# Supplementary Material for
# MultiTalk: Introspective and Extrospective Dialogue for Human-Environment-LLM Alignment

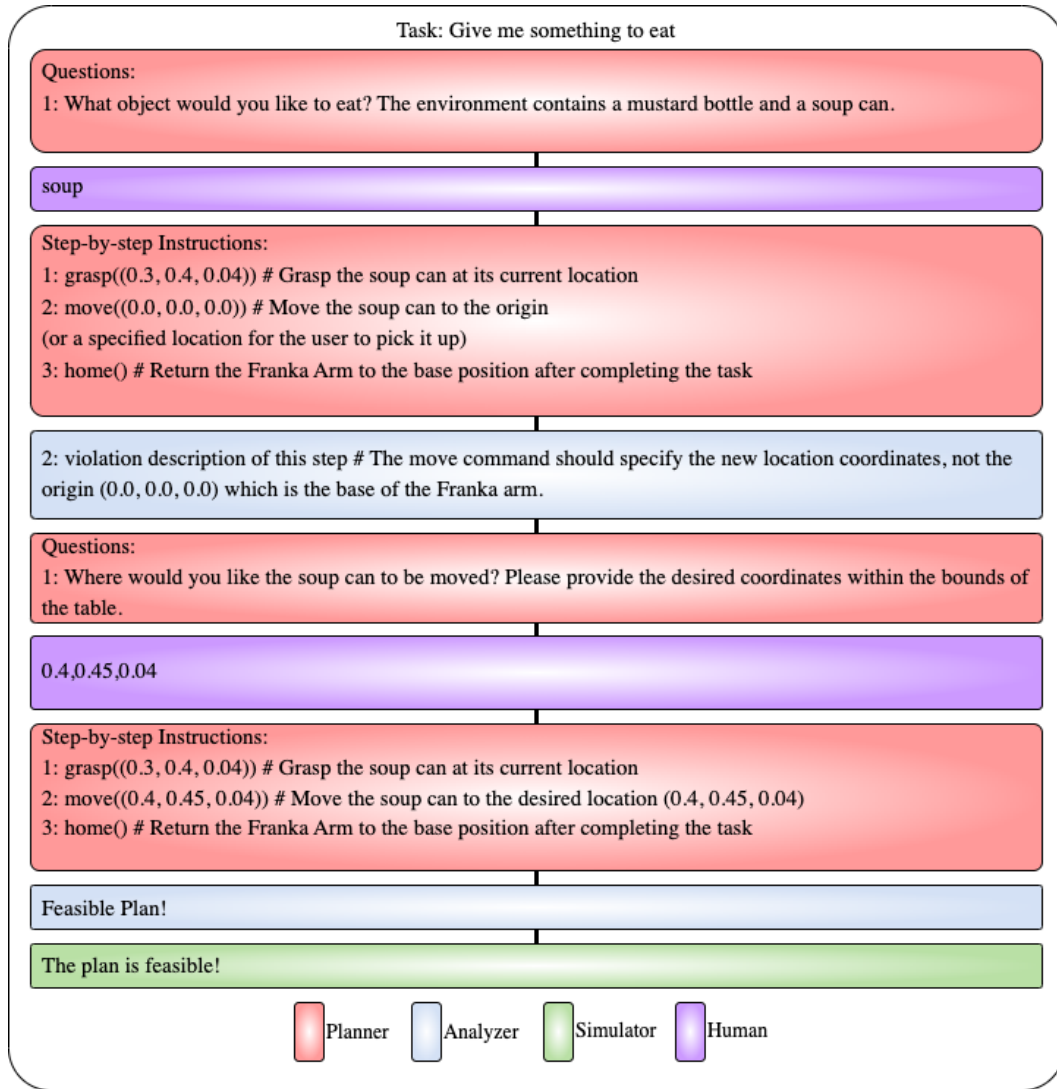## A  Selected Examples of MultiTalk Experiments



Figure 1: An example of the complete pipeline in action to generate a feasible plan. The Planner forgets to ask the target location and the Analyzer reminds it in its feedback. Since there is no execution error in the plan, the Simulator approves the plan in the end.

Task: Interchange the locations of two pairs of objects on the table

1: Which two pairs of objects would you like to interchange the locations of?

The mustard bottles with each other and the soup cans with each other

1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1 at its current location
2: move((0.35, 0.35, 0.04)) # Move mustard bottle 1 to the location of mustard bottle 2
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2 at its current location
4: move((0.32, -0.33, 0.04)) # Move mustard bottle 2 to the location of mustard bottle 1
5: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1 at its current location
6: move((0.4, 0.51, 0.04)) # Move soup can 1 to the location of soup can 2
7: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2 at its current location
8: move((0.38, -0.51, 0.04)) # Move soup can 2 to the location of soup can 1
9: home() # Return the franka arm to base position after completing the task

3: violation description of this step # Attempting to grasp mustard bottle 2 at the location where mustard bottle 1 was moved, which is incorrect.
7: violation description of this step # Attempting to grasp soup can 2 at the location where soup can 1 was moved, which is incorrect.

1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1 at its current location
2: move((0.3, -0.3, 0.04)) # Temporarily move mustard bottle 1 to a valid location on the table
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2 at its current location
4: move((0.32, -0.33, 0.04)) # Move mustard bottle 2 to the location of mustard bottle 1
5: grasp((0.3, -0.3, 0.04)) # Grasp mustard bottle 1 from the temporary location
6: move((0.35, 0.35, 0.04)) # Move mustard bottle 1 to the location of mustard bottle 2
7: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1 at its current location
8: move((0.3, -0.3, 0.04)) # Temporarily move soup can 1 to a valid location on the table
9: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2 at its current location
10: move((0.38, -0.51, 0.04)) # Move soup can 2 to the location of soup can 1
11: grasp((0.3, -0.3, 0.04)) # Grasp soup can 1 from the temporary location
12: move((0.4, 0.51, 0.04)) # Move soup can 1 to the location of soup can 2
13: home() # Return the franka arm to base position after completing the task

Feasible Plan!

Collision detected between mustard bottle 1 and mustard bottle 2

1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1 at its current location
2: move((0.3, -0.6, 0.04)) # Temporarily move mustard bottle 1 to a valid location on the table, far from other objects
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2 at its current location
4: move((0.32, -0.33, 0.04)) # Move mustard bottle 2 to the location of mustard bottle 1
5: grasp((0.3, -0.6, 0.04)) # Grasp mustard bottle 1 from the temporary location
6: move((0.35, 0.35, 0.04)) # Move mustard bottle 1 to the location of mustard bottle 2
7: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1 at its current location
8: move((0.3, -0.6, 0.04)) # Temporarily move soup can 1 to a valid location on the table, far from other objects
9: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2 at its current location
10: move((0.38, -0.51, 0.04)) # Move soup can 2 to the location of soup can 1
11: grasp((0.3, -0.6, 0.04)) # Grasp soup can 1 from the temporary location
12: move((0.4, 0.51, 0.04)) # Move soup can 1 to the location of soup can 2
13: home() # Return the franka arm to base position after completing the task

Feasible Plan!

The plan is feasible!

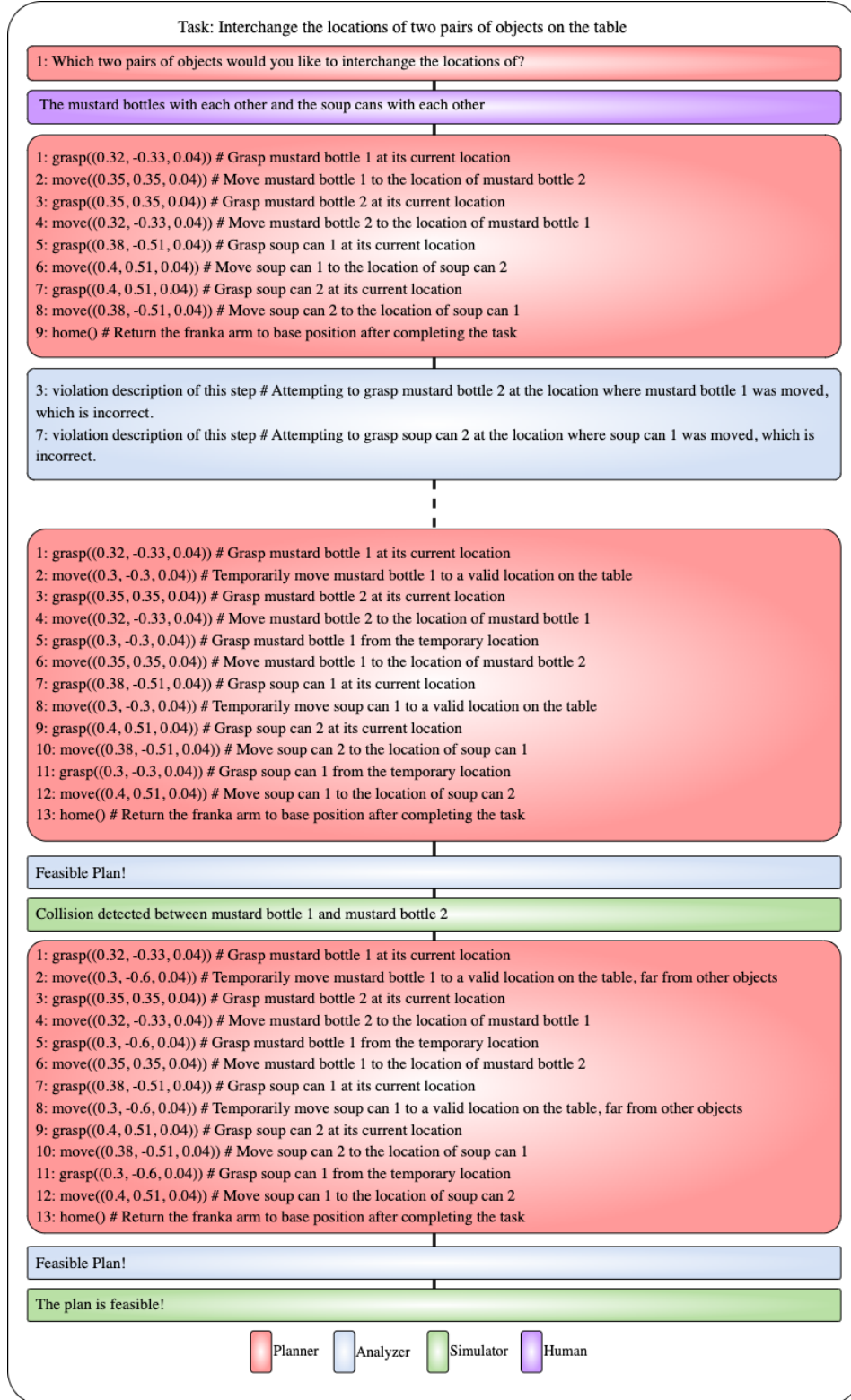Planner    Analyzer    Simulator    Human

Figure 2: An example of the complete pipeline in action to generate a feasible plan. The Analyzer corrects the Planner in the first feedback. In the last part of the execution, the Simulator detects the collision between two objects and the Planner corrects its plan after receiving this feedback.
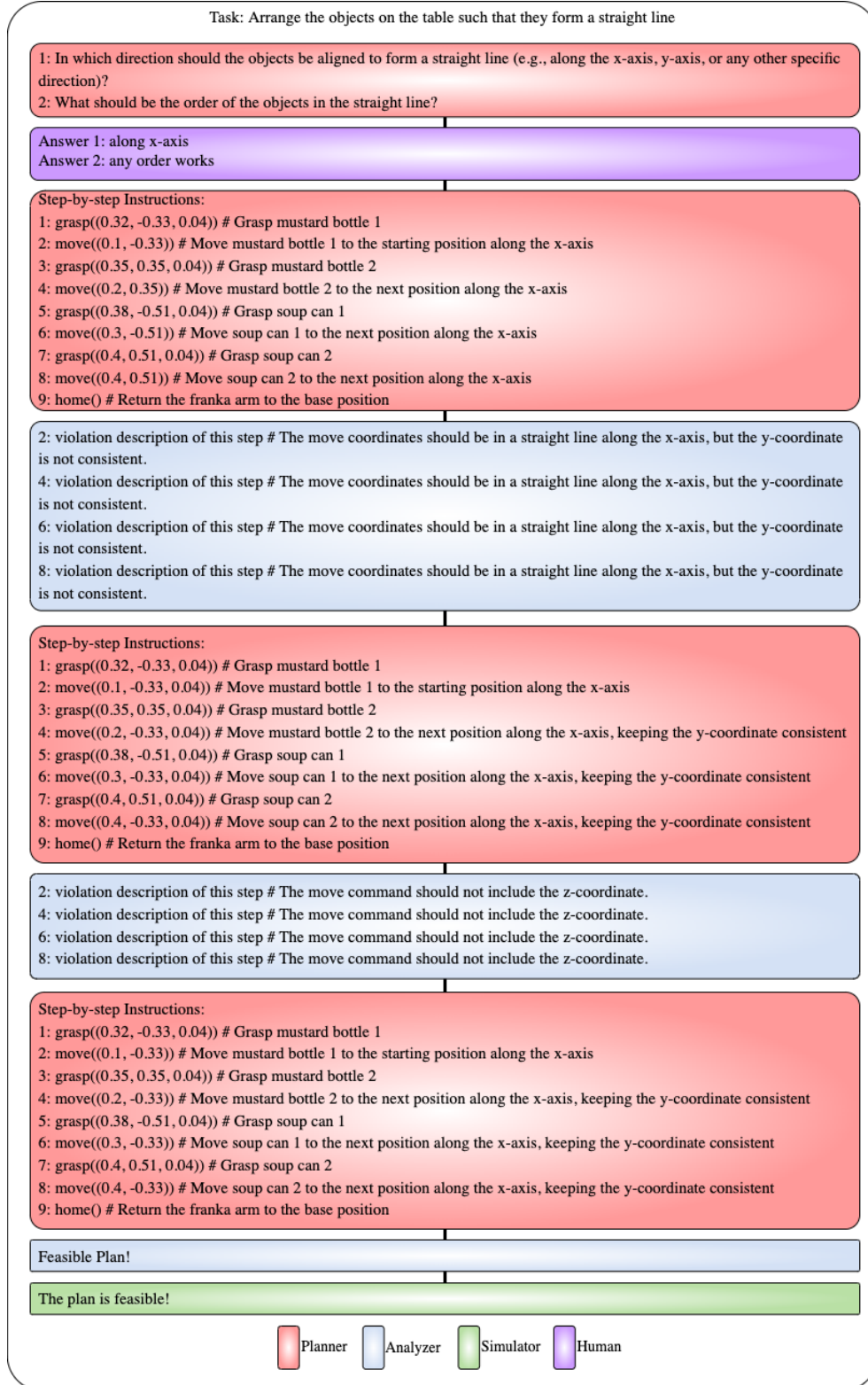
Task: Arrange the objects on the table such that they form a straight line

1: In which direction should the objects be aligned to form a straight line (e.g., along the x-axis, y-axis, or any other specific direction)?
2: What should be the order of the objects in the straight line?

Answer 1: along x-axis
Answer 2: any order works

Step-by-step Instructions:
1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1
2: move((0.1, -0.33)) # Move mustard bottle 1 to the starting position along the x-axis
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2
4: move((0.2, 0.35)) # Move mustard bottle 2 to the next position along the x-axis
5: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1
6: move((0.3, -0.51)) # Move soup can 1 to the next position along the x-axis
7: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2
8: move((0.4, 0.51)) # Move soup can 2 to the next position along the x-axis
9: home() # Return the franka arm to the base position

2: violation description of this step # The move coordinates should be in a straight line along the x-axis, but the y-coordinate is not consistent.
4: violation description of this step # The move coordinates should be in a straight line along the x-axis, but the y-coordinate is not consistent.
6: violation description of this step # The move coordinates should be in a straight line along the x-axis, but the y-coordinate is not consistent.
8: violation description of this step # The move coordinates should be in a straight line along the x-axis, but the y-coordinate is not consistent.

Step-by-step Instructions:
1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1
2: move((0.1, -0.33, 0.04)) # Move mustard bottle 1 to the starting position along the x-axis
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2
4: move((0.2, -0.33, 0.04)) # Move mustard bottle 2 to the next position along the x-axis, keeping the y-coordinate consistent
5: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1
6: move((0.3, -0.33, 0.04)) # Move soup can 1 to the next position along the x-axis, keeping the y-coordinate consistent
7: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2
8: move((0.4, -0.33, 0.04)) # Move soup can 2 to the next position along the x-axis, keeping the y-coordinate consistent
9: home() # Return the franka arm to the base position

2: violation description of this step # The move command should not include the z-coordinate.
4: violation description of this step # The move command should not include the z-coordinate.
6: violation description of this step # The move command should not include the z-coordinate.
8: violation description of this step # The move command should not include the z-coordinate.

Step-by-step Instructions:
1: grasp((0.32, -0.33, 0.04)) # Grasp mustard bottle 1
2: move((0.1, -0.33)) # Move mustard bottle 1 to the starting position along the x-axis
3: grasp((0.35, 0.35, 0.04)) # Grasp mustard bottle 2
4: move((0.2, -0.33)) # Move mustard bottle 2 to the next position along the x-axis, keeping the y-coordinate consistent
5: grasp((0.38, -0.51, 0.04)) # Grasp soup can 1
6: move((0.3, -0.33)) # Move soup can 1 to the next position along the x-axis, keeping the y-coordinate consistent
7: grasp((0.4, 0.51, 0.04)) # Grasp soup can 2
8: move((0.4, -0.33)) # Move soup can 2 to the next position along the x-axis, keeping the y-coordinate consistent
9: home() # Return the franka arm to the base position

Feasible Plan!

The plan is feasible!

Planner    Analyzer    Simulator    Human

Figure 3: An example of the complete pipeline in action to generate a feasible plan. The Analyzer corrects the Planner in the first feedback. Although it hallucinates in its second feedback, it does not affect the executability of the plan, as the z-coordinate is not necessary in our test environment; we use the z-coordinate of the table for all commands during execution.

Figure 4: The **same** example as in Figure 3 but executed by the Planner + Simulator, not by the complete pipeline. Since the incorrect plan is not corrected by the Analyzer as it was in Figure 3, the plan is incorrect although it is executable by the robot.

# B  Implementation of selected examples on FR3



Figure 5: Execution of "*Give me the soup can*" on FR3



Figure 6: Execution of "*Give me something to eat*" on FR3

Figure 7: Execution of "*Give me something to eat*" on FR3



Figure 8: Execution of "*Interchange the object locations in the scene*" on FR3

## C    Prompts to the Planner

The prompts fed to the Planner can be broken down into multiple parts, each conveying different information to the LLM. The following are the exact prompts used in our implementation. The system, actions, environment, and output format are fed to the Planner, in that order, before starting the loop. The feedback prompts are used based on the source of feedback.

**System**

You are in the command of a Franka Arm. Your task is to split a high-level instruction into low-level step-by-step instructions that can be executable by the agent. If there is any ambiguity in the high-level instruction, you should ask for clarification.

## Actions

Necessary and sufficient robot actions are defined as follows:

– move(<desired_coordinates>): Move object to desired location

– grasp(<object_coordinates>): Grab the object at the given location

– home(): return the franka arm to base position after completing the task

Additional Instructions regarding these actions:

– Only one object can be grasped at a time

– An object must be grasped before moving it to another location

– Once the robot returns to home position, no action can be performed

## Output Format

If there is no ambiguity in the high-level task presented to you, your output must be a set of robot actions and their reasoning as numbered bullets with each action and reasoning pair on a new line. For example, if your output has 2 actions, it should look like this:
Step-by-step Instructions:

1: action_1    # reasoning for action_1

2: action_2    # reasoning for action_2

In case of ambiguity, your output must be a set of questions posed to the user that will help you resolve the ambiguity. For example, if you are have two questions, your output should look like this:
Questions:

1: question_1

2: question_2

You can ask questions in the following situations:

– If you are unable to infer the target object/action from the prompt provided.

– If the feedback systems deem that the task cannot be executed, and you need the user to alter it

– If there are multiple ways to interpret a task, go ahead with your interpretation without asking the user. In case your interpretation cannot be executed, then question the user

You must not ask questions in the following situations:

– If the task demands an object that is not present in the environment, you must command the Perception system to rerun and provided the updated environment state rather than asking the user.

To command the perception to rerun, your output should look like this:
"Perception Rerun Required"

## Query

For the following objects, decide if you can give step-by-step instructions or would like to ask questions for clarification or need to command the perception module to re-scan the environment. In either case, use the output prompt format specified earlier.
Enviornment: <objects>
High-level Instruction: <Task>
Output:

# D Prompts to the Analyzer

Similar to the Planner, even the Analyzer's prompts consists of various parts. Apart from the query calls, all other prompts are fed to the LLM before the start of the loop in the same order as enlisted below.

---

**System**

You are a fault finding agent in the command of a Franka arm. Your task is to find faults/violations in a step-by-step plan provided by a planning agent based on a high-level instruction it was given. You are coupled with a visual perception system that reports the observed scene to you. Further, you are provided with information on certain restrictions in the feasible actions and a set of rules that must be followed by the presented plan. Using these, you must specify the violations, if any, made by one or more steps of the plan. Along with this, you will also act as a critic to the planning agent by judging whether the plan generated aligns with the high-level instruction provided. The planning agent will then revise the plan based on your feedback and resubmit it to you for further validation. You must continue this process until the plan is free of any violations.

---

**Input Format**

The Instructions generated by the planning agent are a set of robot actions as numbered bullets with each action on a new line. For example, if the planner instruction has 2 actions, it would look like this:
Step-by-step Instructions:

1: action_1

2: action_2

Necessary and sufficient robot actions are defined as follows:

– move(<desired_coordinates>): Move an object from its current location to the new location defined by <desired_coordinates>

– grasp(<object_coordinates>): Grab an object at its given location defined by <object_coordinates>

– home(): return the franka arm to base position after completing the task

---

**Environment**

The franka arm's base is at origin. The environment is a 2D square top table.
Format of the observed environment as a python dictionary:
`{<object_name>:((<bounding box length>, <bounding box width>, <bounding box height>), <bounding box center>) }`
Observed semantic relations:
`[(<object_1>, <object_2>), <relationship>),(<object_3>, <object_4>), <relationship>), ...]`
Necessary and sufficient object relationships:

– (<object_1>, <object_2>), 'on': object_2 is on object_1

– (<object_1>, <object_2>), 'in': object_2 is in object_1

– (<object_1>, <object_2>), 'near': object_2 is near object_1

## Critic

The planner is given a high level task to generate step-by-step instructions. Along with that, it is also given the opportunity to ask the user questions to clarify the task. You will be provided with both the task and answers to the questions to do your job of a critic. Moreover, if there is an ambiguity in the task and the planner have not resolve this by asking a question to the user, you should warn the planner.

The types of ambiguities include but not limited to the following cases:

– The target location is not specified by the user but directly assumed by the planner.

– The object to be used is not specified by the user but directly assumed by the planner.

– There are more than one possible interpretation of the task but planner has not asked any question to resolve it.

The format of the questions and answers will be as follows:

1: Question 1

2: Answer 1

3: Question 2

4: Answer 2

. . .

If no question was asked by the planner, you will be given 'None'.

## Restrictions

List of rules that must not be violated by any action:

– Each action must be physically feasible

– Each action should have the right arguments, i.e., it should have only the object coordinates

– For the grasp command, also check if the mentioned coordinates actually belong to an object in the environment. Do not check this for the move command.

– The actions should not lead to geometric, physical, or structural conflict in the environment

– Remember that you are providing feedback to a planning agent only. If it chooses to stick to the same plan despite your feedback, you should trust it and accept the plan.

– You are not responsible to detect out of bounds errors.

## Output Format

Your output must be pairs of step number, its violation or infeasibility. If a step has no issues, write 'no violations' for that step. Do not write anything else in the output. For example, if step 3 and 7 violate some rules or constraints, the output must look like this:
Step-by-step Instructions:

3: violation description of this step # Details

7: violation description of this step # Details

Necessary and sufficient robot actions are defined as follows: If all the steps are correct, the output strictly must be the phrase 'Feasible Plan!' only.

**Query for the first call**

High-level instruction given to the planner: <task>
Planner's output:
<plan> Further, the planner also asked the user some questions to clarify the plan. The questions and the user's responses are as follows:
<question_and_answers>
Observed environment:
<objects> Observed semantic relations: "' <semantic_relations> "' The output is:

**Query for subsequent calls**

Based on the feedback you provided, the planning agent has updated its instructions. The new instructions are as follows:
Further, the planner also asked the user some questions to clarify the plan. The questions and the user's responses are as follows:
Check for violation is this updated plan and generate feedback for the same. After your feedback, there might be still problems in the plan. Use the same output format as defined previously. The new feedback is: