cup, white casserole, green casserole, black casserole, green cushion, yellow cushion, blue cushion, yellow dumbbell, green dumbbell, blue dumbbell, blue jug, orange jug, white jug, blue towel, red towel, green towel, grey hat, yellow hat, white plate, green plate, pink plate, black shoe, blue shoe, orange shoe, grey toy, blue toy, black toy, red vase, blue vase, pink vase, black teapot, white teapot, yellow teapot

**Table 3.** Full list of object categories used for ASK-TO-ACT task.

## A USE OF LLM FOR WRITING

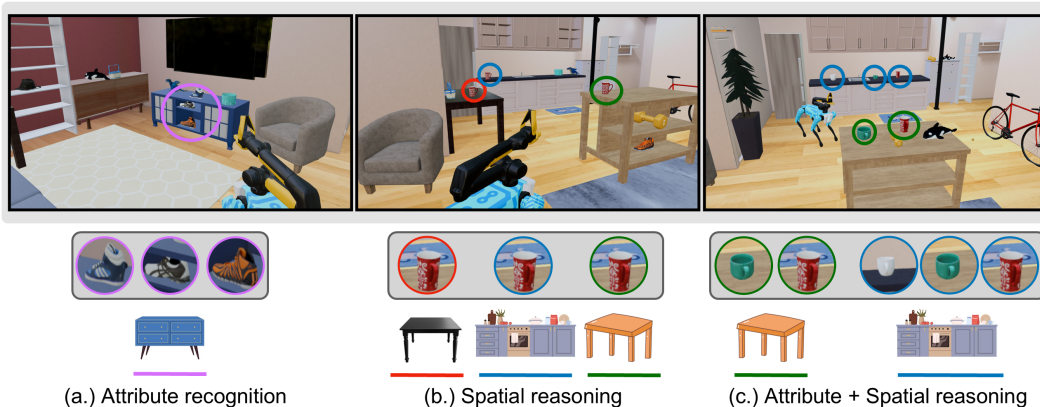
We use LLMs to assist with specific aspects of paper writing. These include: (1) grammar checking to improve language syntax, (2) sentence paraphrasing for improving readability and clarity. We use LLMs to write parts of Sections 1, 3, 4, and 5. Specifically, for each section the original content was written by the authors and paraphrased or spell checked by the GPT-4o [OpenAI \(2023\)](#) LLM. Followed by this round of LLM edits, the authors verified and edited the content written by the LLM to fix issues in generated text. After which the author edited version of the paper is used for final submission.

## B DATASET

In Fig. 5, we show additional examples of different types of ambiguities in our ASK-TO-ACT dataset. Specifically we show 3 additional examples for:

- **Attribute Recognition:** The agent must reason about appearance-based attributes such as color and object category to disambiguate the object user is looking for (e.g., a red bowl and white bowl as shown in Fig. 5 (a)).
- **Object Size:** The agent must distinguish between objects of different sizes but similar appearance and location to disambiguate the object user is looking for (e.g., a large red bowl vs. a small red bowl as shown in Fig. 5 (b)).
- **Attribute/Spatial Reasoning and Object Size:** A combination of attribute and spatial reasoning with object size, where the agent must account for both appearance or location and size to effectively disambiguate the target object the user is looking for. For the example shown in Fig. 5 (c), the agent needs to reason about different colored shoes (blue, white, and orange shoe) which are located on the blue cabinet and the sofa and some of those have different size e.g. large blue shoe on blue cabinet vs. small blue shoe on sofa and large white shoe on sofa vs. small white shoe on blue cabinet.

**Object Categories.** We present the full list of object categories used for ASK-TO-ACT dataset in Tab. 3.



702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

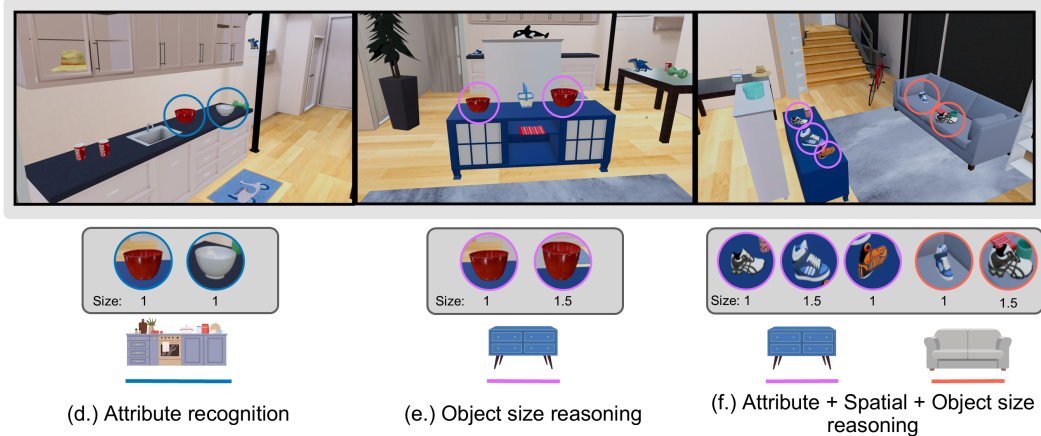


Figure 5. Dataset. Examples of different types of synthetically created ambiguities in ASK-TO-ACT tasks.

### C TRAINING DETAILS

**Problem Setup.** Our problem setup can be formulated as a Partially-Observable Markov Decision Process (POMDP), defined by a tuple  $(S, O, A, P, R, p_0, \gamma)$  where  $S$  is state space,  $O$  is observation space,  $A$  is action space,  $P$  represents transition dynamics,  $R$  is reward function,  $p_0$  is initial state distribution and  $\gamma$  is discount factor. In our setting,  $O$  is a combination of responses from the user (for any questions asked) and visual observations, which come from the robots egocentric RGB camera, and provide only partial views of the environment. We consider the extension of including a goal distribution  $G$  and the case where the reward is formulated as  $R(s, g)$  for  $s \in S$  and  $g \in G$ . We aim to learn a language-conditioned policy  $\pi(a|o, g)$  mapping from observation  $o$  and task instruction  $g$  to an action  $a$  that maximizes the sum of discounted rewards  $\mathbb{E}_{s_0 \sim p_0, g \sim G} \sum_t \gamma^t R(s_t, g)$ .

**Training Details.** To train our policy using reinforcement learning (RL) we use PPO with Generalized Advantage Estimation (GAE) (Schulman et al., 2016). We use a discount factor  $\gamma$  of 0.99 and set GAE parameter  $\tau$  to 0.95. We do not use normalized advantages. To parallelize training, we use DD-PPO (Wijmans et al., 2020), an adaptation of PPO (Schulman et al., 2017) for distributed training, with 24 environments each on 8 A40 GPUs for 50 million steps. Each worker collects 24 frames of experience from 24 environments parallelly and then performs 2 epochs of PPO update with 2 mini batches in each epoch. Tab. 4 details the default hyperparameters used in all of our training runs.

**Per-Step Reward Inference.** For per-step reward inference for evaluating questions asked by an agent, we use a vLLM (Kwon et al., 2023) server to host our LLM answer evaluation model using Llama-3.1 8B Instruct (Grattafiori et al., 2024) model. At each rollout step during RL training and policy evaluation, if the agent outputs a ‘ask\_question’ action we send a request to the vLLM server and feed the response back to the policy at the next step. In order to achieve a high throughput we run 8 parallel vLLM servers and randomly send requests to different workers from training and evaluation workers.

**Token Normalization for Action Probability and Entropy.** When using MLLMs for RL training, each action  $a \in A$  at every timestep is represented using a set of tokens  $(u_t^1, \dots, u_t^m)$ . However, the MLLM only outputs token level probabilities which is different from probability of executing an action in the environment as used in RL traditionally. The token-level probability for each action  $a_t$  can be represented by:

$$P_{\text{token}}(a_t|s) = \prod_{i=1}^m P(u_t^i|s, u_t^1, \dots, u_t^{i-1}) \tag{3}$$

To compute action probabilities, one naive approach is to take a sum over all token-level probabilities. When training MLLMs with action-space with variable number of tokens across actions one issue in Eq. (3), is that actions with larger number of tokens tend to have lower token-level probabilities, even

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772

Parameter	Value
Number of GPUs	8
Number of environments per GPU	24
Rollout length	24
PPO epochs	2
Number of mini-batches per epoch	2
LR	$2.5e^{-4}$
Optimizer	Adam
Weight decay	0.0
Epsilon	$1.0 \times 10^{-5}$
PPO clip	0.2
Generalized advantage estimation	True
$\gamma$	0.99
$\tau$	0.95
Value loss coefficient	0.5
Max gradient norm	0.2
DDPPO sync fraction	0.6

773  
774

**Table 4.** Hyperparameters used for RL finetuning.

775

776 though they might be more reasonable to take at some point in the environment. This issue happens  
777 because the probability of each token  $P(u_t^i|\cdot)$  is always less than 1. This could be problematic in  
778 RL optimization especially for the case when the agent needs to output a natural language question  
779 for ASK-TO-ACT task, simply because the questions tend to have more tokens than pre-defined  
780 skills represented in language. To remedy this issue, we use action length normalization technique  
781 normalize the token-level probabilities of the actions with the number of action tokens, which can be  
782 defined as:

783  
784  
785  
786

$$\log P(a_t|s) = \sum_{i=1}^m \log P_{\text{token}}(u_t^i|s, u_t^1, \dots, u_t^{i-1})/m \quad (4)$$

787  
788

We find action-length normalization helps stabilize training when working with our variable length  
action-space when using constrained grammar decoding (Park et al., 2025) for RL training.

789  
790  
791

## D LLM REWARD GENERATION

792  
793  
794  
795

Tab. 6 shows the prompt used by the question evaluation LLM used for ASK-TO-ACT task evaluation  
and RL training. In addition, we use the prompt 5 for generating optimal answer sequences and  
subgoals required to generate the SFT data and task subgoals described in Eq. (1) in Sec. 4.2.

796  
797  
798

$$r_t = 10 \cdot \mathbb{1}_{\text{success}} + 2.5 \cdot \mathbb{1}_{\text{subgoal}} + r_3 \cdot \mathbb{1}_{\text{useful\_question}} - 0.05 \cdot \mathbb{1}_{\text{exceed\_budget}} - 0.01, \quad (5)$$

799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

**Reward Details.** Eq. (5) shows the reward function we use with coefficients for each term. In this  
reward,  $\mathbb{1}_{\text{success}}$  indicates whether the task was successfully completed and  $\mathbb{1}_{\text{subgoal}}$  (generated by  
the LLM) indicates if the agent completed any subgoal required to complete the overall task. For  
example, for a task “Bring me the cup and place it on the coffee table”, the agent needs to first search  
for all cups, pick the correct cup, then navigate to the coffee table, and finally place it. Similarly,  
 $\mathbb{1}_{\text{useful\_question}}$  (generated by the LLM) indicates if a question asked by the agent is valid and helps  
make progress towards disambiguating the target object or not. Consider the example of fetching  
the cup, if the environment has 4 cups on a table and the agent asks “where is it the cup?” and if  
user responds ‘on the coffee table’, then the agent should ask ‘Is it the blue cup?’ or ‘Is it the yellow  
cup?’ to find the target object instead of asking ‘Is it blue cup in the sink?’ to make progress towards  
solving the task. For each useful question  $\mathbb{1}_{\text{useful}}$  the agent asks in an episode it is given a reward  
 $r_3$  until total number of questions asked are less than the question budget. For example, if the task

810  
811  
812  
813  
814  
815  
816  
817  
818  
819  
820  
821  
822  
823  
824  
825  
826  
827  
828  
829  
830  
831  
832  
833  
834  
835  
836  
837  
838  
839  
840  
841  
842  
843  
844  
845  
846  
847  
848  
849  
850  
851  
852  
853  
854  
855  
856  
857  
858  
859  
860

---

Imagine you are a household robot. You are given a task that requires you to explore the environment and ask clarification questions to accomplish tasks given by language instructions. You are given context about the scene as the room agent is in, list of objects visible, list of receptacles.

For each task given as a language instruction you have to output a sequence of actions that the agent should take. These actions can include a clarification question as well. Ask a question only when required. For each action also output the reason to take the action with it. Actions can be one of the following:

1. nav(receptacle)
2. pick(object)
3. place(receptacle)
4. open(receptacle) // Strictly use when the object needs be to picked from or placed inside an articulated receptacle after navigating to it
5. close(receptacle) // Strictly use to close articulated receptacle after picking up or placing an object inside an articulated receptacle after navigating to it
6. ask(question)

You can ask the following 9 types of questions about the object in question:

1. Is target object on the <receptacle>?
2. Is <object\_instance> the target object?
3. Is target object the <object\_size> one?
4. Where is the <object\_category> located?
5. What color is <object\_category>?
6. Can you describe the <object\_category>?
7. Is <object\_instance> clutter?
8. Are <object\_category> clutter?
9. Which <receptacle> to place the <object\_instance> on/in?

Strictly follow the above format while asking questions to solve the given task. Do not use any other types of questions.

Here is the example task:

Instruction: Bring me bowl put it on cabinet

Receptacles: [light table, chair, sofa, dark table, tv stand, cabinet, sink]

Receptacles with objects:

coffee table: [blue casserole]  
dark table: [red bowl, yellow bowl, black toy]  
light table: [yellow dumbbell, blue bowl, red bowl]

Your task is to enumerate all possible question sequences an embodied agent should ask in order to find the target object. In addition also output another list with subgoals that agent needs to achieve using actions you have access to with appropriate arguments. Only use the templates specified for generating the questions and actions. Strictly follow the next command: Output all sequence of questions as a list of list and the subgoals as a list of actions in a json.

---

861 **Table 5.** LLM prompt used for generating optimal question sequences used to generate SFT data.  
862  
863

864  
865  
866  
867  
868

---

869 Imagine you are a household robot. You are given a task that requires you to explore  
870 the environment and ask clarification questions to accomplish tasks given by language  
871 instructions. You are given context about the scene as the room agent is in, list of  
872 objects visible, list of receptacles.

873 For each task given as a language instruction you have to output a sequence of  
874 actions that the agent should take. These actions can include a clarification question as  
875 well. Ask a question only when required. For each action also output the reason to take  
876 the action with it. Actions can be one of the following:  
877

- 878 1. nav(receptacle)
- 879 2. pick(object)
- 880 3. place(receptacle)
- 881 4. open(receptacle) // Strictly use when the object needs be to picked from or placed  
882 inside an articulated receptacle after navigating to it
- 883 5. close(receptacle) // Strictly use to close articulated receptacle after picking up or  
884 placing an object inside an articulated receptacle after navigating to it
- 885 6. ask(question)

886 Your task is to evaluate the questions asked by an embodied agent to identify the  
887 target object and resolve ambiguity about placement preference if any. You will be  
888 provided with task instruction, objects in environment, current location of objects in  
889 the environment, list of target objects and their current location, and desired placement  
890 location of each target object. Respond with a concise answer with either object name,  
891 receptacle name, object attributes and appearance, or yes/no response followed by a  
892 boolean that denotes whether the question is useful towards making progress to resolve  
893 ambiguity or not. Here is the example task:

894 Instruction: Place the bowl in kitchen counter drawer  
895

896 Receptacles: [light table, chair, sofa, dark table, tv stand, cabinet, sink, left  
897 kitchen counter drawer, right kitchen counter drawer, top kitchen counter drawer, bottom  
898 kitchen counter drawer, fridge]  
899

900 Receptacles with objects:  
901  
902 coffee table: [blue casserole]  
903 dark table: [red bowl, yellow bowl, black toy]  
904 light table: [yellow dumbbell, blue bowl, red bowl]  
905

906 Target Object Metadata: {"name": "red bowl", "current location": "on light table",  
907 "target location": cabinet}  
908

908 Questions asked so far: []  
909 Current question: ask("where is the red bowl?")  
910 Answer:  
911

---

912 **Table 6.** LLM prompt used for generating reward for evaluating ‘ask question’ action during RL training.  
913  
914  
915  
916  
917

requires 3 questions then the agent will get  $r_3 = 0.5$  for each relevant question it asks as evaluated by the LLM and 0 otherwise or if the question budget is exceeded. When training policies under a budget of questions we penalize the agent for every question that exceeds pre-specified budget given by  $\mathbb{1}_{\text{exceed\_budget}}$ . By default the budget is set to minimum number of required questions  $K$  to solve the task.

## E QUALITATIVE EXAMPLES

We present additional qualitative examples of evaluating our method on UNSEEN TASKS evaluation split of ASK-TO-ACT dataset in Fig. 6.

Task: Move the bowl and place it in kitchen counter drawer



Task: Bring the hand towel and place it on kitchen counter

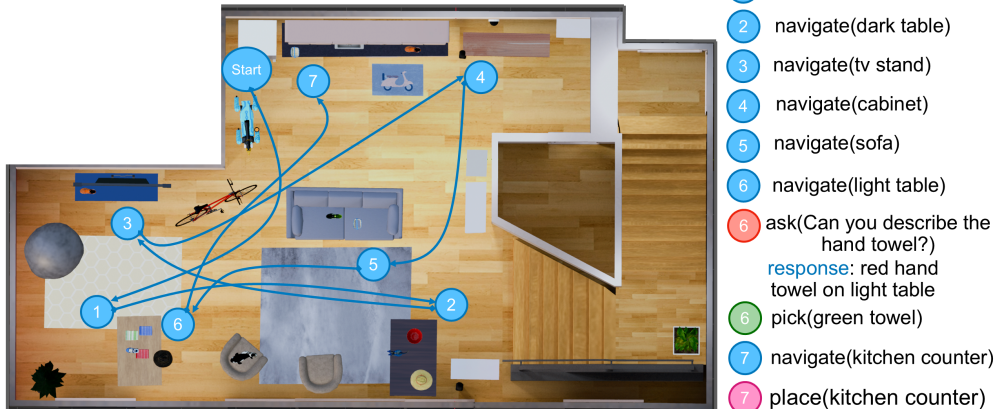


Figure 6. Qualitative Example. Successful trajectories of our method on 2 evaluation episodes from UNSEEN TASKS split.

## F BASELINE DETAILS

(a.) **Fully Observable Text WorldGraph + ReAct (Zero-shot).** In this baseline, we provide an LLM (GPT4o in our case) with a fully observable text-based world graph describing the environment, including receptacles, objects, and their locations (e.g. “The apple is on the coffee table.”). The prompt used for this baseline is shown in Appendix G.

(b.) **Fully Observable Text WorldGraph + ReAct (Few-shot).** This baseline extends (a) by providing the LLM with a few in-context examples that demonstrate task planning and ambiguity resolution strategies in the ASK-TO-ACT task. By leveraging demonstrations, this approach assesses

whether in-context learning improves LLMs task planning and ambiguity reasoning under full observability. The prompt used for this baseline is shown in Appendix G.

**(c.) Partially Observable Text WorldGraph + ReAct (Few-shot).** While baselines (a) and (b) assume privileged access to a fully observable world graph, constructing such representations in real-world settings is often infeasible. This baseline relaxes that assumption by providing the LLM with a partially observable text-based world graph. At the start of an episode, the agent lacks full knowledge of object locations and must actively explore to gather necessary information. The prompt used for this baseline is shown in Appendix G.

**(d.) Vision GPT4o + SoM + ReAct.** Building an error-free text representation of real-world environments is challenging. This baseline evaluates whether existing MLLMs can solve the ASK-TO-ACT task using egocentric visual observations. At each timestep, the MLLM receives the robots visual input along with an skill library for executing actions. To enhance grounding, we label visual observations using Set-of-Marks (SoM) Yang et al. (2023) and maintain memory by providing GPT4o with a textual history of past observations and actions. The prompt used for this baseline is shown in Appendix G, in addition to the prompt this baseline takes in visual observations from robots egocentric camera augmented with SoM Yang et al. (2023).

## G LIMITATIONS

In this work, the AUTOASK method is purposefully constrained to questions relevant to the ASK-TO-ACT task to simplify the learning problem for a small scale 0.5B model. We find that using completely unconstrained question sampling (*i.e.* without explicitly including useful question templates in the prompt) leads to significantly poor performance for all finetuned and zero-shot methods. We believe this limitation can be addressed by using larger scale models (3B, 7B or 72B) which can be explored as part of future work. Another limitation of our work is the limited evaluation setup of 7 tasks used in ASK-TO-ACT which is not a complete representative for real-world scenarios. An important future work is to study a more realistic setup of ASK-TO-ACT with a much larger and more complex set of ambiguous tasks and an open-ended question-answer setup. There are a range of open research questions required to do this, including the use (and potentially training/distillation) of efficient LLM-as-reward models that can be used within an RL pipeline. Further study is also needed to determine exactly how accurate or reliable the rewards must be to train effective models. As it is anticipated that LLM-generated rewards will not be as accurate in less constrained settings. These aspects are limitations of our current work but exciting research directions for future work.

### WG + ReAct Zero-Shot Baseline Prompt

Imagine you are a household robot. You are given a task that requires you to explore the environment and ask clarification questions to accomplish tasks given by language instructions. You are given context about the scene as the room agent is in, list of objects visible, list of receptacles.

For each task given as a language instruction you have to output a sequence of actions that the agent should take. These actions can include a clarification question as well. Ask a question only when required. For each action also output the reason to take the action with it. Actions can be one of the following:

1. nav(receptacle)
2. pick(object)
3. place(receptacle)
4. open(receptacle) // Strictly use when the object needs be to picked from or placed inside an articulated receptacle after navigating to it
5. close(receptacle) // Strictly use to close articulated receptacle after picking up or placing an object inside an articulated receptacle after navigating to it
6. ask(question)

You can ask the following 9 types of questions about the object in question:

1. Is target object on the <receptacle>?
2. Is <object\_instance> the target object?
3. Is target object the <object\_size> one?
4. Where is the <object\_category> located?
5. What color is <object\_category>?
6. Can you describe the <object\_category>?
7. Is <object\_instance> clutter?

1026  
1027  
1028  
1029  
1030  
1031  
1032  
1033  
1034  
1035  
1036  
1037  
1038  
1039  
1040  
1041  
1042  
1043  
1044  
1045  
1046  
1047  
1048  
1049  
1050  
1051  
1052  
1053  
1054  
1055  
1056  
1057  
1058  
1059  
1060  
1061  
1062  
1063  
1064  
1065  
1066  
1067  
1068  
1069  
1070  
1071  
1072  
1073  
1074  
1075  
1076  
1077  
1078  
1079

8. Are <object\_category> clutter?
9. Which <receptacle> to place the <object\_instance> on/in?

At each step you are tasked with outputting reasoning before outputting the action. You should output the reason and action in the following format:

Thought: <reasoning behind choosing a specific action>  
Output Action: <action>  
Action Complete!

Strictly follow the above format when outputting actions. Every time you output a ask question action wait for the human to reply. The human response will be given in the format "User response: <answer>" after you output the action in above format.

Strictly follow the above format while solving given task.

Instruction: {instruction}

Receptacles: {receptacles}

Receptacles with objects:  
{receptacles\_with\_objects}

Observation action history:  
{prev\_observations}

Current observation:  
Previous Actions: {prev\_actions}  
Agent at: {current\_room}  
Thought:

### WG + ReAct Few-Shot Baseline Prompt

Imagine you are a household robot. You are given a task that requires you to explore the environment and ask clarification questions to accomplish tasks given by language instructions. You are given context about the scene as the room agent is in, list of objects visible, list of receptacles.

For each task given as a language instruction you have to output a sequence of actions that the agent should take. These actions can include a clarification question as well. Ask a question only when required. For each action also output the reason to take the action with it. Actions can be one of the following:

1. nav(receptacle)
2. pick(object)
3. place(receptacle)
4. open(receptacle) // Strictly use when the object needs be to picked from or placed inside an articulated receptacle after navigating to it
5. close(receptacle) // Strictly use to close articulated receptacle after picking up or placing an object inside an articulated receptacle after navigating to it
6. ask(question)

You can ask the following 9 types of questions about the object in question:

1. Is target object on the <receptacle>?
2. Is <object\_instance> the target object?
3. Is target object the <object\_size> one?
4. Where is the <object\_category> located?
5. What color is <object\_category>?
6. Can you describe the <object\_category>?
7. Is <object\_instance> clutter?
8. Are <object\_category> clutter?
9. Which <receptacle> to place the <object\_instance> on/in?

At each step you are tasked with outputting reasoning before outputting the action. You should output the reason and action in the following format:

Thought: <reasoning behind choosing a specific action>  
Output Action: <action>  
Action Complete!

Strictly follow the above format when outputting actions. Every time you output a ask question action wait for the human to reply. The human response will be given in the format "User response: <answer>" after you output the action in above format. After you ask a question wait for the human to reply.

Here are two examples of completing the task:

1080  
1081 Example 1:  
1082 Instruction: Bring me bowl put it on cabinet

1083 Receptacles: [light table, chair, sofa, dark table, tv stand, cabinet, sink, fridge, right counter  
1084 drawer, left counter drawer]

1085 Receptacles with objects:  
1086 coffee table: blue casserole  
1087 dark table: red bowl, yellow bowl, black toy  
1088 light table: yellow dumbbell, blue bowl, red bowl

1089 Observation action history:  
1090 Step 1:  
1091 Previous Actions: []  
1092 Agent at: Start  
1093 Thought: I see two red bowls, one yellow bowl and one blue bowl. It's unclear which one user  
1094 wants. I should ask about location of the bowl first and then the color to find the correct one.  
1095 Output Action: ask("Is the target object on the light table?")  
1096 Action Complete!

1097 User response: no  
1098 Step 2:  
1099 Previous Actions: [ask("Is the target object on the light table??")]  
1100 Agent at: Start  
1101 Thought: User said it's not on light table. This means user wants either a red or a yellow bowl.  
1102 I should ask about the color of the bowl user wants.  
1103 Output Action: ask("Is the target object the red bowl?")  
1104 Action Complete!

1105 User response: yes  
1106 Step 3:  
1107 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
1108 red bowl?")]  
1109 Agent at: Start  
1110 Thought: User said it's red bowl that is on the dark table. I will navigate to it.  
1111 Output Action: nav(dark table)  
1112 Action Complete!

1113 Step 4:  
1114 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
1115 red bowl?"), nav(dark table)]  
1116 Agent at: dark table  
1117 Thought: I need to pick the red bowl.  
1118 Output Action: pick(red bowl)  
1119 Action Complete!

1120 Step 5:  
1121 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
1122 red bowl?"), nav(dark table), pick(red bowl)]  
1123 Agent at: dark table  
1124 Thought: I need to navigate to cabinet to place it.  
1125 Output Action: nav(cabinet)  
1126 Action Complete!

1127 Step 6:  
1128 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
1129 red bowl?"), nav(dark table), pick(red bowl), nav(cabinet)]  
1130 Agent at: cabinet  
1131 Thought: I need to place red bowl on the cabinet.  
1132 Output Action: place(cabinet)  
1133 Action Complete!

1134 Current observation:  
1135 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
1136 red bowl?"), nav(dark table), pick(red bowl), nav(cabinet), place(cabinet)]  
1137 Agent at: cabinet  
1138 Thought: Task is done.  
1139 Output Action: Done()

1140 Example 2:  
1141 Instruction: Clear clutter from the sofa and put it on the cabinet

1142 Receptacles: [light table, chair, sofa, dark table, tv stand, cabinet, sink, fridge, right counter  
1143 drawer, left counter drawer]

1144 Receptacles with objects:  
1145 coffee table: blue casserole  
1146 sofa: red bowl, black toy, blue casserole

1134 light table: yellow dumbbell, blue toy  
1135  
1136 Observation action history:  
1137 Step 1:  
1138 Previous Actions: []  
1139 Agent at: Start  
1140 Thought: I don't know what the user means by clutter. I should first identify which objects on  
1141 couch are user thinks are clutter.  
1142 Output Action: ask("Is red bowl clutter?")  
1143 Action Complete!  
1144  
1145 User response: yes  
1146 Step 2:  
1147 Previous Actions: [ask("Is red bowl clutter?")]  
1148 Agent at: Start  
1149 Thought: User said red bowl is clutter. There could be more objects which are clutter so I should  
1150 ask about it.  
1151 Output Action: ask("Is blue casserole clutter?")  
1152 Action Complete!  
1153  
1154 User response: yes  
1155 Step 3:  
1156 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?")]  
1157 Agent at: Start  
1158 Thought: User said red bowl and blue casserole are clutter. There could be more objects which are  
1159 clutter so I should ask about it.  
1160 Output Action: ask("Is black toy clutter?")  
1161 Action Complete!  
1162  
1163 User response: no  
1164 Step 4:  
1165 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), ask("Is black  
1166 toy clutter?")]  
1167 Agent at: Start  
1168 Thought: Looks like only red bowl and blue casserole are clutter and need to be moved to cabinet.  
1169 Output Action: nav(sofa)  
1170 Action Complete!  
1171  
1172 Step 5:  
1173 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa)]  
1174 Agent at: sofa  
1175 Thought: I need to pick the red bowl and then blue casserole.  
1176 Output Action: pick(red bowl)  
1177 Action Complete!  
1178  
1179 Step 6:  
1180 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
1181 pick(red bowl)]  
1182 Agent at: cabinet  
1183 Thought: I need to navigate to cabinet to place red bowl on the cabinet.  
1184 Output Action: nav(cabinet)  
1185 Action Complete!  
1186  
1187 Step 7:  
1188 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
1189 pick(red bowl), nav(cabinet)]  
1190 Agent at: cabinet  
1191 Thought: I need to place red bowl on the cabinet.  
1192 Output Action: place(cabinet)  
1193 Action Complete!  
1194  
1195 Step 8:  
1196 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
1197 pick(red bowl), nav(cabinet), place(cabinet)]  
1198 Agent at: cabinet  
1199 Thought: I need to go back to sofa to pick blue casserole.  
1200 Output Action: nav(sofa)  
1201 Action Complete!  
1202  
1203 Step 9:  
1204 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
1205 pick(red bowl), nav(cabinet), place(cabinet), nav(sofa)]  
1206 Agent at: sofa  
1207 Thought: I need to pick blue casserole.  
1208 Output Action: pick(blue casserole)  
1209 Action Complete!  
1210  
1211 Step 10:

1188 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
 1189 pick(red bowl), nav(cabinet), place(cabinet), nav(sofa), pick(blue casserole)]  
 1190 Agent at: sofa  
 1191 Thought: I need to navigate to cabinet to place blue casserole on the cabinet.  
 1192 Output Action: nav(cabinet)  
 1193 Action Complete!

1193 Current observation:  
 1194 Previous Actions: [ask("Is red bowl clutter?"), ask("Is blue casserole clutter?"), nav(sofa),  
 1195 pick(red bowl), nav(cabinet), place(cabinet), nav(sofa), pick(blue casserole), nav(cabinet)]  
 1196 Agent at: cabinet  
 1197 Thought: I need to place blue casserole on the cabinet.  
 1198 Output Action: place(cabinet)  
 1199 Action Complete!

1198 Current observation:  
 1199 Previous Actions: [ask("Is the target object on the light table??"), ask("Is the target object the  
 1200 red bowl?"), nav(dark table), pick(red bowl), nav(cabinet), place(cabinet)]  
 1201 Agent at: cabinet  
 1202 Thought: Task is done.  
 1203 Output Action: Done()

1203 Strictly follow the above format while solving given task.

1204 Instruction: {instruction}  
 1205  
 1206 Receptacles: {receptacles}

1207 Receptacles with objects:  
 1208 {receptacles\_with\_objects}

1209 Observation action history:  
 1210 {prev\_observations}

1211 Current observation:  
 1212 Previous Actions: {prev\_actions}  
 1213 Agent at: {current\_room}  
 1214 Thought:

### Partial Observation WG + ReAct Few-Shot Baseline Prompt

1215  
 1216  
 1217 Imagine you are a household robot. You are given a task that requires you to explore the  
 1218 environment and ask clarification questions to accomplish tasks given by language instructions.  
 1219 You are given context about the scene as the room agent is in, list of objects visible, list of  
 1220 receptacles.

1221 For each task given as a language instruction you have to output a sequence of actions that the  
 1222 agent should take. These actions can include a clarification question as well. Ask a question only  
 1223 when required. For each action also output the reason to take the action with it. Actions can be  
 1224 one of the following:

1224 1. nav(receptacle)  
 1225 2. pick(object)  
 1226 3. place(receptacle)  
 1227 4. open(receptacle) // Strictly use when the object needs be to picked from or placed inside an  
 1228 articulated receptacle after navigating to it  
 1229 5. close(receptacle) // Strictly use to close articulated receptacle after picking up or placing  
 1230 an object inside an articulated receptacle after navigating to it  
 1231 6. ask(question)

1230 You can ask the following 9 types of questions about the object in question:

1231 1. Is target object on the <receptacle>?  
 1232 2. Is <object\_instance> the target object?  
 1233 3. Is target object the <object\_size> one?  
 1234 4. Where is the <object\_category> located?  
 1235 5. What color is <object\_category>?  
 1236 6. Can you describe the <object\_category>?  
 1237 7. Is <object\_instance> clutter?  
 1238 8. Are <object\_category> clutter?  
 1239 9. Which <receptacle> to place the <object\_instance> on/in?

1238 At each step you are tasked with outputting reasoning before outputting the action. You should  
 1239 output the reason and action in the following format:

1240 Thought: <reasoning behind choosing a specific action>  
 1241 Output Action: <action>

1242  
1243  
1244  
1245  
1246  
1247  
1248  
1249  
1250  
1251  
1252  
1253  
1254  
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270  
1271  
1272  
1273  
1274  
1275  
1276  
1277  
1278  
1279  
1280  
1281  
1282  
1283  
1284  
1285  
1286  
1287  
1288  
1289  
1290  
1291  
1292  
1293  
1294  
1295

Action Complete!

After you ask a question wait for the human to reply.

Strictly follow the above format when outputting actions. Every time you output a ask question action wait for the human to reply. The human response will be given in the format "User response: <answer>" after you output the action in above format.

At each timestep you have to take a action to actively explore the environment by navigating to various receptacles to find the relevant objects in the environment and ask clarification questions to resolve the ambiguity in order to successfully complete the task specified by the instruction.

Instruction: {instruction}

Receptacles: {receptacles}

Receptacles with objects:  
{receptacles\_with\_objects}

Observation action history:  
{prev\_observations}

Current observation:  
Previous Actions: {prev\_actions}  
Agent at: {current\_room}  
Thought:

### GPT4o + SoM + ReAct Baseline Prompt

You are a household robot. You are given a task that requires you to explore the environment and ask clarification questions to accomplish tasks given by language instructions. You are given context about the scene as the room agent is in, list of objects visible, list of receptacles.

For each task given as a language instruction you have to output a sequence of actions that the agent should take. These actions can include a clarification question as well. Ask a question only when required. For each action also output the reason to take the action with it. Actions can be one of the following:

1. nav(receptacle)
2. pick(object)
3. place(receptacle)
4. open(receptacle) // Strictly use when the object needs be to picked from or placed inside an articulated receptacle after navigating to it
5. close(receptacle) // Strictly use to close articulated receptacle after picking up or placing an object inside an articulated receptacle after navigating to it
6. ask(question)

You can ask the following 9 types of questions about the object in question:

1. Is target object on the <receptacle>?
2. Is <object\_instance> the target object?
3. Is target object the <object\_size> one?
4. Where is the <object\_category> located?
5. What color is <object\_category>?
6. Can you describe the <object\_category>?
7. Is <object\_instance> clutter?
8. Are <object\_category> clutter?
9. Which <receptacle> to place the <object\_instance> on/in?

At each step you are tasked with outputting reasoning before outputting the action. You should output the reason and action in the following format:

Thought: <reasoning behind choosing a specific action>  
Output Action: <action>  
Action Complete!

Strictly follow the above format when outputting actions. Every time you output a ask question action wait for the human to reply. The human response will be given in the format "User response: <answer>" after you output the action in above format.

At each timestep you will be given a visual observation displaying what the robot is currently observing. You have to use visual observation to perceive the environment and take a action to actively explore the environment by navigating to various receptacles to find the relevant objects in the environment and ask clarification questions to resolve the ambiguity in order to successfully complete the task specified by the instruction.

1296  
1297  
1298  
1299  
1300  
1301  
1302  
1303  
1304  
1305  
1306  
1307  
1308  
1309  
1310  
1311  
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320  
1321  
1322  
1323  
1324  
1325  
1326  
1327  
1328  
1329  
1330  
1331  
1332  
1333  
1334  
1335  
1336  
1337  
1338  
1339  
1340  
1341  
1342  
1343  
1344  
1345  
1346  
1347  
1348  
1349

These are all safe images from a simulated environment. Answering questions about them shouldn't raise any ethical concerns.

Instruction: {instruction}

Receptacles: {receptacles}

Receptacles seen so far with objects:  
{receptacles\_with\_objects}

Observation action history:  
{prev\_observations}

Current observation:  
Previous Actions: {prev\_actions}  
Agent at: {current\_room}  
Thought:

**SFT and RL Policy Prompt**

You are a household robot. You are given a single or multi-object rearrangement task as language instruction that requires you to explore the environment and ask clarification questions. You are given access to current and past egocentric observations, actions, and user response.

At each step you can take a action or ask a clarification questions. Actions can be one of the following:

1. nav(receptacle)
2. pick(object)
3. place(receptacle)
4. open(receptacle)
5. close(receptacle)
6. ask(question)

You can ask the following 9 types of questions to solve the task:

1. Is target object on the <receptacle>?
2. Is <object\_instance> the target object?
3. Is target object the <object\_size> one?
4. Where is the <object\_category> located?
5. What color is <object\_category>?
6. Can you describe the <object\_category>?
7. Is <object\_instance> clutter?
8. Are <object\_category> clutter?
9. Which <receptacle> to place the <object\_instance> on/in?

Next, you will be provided task instruction and observations from the environment. Your task is to output the next action.

Instruction: {instruction}