

# Lifelong Robot Learning with Human Assisted Language Planners

Meenal Parakh<sup>\*, $\alpha,\gamma$</sup> , Alisha Fong<sup>\*, $\alpha,\gamma$</sup> , Anthony Simeonov <sup>$\alpha,\gamma$</sup> , Tao Chen <sup>$\alpha,\gamma$</sup> , Abhishek Gupta <sup>$\alpha,\beta,\gamma$</sup> , Pulkit Agrawal <sup>$\alpha,\gamma$</sup>

<sup>$\alpha$</sup> Improbable AI Lab

<sup>$\beta$</sup> University of Washington

<sup>$\gamma$</sup> Massachusetts Institute of Technology

\*Authors contributed equally

**Abstract**—Large Language Models (LLMs) have been shown to act like planners that can decompose high-level instructions into a sequence of executable instructions. However, current LLM-based planners are only able to operate with a fixed set of skills. We overcome this critical limitation and present a method for using LLM-based planners to query new skills and teach robots these skills in a data and time-efficient manner for rigid object manipulation. Our system can re-use newly acquired skills for future tasks, demonstrating the potential of open world and lifelong learning. We evaluate the proposed framework on multiple tasks in simulation and the real world. Videos are available at: <https://sites.google.com/mit.edu/halp-robot-learning>

## I. INTRODUCTION

A dream shared by many roboticists is to instruct robots using simple language commands such as “clean up the sink.” Large language models (LLMs) can support this dream by decomposing an abstract task into a sequence of executable actions or “skills” [15]. Several LLM-based works use a *fixed* set of skills (i.e., *skill library*) for planning [1, 14]. However, the available skills may not suffice in certain task scenarios. For instance, given the task, “clean up the sink”, an LLM may plan a sequence of picks and places that move all the dishes to a dishrack. Suppose one cup contains water which must be emptied before the robot puts it away. Without access to an “empty cup” skill, the system is fundamentally incapable of achieving this task variation. On detecting failure, LLM planners may attempt to expand their abilities – the system could *request* a new skill for “pouring” if it detects water in the cup. However, unless the robot can also *execute* new skills, the problem remains unsolved.

Based on the tasks and scenarios the robot encounters, the planner must have the capacity to request and acquire *new skills*. Further, such skill acquisition ought to be *quick* – a system that requires days, weeks, or months to acquire the new skill is of little utility. Concurrent to our work, the ability of an LLM-based planner to acquire new skills has been demonstrated in the virtual domain of Minecraft [35]. However, in virtual domains, new skills can be simply represented as code that can execute high-level and abstract actions. In contrast, learning a new skill for a robot also involves finding low-level actions that can affect the physical world. To the best of our knowledge, the ability to add skills to a skill library in a time and data-efficient manner and utilize them for future tasks, especially in the context of LLM-based planners, has not been demonstrated.

Existing LLM-based robotic systems struggle with online skill acquisition because common mechanisms for learning

skills (e.g., end-to-end behavior cloning or reinforcement learning) typically require a large amount of data and/or training time. Some methods are able to acquire new skills in a more data-efficient manner in limited scenarios such as in-plane manipulation (e.g., TransporterNets [38]), but these skills are insufficient for 6-DoF actions (e.g., “grasp the mug from the side”, “hang the mug on a rack” or “stack a book in a bookshelf”). Another body of work such as in few-shot imitation learning can efficiently solve new instances from a task family but requires large amounts of pre-training data [10, 26] which is seldom available for new skills. We first present a method that allows LLMs to request new skills to complete the given task. Second, we propose to use Neural Descriptor Fields (NDFs) [30] to realize these new skills. We choose NDFs as they require only 5-10 demonstrations to perform rigid body manipulation in the full space of 3D translations and rotations.

Our system works by prompting an LLM with a textual scene description obtained by a perception system, a library of skills expressed as Python functions, and a natural language task specification. With this information, the LLM plans and produces a sequence of skills (in the form of code) that achieves the task. Along with the skills in the skill library, we also provide the LLM with a special function for requesting a new skill to be added to the library. When the LLM plans call this `learn_skill` function, it returns a new skill name and a docstring description of the skill. However, such a skill is abstract and is not mapped to actions. NDFs allow the user to quickly realize this new skill by providing a few demonstrations, after which the skill is added to the skill library so that it can be re-used on future tasks. In summary, this work demonstrates a proof-of-concept implementation of an LM-powered robotic planning agent that can interactively grow its skill library based on the needs of the task. We show an instance of such a system using NDFs and perform experiments that highlight the abilities of our system.

## II. RELATED WORK

a) *LLMs as Zero-Shot Planners*: Prior work that uses large language models (LLMs) as planners include SayCan [1], InnerMonologue [14], NLMap-SayCan [4] and Socratic Models [37]. These methods make significant contributions: [1] and [37] using LLMs as planners; [14] emphasizes the importance of feedback; and [4] improves upon [1] by introducing the ability of open-vocabulary detection for grounding using CLIP and ViLD features [27] [13]. The planners in these methods either generate the plan in textual

format or choose the next step based on a given set of actions described through text. Another set of methods [21] [34] [32] [22] using LLM as planners chose to output the plans directly using a Python or symbolic API, given the function documentation and sufficiently expressive function names.

*b) End-to-End Language Conditioned Manipulation:*

Another class of methods processes inputs from different modalities such as visual, textual, and sound, and train an LLM to use these inputs to output robot actions end-to-end (e.g., CLIPort [28], Interactive Language [23], RT1 [2], PerAct [29] and VIMA [17]). Another end-to-end approach is Palm-e [9] that generates textual steps as output, and are assumed to map to a small set of low level policies. One main advantage these offer is more faithful LLM grounding, in contrast to modular approaches that only list the objects in the scene and sometimes fail due to partial scene descriptions. However, they each suffer from requiring a large amount of data for training or fine-tuning. Such large data requirements also make it difficult to achieve generalization. Finally, many of these works are limited to performing 3-DoF (top-down) manipulation actions.

*c) Low-Level Robot Primitives.:* The modular approaches [1] [14] [4] [37] [21] [34] use a predefined set of primitive skills, often hardcoded or learned from behavior cloning. These low-level primitives can also be learned through methods such as [16], [11], [6], [38]. While these skills can be composed to perform a wide range of actions, many times a required skill cannot be composed from the primitive set and adding a new primitive may require careful engineering, or large number of demonstrations. Thus, we employ [30] to incorporate new skills at runtime using only a few demonstrations, with the only drawback of limiting the skills to known object categories.

### III. METHOD

In the spirit of prior work on performing long-horizon tasks wherein a high-level planning algorithm chains together different low-level skills [12, 20, 24, 37], our system has explicit modules for perception, planning, and control (Fig. 1). The modularity of our system allows us to take advantage of state-of-the-art (SOTA) models like SAM [18] for segmentation and GPT-4 [25] for planning skill sequences. At a high level, our perception module describes the scene from RGB and depth observations, generating a language-based scene description containing information about the objects in the scene and the spatial relationship between them. Given the scene description and a library of skills, the planning module plans a sequence of steps to solve the task based on the scene description and task requirements. The skill sequence corresponds to a set of executable behaviors on the robot.

In contrast to previous work that uses LLMs in robotics, our planning module can request to learn a new skill when it determines that the existing skills are insufficient, and a data-efficient skill learning method can be used to extend the skill library with this new executable behavior. With an expanded skill library, the planner can utilize both the original primitive skills and the newly learned skills when

completing subsequent tasks. Thus, our approach endows the system with a form of continual learning. In the following subsections, we describe each module in detail.

#### A. Perception

The perception module (Fig. 2a) processes RGBD images to obtain and store information about the scene objects. First, the module identifies objects using an open-vocabulary object detector [39]. We also perform segmentation to obtain object masks using SAM [18] and combine them with the depth images to obtain object point clouds. In addition to object labels and segmentation masks, the planner may require additional information about the spatial arrangement of the scene. For example, if a robot needs to empty a mug, it first needs to know whether there is an object *in* the mug, and only execute the skill of emptying it if there is. We generate spatially-grounded scene descriptions automatically by computing positional relationships between objects using the object point clouds. A scene description that is not spatially grounded only describes the objects present in the scene, without specifying the spatial relationship between them. Lastly, to enable open-vocabulary language commands that target specific object instances, we extract CLIP embeddings of each segmented object in the scene. In this way, given a scene with multiple mugs, if the task is to “pick up the red mug,” we are able to identify the object that corresponds to the description of a “red mug” (additional examples in Appendix). Overall, our perception components output segmented object point clouds with associated detection labels, inter-object relations, and CLIP embeddings.

*a) Spatially-grounded Textual Scene Description:* To inform the planner about the environment state, we format the perception outputs into a language-based scene description with information about the scene objects and their inter-object relations. This involves constructing a string with the names of the objects along with the relations that hold between them. The description is akin to a textual description of “scene graph”. Please see Appendix for further details. Note that the particular method of describing the scene is not critical to our work and in the future vision-language models capable of describing objects and the relationship between them can replace this system.

#### B. Planning and Control

Given the language command and the textual scene description from the perception system, GPT-4 is used to plan a sequence of the steps to be executed. The inputs and outputs of the LLM are structured as follows:

*a) Skill Definitions via Code API:* One way to design a planner is to output a plan in natural language. However, a more machine-friendly alternative is to have the planner output programming code [21, 32]. Having an LLM planner directly produce code avoids the need to map a textual plan to a robot-executable plan. In addition, communicating with LLM in a programming language allows a human to give prompts in the form of comments, docstrings, and usage examples, which helps the planner understand how each skill

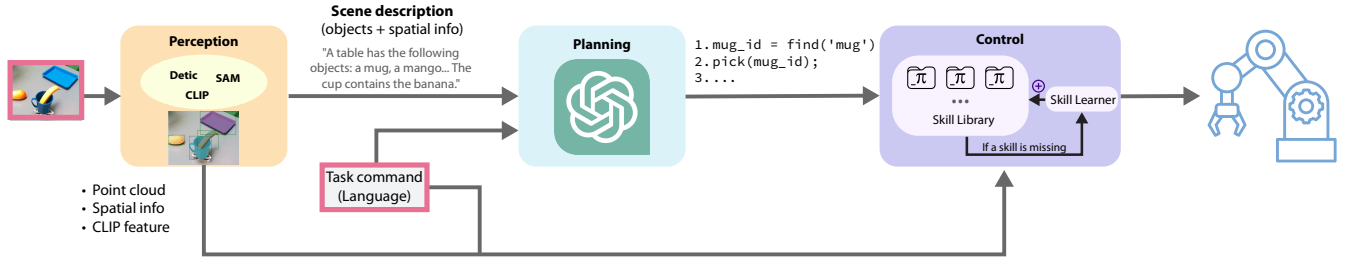


Fig. 1: Our system consists of three modules: *perception*, *planning*, and *control*. The *perception* module processes RGB-D images and outputs a textual scene description that identifies objects and their spatial relationships. The *planning* module uses GPT-4 to plan a sequence of steps based on the available skills and the task command. We added a `learn_skill(skill_name)` function to the planner so that it can plan to learn a new skill if such learning is necessary for completing the task. Finally, the *control* module executes the planned steps using the available skills or starts learning a new skill.



Fig. 2: (a) From RGB-D images, our perception module obtains information about the objects and their relations, creates an object information dictionary, and generates a scene description (detection, object pairs corresponding to given object relations, and the template is in black). (b) An example showing the interaction between the robot, the user, and the planner.

operates. To take advantage of these benefits, we define each skill as a Python function that takes input arguments such as object identifiers and environment locations. We provide the planner with a description and set of input/output examples for each function. The code API is initialized with a skill library  $S_0$  containing five primary functions: `find`, `pick`, `get_place_position`, `place`, and `learn_skill`:

- `find(object_label=None, visual_description=None, location_description=None)`: searches with the perception system for an object based on category,

visual property, or location. Returns an `object-id`.

- `pick(object_id)`: uses Contact-GraspNet [33] to find a 6-DoF grasp for the object point cloud associated with the `object_id` and executes the grasp.
- `get_place_position(object_id, reference_id, relation)`: for the object given by `object_id`, returns the  $(x, y, z)$  location determined by the text description `relation` relative to `reference_id`.
- `place(object_id, place_position)`: places the object at the  $(x, y, z)$  value given in `place_position`.
- `learn_skill(skill_name)`: returns a new executable skill function and a docstring describing the skill behavior.

The above API functions also output a signal indicating whether or not the function executes properly (i.e., to catch and correct runtime errors due to syntax mistakes). If new skills are learned (discussed in Sec. III-C), the library is updated  $S_i = S_{i-1} \cup \{\pi_i\}$  where  $\pi_i$  denotes the new skill.

*b) Full Planner Input/Output and Skill Execution:* The planner is prompted to produce the plan in two steps. First, given the scene and task description, the planner generates a sequence of steps described in natural language. Next, the planner is provided with the code API of skills as discussed above and tasked to write code for executing the task using the given skills. For example, if the first step in the plan is to “find” a mug with the `find` function, the planner may output `object_id = find("mug")`. Since our system uses a LLM planner, the human user can interact with the planner at either stage of the planning to further refine the plan or correct mistakes. An example of the interaction between the user, planner, and robot is shown in Fig. 2. We qualitatively observe this two-step process helps the model generate higher-quality plans, as compared to producing the full plan directly. The two-step breakdown potentially helps in the same way “chain-of-thought” prompting has helped LLM find better responses [36].

The code returned by the LLM is executed using the `exec` construct in Python. For skills involving robot actions, the skill function calls a combination of inverse kinematics (IK), motion planning, and trajectory following using a joint-level PD controller.

### C. Learning New Skills and Expanding the Skill Library

a) *Requesting New Abilities with `learn_skill` function:* The code API for the `learn_skill`, contains a docstring detailing the role of the function and also includes a few examples of the desired output of using the `learn_skill` function. The reason for providing examples is to exploit the in-context learning ability of LLMs – these examples help the LLM figure out how to use the `learn_skill` function. More details are in the Appendix. The `learn_skill(skill_name)` returns the handle to a new executable skill function along with a docstring that describes the behavior of the function. The function is parameterized by either one or two `object_ids` - one for specifying which object `skill_name` acts upon, and another for specifying a reference object for relational skills (e.g., `pick(bottle_id)` vs. `insert(peg_id, hole_id)`). The exact parameterization is decided by the LLM. When `learn_skill` is called, the returned function is added to the skill library so that the new skill can be reused in the future.

b) *Data- and time-efficient skill grounding with NDFs:* Our framework is agnostic to the specific method used to ground newly learned skills into actions. It can be end-to-end learning with reinforcement learning, or behavior cloning from demonstrations. In this work, we choose to use NDFs [30] to learn new skills because it allows efficiently learning a skill from just a few ( $\leq 10$ ) demonstrations. NDFs also facilitate a degree of category-level generalization across novel object instances, as well as generalization to novel object poses due to built-in rotation equivariance. More information on NDFs can be found in [30, 31].

c) *Learning from Feedback:* If we specify a task the system cannot solve using the available skills (such as “pick up the mug by the handle”, when the available “pick” skill grasps the mug from the rim), we would expect the LLM to directly request a new skill with `learn_skill`. While this occurs the majority of the time (see Experiments Section), the planner sometimes directly attempts the task using a skill that does not satisfy the task requirements. In these cases, if a user provides the *outcome* of a task attempt (e.g., “the mug was grasped by the rim”), the planner can use this information to register its usage of an incorrect skill and subsequently call `learn_skill` to expand its abilities. The system can then attempt the task with the newly learned skill.

This highlights the need for feedback mechanisms that, in addition to detecting runtime errors, also inform the planner about the state of the environment after skill execution. To achieve this, we allow a human operator to manually but *optionally* provide feedback before and after code execution. We allow the human to provide feedback after the execution of every step in the code. The combination of *outcome* feedback from the user and the *execution* feedback from the skill functions enables the system to detect failures, replan and if necessary expand its skillset using `learn_skill`.

d) *Continual Learning:* Learning new skills allows one to execute a task that was previously not possible. However, the full potential of learning new skills is realized when we

allow the system to *continually* acquire and *re-use* skills to solve future tasks. This creates a system with ever-expanding capabilities. There are many ways this can be achieved – our implementation involves simply adding a new skill function expressed as a code API to the skill library, and using the updated library for future tasks.

## IV. EXPERIMENTS

a) *Environment Design and Setup:* We design our experiments to achieve three goals: (1) Show a proof-of-concept implementation of LLM-based task planning and execution with interactive skill learning in the real world, (2) Evaluate the abilities of current LLMs to appropriately request and re-use new skills based on the needs of different manipulation tasks, and (3) Compare the performance of the system when different components (such as object relations) are included vs. removed.

In the real world, we tested our system on the Franka Panda robot with a Robotiq 2F-140 parallel jaw gripper. We used four calibrated RealSense cameras to obtain RGB-D images and point clouds. We also evaluated the LLM planner in isolation with a set of manually crafted tasks, scene descriptions, and success criteria. To perform additional system ablations, we evaluate our approach in simulation using PyBullet [7] and the AI Robot library [5]. Our environment includes a tabletop-mounted Panda with the default gripper, and synthetic cameras for obtaining RGB-D images and segmentation masks. We use a combination of ShapeNet [3] and manually-generated objects for experiments in simulation.

### A. Real-world tasks requiring `learn_skill`

We first showcase the benefits of incorporating `learn_skill`. The system is deployed to perform three tasks in the real world: (1) grasping a mug by a specific part, such as the handle, (2) placing a bottle in a container that must fit on a small shelf, and (3) emptying a mug from a “sink”. Each task can be completed in multiple ways, some of which do not fulfill the full set of task requirements. Our reference point for comparison is the overall system with no feedback mechanism and no `learn_skill` capability. This version directly attempts each task using the base set of primitive skills. Below, we discuss the differences between this baseline and the full version of our system. The full set of planner inputs/outputs for these tasks can be found in the Appendix.

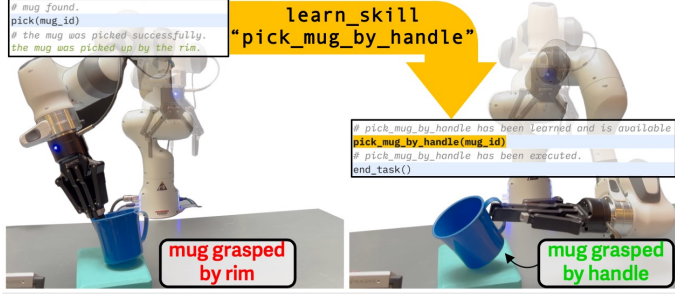
#### 1) Learning and requesting new pick and pick-place skills:

**Task 1: Grasp mug by handle** Our warm-up task that highlights how learning new skills can benefit our system is to perform grasping by a specific part. In this case, we ask the system to “grasp the mug by the handle” (see Fig. 3A). Without `learn_skill`, the planner directly calls `pick` on the mug. This triggers a grasp detector [33] to output a set of grasps on the corresponding mug point cloud. Since most of these grasps are along the rim of the mug, the robot executes a grasp along the rim of the mug, and the task finishes.

If an incorrect skill is used, the human can prompt the system with feedback. By telling the system “the mug was



### (A) Task: Grasp the mug by the handle



### (B) Task: Place bottle in container on its side

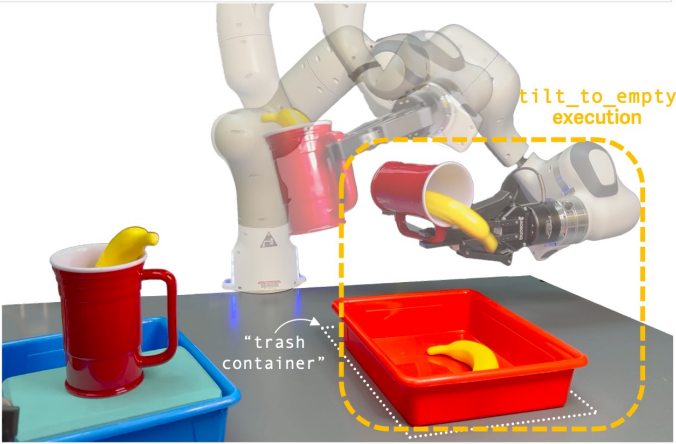


### (C) Task: Empty the sink

# A table has the following objects: a sink, a mug, a trash container, and a banana. At the right of all the objects on the table lies the sink. The sink contains the mug. The mug contains the banana. To the left of the sink lies a trash container.

#### Initial Output:

```
start_task()
sink_id = find(object_label='sink')
mug_id = find(place_description='mug that is inside the sink')
# sink found. mug found.
```



#### Output without Learning

```
pick(mug_id)
# sink found. mug that is inside the sink found. the mug was picked successfully
banana_id = find(place_description='banana that is inside the mug')
table_position = get_place_position(mug_id, sink_id, "above")
place(mug_id, table_position)
# banana that is inside the mug was found. placing mug was successful.
Task was terminated early: the banana is still in the mug
```

#### Output with Learning

```
grasp_mug_by_handle(mug_id)
# grasp_mug_by_handle executed successfully.
tilt_to_empty = learn_skill("tilt_to_empty")
# tilt_to_empty has been learned and is available to use.
trash_container_id = find(object_label='trash container')
tilt_to_empty(mug_id, trash_container_id)
# trash container found. tilt_to_empty has been executed
place_position = get_place_position(mug_id, sink_id, "inside")
place(mug_id, place_position)
# placing mug was successful.
the sink is not empty
pick(mug_id)
# the mug was picked successfully.
put the mug to the right of the sink
place_position = get_place_position(mug_id, sink_id, "to the right")
place(mug_id, place_position)
# placing mug was successful.
```

Fig. 3: High-level plan and images for three tasks requiring a new skill: (A) Grasp mug by the handle, (B) Place bottle in container on its side, and (C) Empty the sink. The gray comments represent execution feedback while the green text is human feedback. When `learn_skill` is not available, the robot fails to complete the tasks. However, by learning new skills, the planner expands its abilities and satisfies each task requirement.

picked up by the rim”, the planner puts the mug back down and requests to learn a new `pick_mug_by_handle` skill. We teach this as a side-grasp at the handle using NDFs with five demonstrations. After collecting the demos, we add `pick_mug_by_handle` to the skill library. Finally, the LLM directly calls `pick_mug_by_handle` and finishes the task successfully.

**Task 2: Place bottle in flat tray** Our next task is to place a bottle in a container that must eventually fit in a small shelf. Here, we prompt the system to “place the bottle sideways in the container” (see Fig. 3B). When the pipeline runs using the base set of skills, the robot uses the only available “place” skill, which places the bottle upright in the tray.

Instead, when we provide the feedback “the bottle was placed upright in the tray”, the LLM calls `learn_skill` to acquire a `place_bottle_sideways_in_tray` skill. This is implemented via NDFs as a side grasp on the bottle along with a reorientation and placement inside of the tray. Once this new skill has been added, the robot is able to successfully complete the task.

2) *Continual learning by re-using previously-learned skills:* **Task 3: Empty mug from sink** Finally, we prompt the system with the abstract objective of emptying a “sink” by removing a mug from the container and placing it on the table (see Fig. 3C). This task implicitly requires *emptying* the mug before placing it. We test the LLM’s ability to satisfy this requirement by placing an additional small object (banana) inside the mug (ensuring the object is at least visible by the cameras, but difficult to pick up directly). The baseline system directly calls a combination of `pick` on the mug and `place` to put the mug down on the table.

However, with access to `learn_skill` and the dynamic skill library, the planner *reuses* `pick_mug_by_handle` learned in Task 1 and immediately requests to learn `tilt_mug` so it can first move any objects in the mug to the trash container. We again use NDFs to teach `tilt_mug`, which reorients the mug above the tray. After emptying, the system plans to place the mug back *into* the sink. The user tells the system “the sink is not empty, put the mug to

the right of the sink”. Finally, the LLM re-plans with this feedback and achieves the final placement on the table.

### B. LLM-only skill learning evaluation

In this section, we examine the isolated ability of the LLM-planner to utilize the `learn_skill` function and to appropriately re-use and/or *not* re-use newly-learned skills on subsequent runs. This enables further analysis of GPT-4’s ability to interpret manipulation scenarios represented via textual scene descriptions and correctly use the available skills provided in the code API. For each task in the following subsections, we provide a manually-constructed scene description (that does not correspond to any particular real-world scene) along with a task prompt and the skill API. We ask the planner to output code that completes the task using the API functions. The code output is manually evaluated as correct/incorrect by a human.

**Requesting new skills when needed** First, we study the ability to either (i) properly call `learn_skill` or (ii) properly *not* call `learn_skill`, for a variety of tasks where either the base skill set is (ii) or is not (i) insufficient for the task, respectively. We report the fraction of attempts that correctly use or ignore `learn_skill` in a scenario where human feedback is not provided. The results are shown in the top two sections of Table I. The 91% success rate for using `learn_skill` without feedback indicates GPT-4 can be used for requesting an expanded skill set in a purely feed-forward fashion. Similarly, the LLM usually does not call `learn_skill` when it is not needed (87% success). However, some performance gap remains in both settings.

**Re-using new skills with varying level-of-detail skill descriptions** Next, we focus on the ability to properly re-use the previously-learned skills on subsequent runs, when they can either be applied or when they specifically should *not* be applied (e.g., in scenarios where they are inappropriate or infeasible). We consider varying levels of detail in the description that accompanies the newly-learned skill as it is added to the code API. For instance, we can provide minimal information and only add the name of the new skill, or we can modify the return values of `learn_skill` so that the LLM writes its own docstring/function description to accompany the new skill when we add it to the API. The results are shown in the last two rows of Table I. The success rates indicate that the language model correctly uses the newly-learned skills with higher frequency when the skill descriptions also include docstrings. This makes intuitive sense, as it provides extra context for both the ability and applicability of the newly learned skill, which the LLM can attend to when generating the output code for executing the task (mimicking the chain-of-thought and “let’s think step-by-step” improvements observed in prior work [19, 36]).

Despite the performance increase when describing newly-added skills in more detail, the LLM only achieves moderate overall performance (75% success rate). We observe this is due to a combination of sometimes using new skills when they should not be used (e.g., calling a `side_pick_bottle` skill even when the scene description

Eval Metric	Variation	Success Rate
Correct use of <code>learn_skill</code>	–	0.91
Correctly did <i>not</i> use <code>learn_skill</code>	–	0.87
Correct re-use of new skill (varying skill description)	Name only Name + docstring	0.50 0.75

TABLE I: Success rates for evaluations LLM-only `learn_skill` evaluation.

says “the bottle *cannot* be reached from the side”) and re-learning the same skill multiple times (while occasionally calling it a very similar name) rather than directly utilizing the function that is already available in the API. We deem this as a somewhat negative result which points to potential gaps in such a method of LLM-based task planning. Namely, directly outputting a sequence of high-level skills (or exhaustively scoring them with a language model) does not allow more information about the operation of high-level skills (such as scenarios when they are or are *not* applicable) to be provided or utilized during planning/reasoning.

## V. LIMITATIONS

While our system takes advantage of SOTA components, they sometimes fail and trigger compounding inaccuracies in the downstream pipeline. For example, the LLM heavily depends on an accurate description of the scene, which can sometimes contain erroneous detections and incorrect object relations. We also leverage human feedback to obtain environment descriptions that inform task success and skill acquisition. Humans can provide accurate descriptions that inform when to learn new skills, but repeated user interaction makes the system less autonomous and slower to execute tasks. Leveraging learned success detectors would make the system more autonomous and self-sufficient. Similarly, human verification is typically needed to confirm the overall success or failure of a task, making it difficult to run system evaluation experiments at scale and limiting our evaluations primarily to qualitative demonstrations.

## VI. CONCLUSION

This paper presents a modular system for achieving high-level tasks specified via natural language. Our framework can actively request and learn new manipulation capabilities, leading to an ever-expanding set of available skills to use during planning. We show how an LLM planner can use this ability to adapt its skill set to the demands of real-life task scenarios via both feed-forward reasoning and environmental feedback. In conjunction with perceptual scene representations obtained from off-the-shelf components and a data-efficient method for learning 6-DoF manipulation skills, we provide an example of a complete system. Our results demonstrate how this combination of full-stack modularity, spatially-grounded scene description, and online learning enables a qualitatively improved ability to perform manipulation tasks specified at a high level.

## VII. ACKNOWLEDGEMENT

We thank the members of Improbable AI for their feedback on the project. This work is supported by Sony, Amazon Robotics Research Award, and MIT-IBM Watson AI Lab. Anthony Simeonov is supported in part by an NSF Graduate Research Fellowship.

### A. Author Contributions

**Meenal Parakh** Co-led the project, developed the core LLM-planning framework and full-stack system, set up and ran experiments in simulation and the real world, and drafted the paper.

**Alisha Fong** Co-led the project, integrated NDF-based skill learning into the LLM-planning framework, set up and conducted experiments in the real world and simulation, helped evaluate the LLM in isolation, and drafted the paper. **Anthony Simeonov** helped integrate NDF-based skills into the framework, supported real robot experiments and LLM-only evaluation, and helped revise the paper.

**Tao Chen** engaged in brainstorming and discussion about system implementation and experiment design, mentored Meenal Parakh, and helped draft the paper.

**Abhishek Gupta** was involved with technical discussions, advised Meenal Parakh, and helped with project brainstorming in the early phases.

**Pulkit Agrawal** advised the project and facilitated technical discussions throughout, helped refine the project focus on interactive skill learning with LLMs, and revised the paper.

## REFERENCES

- [1] Michael Ahn et al. “Do As I Can and Not As I Say: Grounding Language in Robotic Affordances”. In: *arXiv preprint arXiv:2204.01691*. 2022.
- [2] Anthony Brohan et al. “Rt-1: Robotics transformer for real-world control at scale”. In: *arXiv preprint arXiv:2212.06817* (2022).
- [3] Angel X. Chang et al. *ShapeNet: An Information-Rich 3D Model Repository*. Tech. rep. arXiv:1512.03012 [cs.GR]. Stanford University — Princeton University — Toyota Technological Institute at Chicago, 2015.
- [4] Boyuan Chen et al. “Open-vocabulary queryable scene representations for real world planning”. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [5] Tao Chen, Anthony Simeonov, and Pulkit Agrawal. *AI Robot*. <https://github.com/Improbable-AI/airobot>. 2019.
- [6] Cheng Chi et al. “Diffusion Policy: Visuomotor Policy Learning via Action Diffusion”. In: *Proceedings of Robotics: Science and Systems (RSS)*. 2023.
- [7] Erwin Coumans and Yunfei Bai. “Pybullet, a python module for physics simulation for games, robotics and machine learning”. In: (2016).
- [8] Jacob Devlin et al. “BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding”. In: *ArXiv abs/1810.04805* (2019).
- [9] Danny Driess et al. “PaLM-E: An Embodied Multimodal Language Model”. In: *arXiv preprint arXiv:2303.03378*. 2023.
- [10] Yan Duan et al. “One-shot imitation learning”. In: *Advances in neural information processing systems* 30 (2017).
- [11] Pete Florence et al. “Implicit behavioral cloning”. In: *Conference on Robot Learning*. PMLR, 2022.
- [12] Caelan Reed Garrett et al. “Integrated task and motion planning”. In: *Annual review of control, robotics, and autonomous systems* 4 (2021).
- [13] Xiuye Gu et al. “Open-vocabulary Object Detection via Vision and Language Knowledge Distillation”. In: *International Conference on Learning Representations*. 2021.
- [14] Wenlong Huang et al. “Inner Monologue: Embodied Reasoning through Planning with Language Models”. In: *Conference on Robot Learning*. PMLR, 2023.
- [15] Wenlong Huang et al. “Language models as zero-shot planners: Extracting actionable knowledge for embodied agents”. In: *International Conference on Machine Learning*. PMLR, 2022.
- [16] Eric Jang et al. “BC-Z: Zero-Shot Task Generalization with Robotic Imitation Learning”. In: *5th Annual Conference on Robot Learning*. 2021.
- [17] Yunfan Jiang et al. “VIMA: General Robot Manipulation with Multimodal Prompts”. In: *arXiv preprint arXiv: Arxiv-2210.03094* (2022).
- [18] Alexander Kirillov et al. “Segment anything”. In: *arXiv preprint arXiv:2304.02643* (2023).
- [19] Takeshi Kojima et al. “Large language models are zero-shot reasoners”. In: *Advances in neural information processing systems* 35 (2022).
- [20] John Leonard et al. “Team MIT urban challenge technical report”. In: (2007).
- [21] Jacky Liang et al. “Code as policies: Language model programs for embodied control”. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [22] Kevin Lin et al. “Text2Motion: From Natural Language Instructions to Feasible Plans”. In: *arXiv preprint arXiv:2303.12153* (2023).
- [23] Corey Lynch et al. *Interactive Language: Talking to Robots in Real Time*. 2022. arXiv: 2210.06407 [cs.RO].
- [24] Michael Montemerlo et al. “Junior: The stanford entry in the urban challenge”. In: *Journal of field Robotics* 25.9 (2008).
- [25] R OpenAI. “GPT-4 technical report”. In: *arXiv* (2023).
- [26] Deepak Pathak\* et al. “Zero Shot Visual Imitation”. In: *International Conference on Learned Representations* (2018 (\*equal contribution)).
- [27] Alec Radford et al. *Learning Transferable Visual Models From Natural Language Supervision*. 2021. arXiv: 2103.00020 [cs.CV].
- [28] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. “CLIPort: What and Where Pathways for Robotic Manipulation”. In: *Proceedings of the 5th Conference on Robot Learning (CoRL)*. 2021.
- [29] Mohit Shridhar, Lucas Manuelli, and Dieter Fox. “Perceiver-Actor: A Multi-Task Transformer for Robotic Manipulation”. In: *Proceedings of The 6th Conference on Robot Learning*. Ed. by Karen Liu, Dana Kulic, and Jeff Ichnowski. Vol. 205. Proceedings of Machine Learning Research. PMLR, 2023, pp. 785–799.
- [30] Anthony Simeonov et al. “Neural descriptor fields: Se (3)-equivariant object representations for manipulation”. In: *2022 International Conference on Robotics and Automation (ICRA)*. IEEE, 2022.
- [31] Anthony Simeonov et al. “Se (3)-equivariant relational rearrangement with neural descriptor fields”. In: *Conference on Robot Learning*. PMLR, 2023.
- [32] Ishika Singh et al. “Progprompt: Generating situated robot task plans using large language models”. In: *2023 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2023.
- [33] Martin Sundermeyer et al. “Contact-graspnet: Efficient 6-dof grasp generation in cluttered scenes”. In: *2021 IEEE International Conference on Robotics and Automation (ICRA)*. IEEE, 2021.
- [34] Sai Vemprala et al. *ChatGPT for Robotics: Design Principles and Model Abilities*. Tech. rep. MSR-TR-2023-8. Microsoft, 2023.
- [35] Guanzhi Wang et al. “Voyager: An Open-Ended Embodied Agent with Large Language Models”. In: *arXiv preprint arXiv: Arxiv-2305.16291* (2023).
- [36] Jason Wei et al. “Chain-of-thought prompting elicits reasoning in large language models”. In: *Advances in Neural Information Processing Systems* 35 (2022).
- [37] Andy Zeng et al. *Socratic Models: Composing Zero-Shot Multimodal Reasoning with Language*. 2022. arXiv: 2204.00598 [cs.CV].
- [38] Andy Zeng et al. “Transporter Networks: Rearranging the Visual World for Robotic Manipulation”. In: *Conference on Robot Learning (CoRL)* (2020).
- [39] Xingyi Zhou et al. *Detecting Twenty-thousand Classes using Image-level Supervision*. 2022. arXiv: 2201.02605 [cs.CV].

## APPENDIX

### A1. ADDITIONAL EXPERIMENTS

#### A. Benefiting from spatially-grounded scene description

The next set of tasks we consider evaluates the benefit of providing a spatially-grounded scene description. We also consider basic tasks of picking objects and placing them in a described position, and the stacking task that involves a sequence of three pick-place actions. The main results are shown in Table II, where each row in the tasks column represent a set of tasks involving different object categories (for example, mugs, bowls, bottles) and receptacles (for example, containers, and baskets). Each task in a row is performed for  $\geq 10$  runs with varying object instances, and success or failure is assigned to each run, which is then used to find the success rate.

The first task category in scene description requires grasping a target object that lies beneath a second object. Without a relational scene description, the planner directly picks the target object. On the other hand, if the planner is informed that the second object is *above* the target, it calls `pick` and `place` on the second object before fetching the target object.

The second example showing the advantage of our scene description involves detecting task progress to minimize the number of actions that are used. For instance, to “put all objects of a specified category into a basket”, the generated plan completed the task in fewer steps if the planner knows some objects are already in the tray.

The final category of tasks calls for satisfying relational constraints with an *indirectly specified* reference object. One example is “place the apple in a tray *without mugs*”. When provided with the scene description, the planner detects the tray without any mugs and uses it as the placing target. In contrast, without the scene description, the system fails to place in the correct tray.

**Task: Fetch the mug**

**Without Scene Description**

A table has the following objects: a mug and a box.

**Output:**

```
start_task()
mug_id = find(object_label="mug")
pick(mug_id)
end task()
```

**With Scene Description**

A table has the following objects: a mug and a box.  
At the right of all the objects on the table lies the mug. Over the mug lies the box.

**Output:**

```
start_task()
object_ids = get_all_object_ids()
mug_id = find(object_label='mug', object_ids=object_ids)
objects_over_mug = \
    get_objects_contained_and_over(mug_id)
if objects_over_mug:
    for object_id in objects_over_mug:
        container_id = get_container_id(object_id)
        if container_id == mug_id:
            box_id = object_id
            pick(box_id)
            safe_spot_position = \
                get_place_position(box_id, mug_id,
                                   "to the left")
            place(box_id, safe_spot_position)
            break
pick(mug_id)
end task()
```

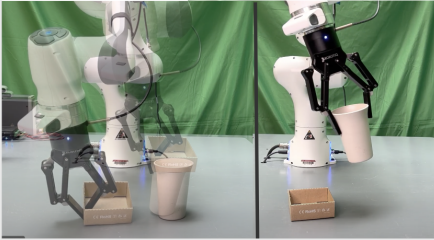


Fig. A1: Example plan for “box over a mug” task, with and without spatially-grounded scene description information in the input. If the scene description lacks spatial information, the planner fails to communicate that the box must be removed before picking up the mug.

	Tasks	Success
<b>Basic Tasks</b>	Pick object	91%
	Place object-1 to the left of object-2	80%
	Place object into receptacle	76%
	Stack three bowls	70%
<b>Scene Description</b>	Fetch object-1 (when object-2 lies over it)	79%
	Place all objects into the receptacle (partial progress)	75%
	Place object-1 into receptacle that has no object-2	75%

TABLE II: Success rates for evaluations in simulation.

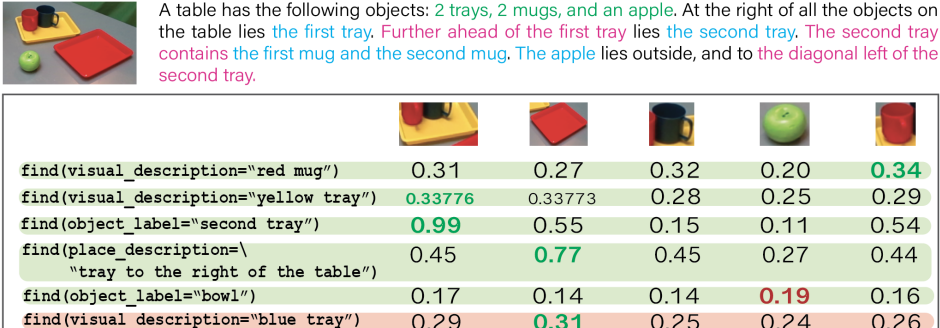
## A2. TUNABLE SYSTEM PARAMETERS

Parameter	Value	Description
Detection thd.	0.3	Object detection threshold in Detic
Mask Erosion	10	The number of pixels removed from the mask’s boundary to lessen the impact of high depth noise near object’s edges.
PCD Merging thd.	0.03	To find correspondence between objects in images from different views. If the change in std deviation for two object point clouds in the two images is less than the threshold, they are merged to represent one object.
Overlap thd.	20	In the top down projection of two objects convex hull intersection contains atleast the threshold points, then the objects belong to either the “above” case or the “contained” case.
Contained thd.	0.1	If the percentage of an object’s points that lie inside the convex hull of another object is less than this threshold, then the object is said to be “contained-in” in the other object.
Grasp thd.	0.8	Grasps (generated from [33]) with predicted success value higher than the threshold are considered.
Place description thd.	0.6	The threshold for <code>location-description</code> option in <code>find primitive</code> . Uses BERT [8] for the finding the score.
Visual description thd.	0.3	The threshold for <code>visual-description</code> option in <code>find primitive</code> . Uses CLIP [27] for the finding the score.
Object label thd.	0.3	The threshold for <code>object-label</code> option in <code>find primitive</code> . Uses BERT [8] for the finding the score.



### A3. OPEN VOCABULARY OBJECT DETECTION

The primitive function `find` performs the task of detecting any object in the scene, either through (a) an object label that comes from the scene description, (b) a visual description that describes an object’s visual properties (identifiable through CLIP features), and (c) based on the object’s location in the environment with respect to another object, as described through the scene description. Examples for different ways in which a call to `find` can be made is shown in Figure A2



A table has the following objects: 2 trays, 2 mugs, and an apple. At the right of all the objects on the table lies the first tray. Further ahead of the first tray lies the second tray. The second tray contains the first mug and the second mug. The apple lies outside, and to the diagonal left of the second tray.






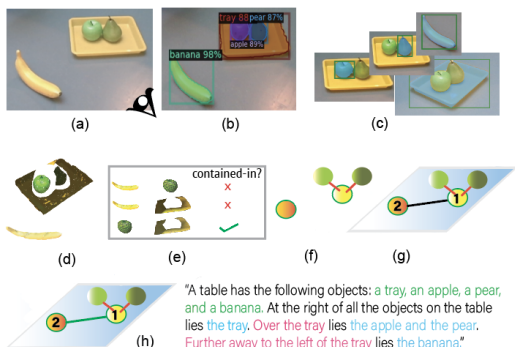
					
<code>find(visual_description="red mug")</code>	0.31	0.27	0.32	0.20	<b>0.34</b>
<code>find(visual_description="yellow tray")</code>	<b>0.33776</b>	0.33773	0.28	0.25	0.29
<code>find(object_label="second tray")</code>	<b>0.99</b>	0.55	0.15	0.11	0.54
<code>find(place_description="tray to the right of the table")</code>	0.45	<b>0.77</b>	0.45	0.27	0.44
<code>find(object_label="bowl")</code>	0.17	0.14	0.14	<b>0.19</b>	0.16
<code>find(visual_description="blue tray")</code>	0.29	<b>0.31</b>	0.25	0.24	0.26

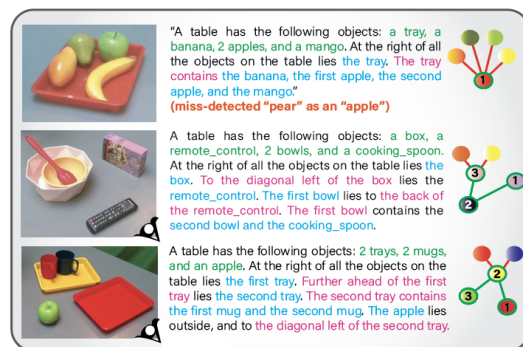
Fig. A2: Examples for `find` primitive in the API. Numbers in green pass the threshold and the object is found. Number in red are below threshold and no object is found. The green highlighted are true positives, while the red highlighted is an example of a false positive.

### A4. GENERATING SCENE DESCRIPTION

Generating scene description involves the following main steps: finding object relations between pair of objects, generating an object scene graph and defining relations over the edges and finally filling in the information in a template. These steps from (a) to (h) are shown in Figure A3a. First, multiple RGBD images are obtained the camera sensors followed by detection (b) and segmentation (c) in those images. For each image, we obtain an object instance’s partial point cloud which are merged with the partial point clouds from other camera images (d). Now each object’s point cloud is compared with another object’s; if they overlap and the convex hull of one covers some part of the object, then one is identified as containing the other object. Similar heuristic is used for identifying when an object lies above other object (e). Using these relations, objects are grouped together: a tray containing an apple will be treated as one vertex (f). Now all the grouped vertices are connected with one of their nearest vertex, and assigned one of the directional relations (left, right, front, down)(f). Next we traverse through the object graph and convert each edge, or relation into a sentence. While traversing the graph, on each grouped object node, the description of what objects are contained, or lies above within the node is also added to the description. Finally, we list all objects (and sum all objects of belonging to same category) and add the sentence to the beginning of the description, thus producing the final description (h). A few examples are shown in Figure A3b.



(a) Generating scene description involves steps (a) to (h): finding object point clouds, then finding spatial relations between pair of objects, followed by generating an object scene graph and defining relations over the edges and finally filling in the information in a template.



(b) Scene Description Examples: on the left are the example scene images, in which the user is standing on the right side of the image; at the center are the scene descriptions; and at the right side shows the object graphs that were traversed to produce the scene description.

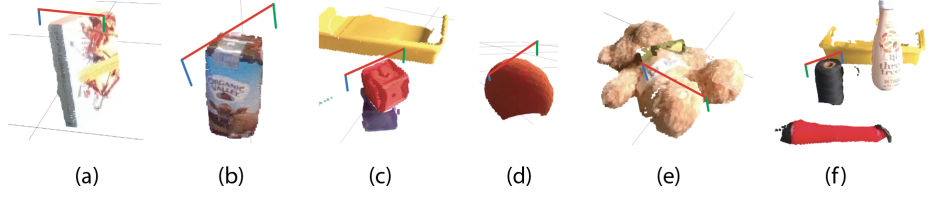


## A5. PICK AND PLACE PRIMITIVES

Pick and Place primitives are main skills that the library has by default, prior to new skills being added. The pick function takes in an object id and access the object's point cloud, runs Contact Graspnet [33] which is then executed. Some example objects we use for real world experiments are shown in Figure A4a and some of the computed grasps are on those objects are shown in Figure A4b.



(a) Example real objects



(b) Real World Grasp Examples. From left to right, showcasing the grasp for a book, a milk carton, a cube, a ball, a toy and a spool of thread.

The place function takes in the object id to be placed, and the place position. The place position usually comes from another function called `get-place-position` which takes in the object id of the object to be placed, the object id of a reference object and a description of how the object has to be placed with respect to the reference. This description provided as an argument to `get-place-position` function is matched using BERT embeddings, with a predefined set of descriptions of place positions with respect to the reference object. These predefined descriptions are same as the relations defined between objects in scene description (Section A4). Figure A5 shows some of the example place positions and their descriptions with respect to the object in blue.

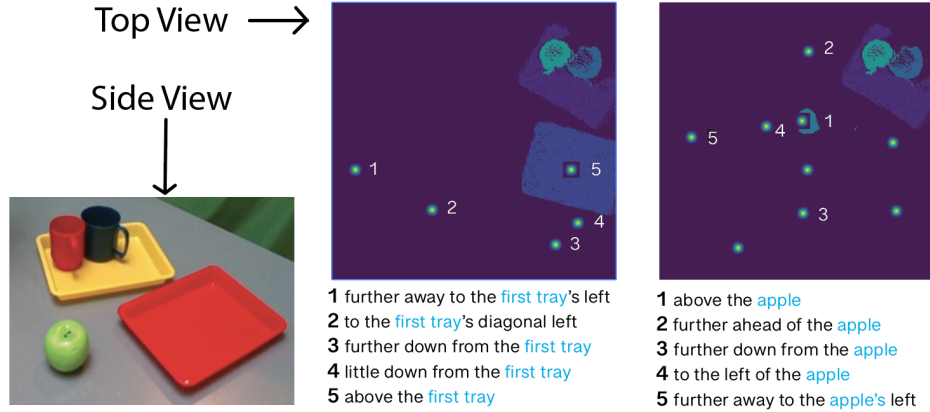


Fig. A5: Placement possibilities for an example scene. The description in `get-place-position` function argument is matched with these possibilities of the location descriptions.

## A6. PYTHON API

The Python code API provided in the input prompt to the language planner is seen here:

```
start_task()
    must be called at the start of any task. It starts the robot.

end_task()
    must be called when a task completes. It stops the robot.

get_all_object_ids()
    returns a list of all integer object ids present in the scene;
    Returns:
        ids: list(int)

get_container_id(object_id)
    gives the id of the object that contains 'object_id'
    Arguments:
        object_id: int
            id of the object that is contained in some container
    Returns:
        container_id: int or None
            the id of the container that contains 'object_id'
            None is returned when the object_id is not contained in
            any container.

get_objects_contained_and_over(object_id)
    gives the ids of all the objects that lie inside or over 'object_id'
    Arguments:
        object_id: int
            id of an object that contains something or over which lie
            other objects
    Returns:
        ids: list(int)
            the ids of all the objects that lie either inside or over
            the 'object_id' an empty list is returned when nothing lies
            over or inside.

find(object_label=None, visual_description=None, place_description=None,
      object_ids=None)
    Finds an object in the scene given atleast one of object_label,
    visual description or place description.
    Arguments:
        object_label: str
            The name with which the object has been referred to earlier
            For example, "the second tray", "the third bowl" etc
            By default, this argument is None

        visual_description: str
            object with some visual description of what it is.
            For example, "the red mug", "the blue tray", "the checkered box"
            By default, this argument is None

        place_description: str
            a string that describes where the object is located.
            For example, to find a bowl that is on the right of the
            tray, the function call will be
            'find("bowl that lies to the right of the tray")',
            or to get the mug that is contained in the second bowl,
            the call would be
            'find("mug that is inside the second bowl") and so on.
            By default, this argument is None

        object_ids: list(int)
            A list of 'int' object ids in which the object should
            be found, when specified.
            By default when this argument is None, all the objects
            are considered for
            finding the best matching

    Atleast one of the first three arguments must be specified.
    Typically, the use of the first and
    the second argument is enough but third can be used whenever needed.

    Returns: int
        object_id, an integer representing the object that best
```

matched with the description

`get_location(object_id)`

gives the location of the object 'object\_id'

Arguments:

object\_id: int  
Id of the object

Returns:

position: 3D array  
the location of the object

`pick(object_id)`

Picks up an object that 'object\_id' represents. A 'place' needs to occur before another call to pick, i.e. two picks cannot occur one after the other

Arguments:

object\_id: int  
Id of the object to pick

Returns: None

`get_place_position(object_id, reference_object_id, place_description)`

Finds the position to place the object 'object\_id' at a location described by 'place\_description' with respect to the 'reference\_object\_id'.

Arguments:

object\_id: int  
Id of the object to place  
reference\_object\_id: int  
id of the object relevant for placing the object\_id  
place\_description: str  
a string that describes where with respect to the reference\_object\_id the object\_id should be placed.

Returns: 3D array

the [x, y, z] value for the place location is returned.

For example,

to place a mug to the left of a bowl, the following function call should be used

`get_place_position(mug_id, bowl_id, "to the left")`

to place a mug into a bowl:

`get_place_position(mug_id, bowl_id, "inside")`

to place a mug above a box:

`get_place_position(mug_id, box_id, "above")`

`place(object_id, position)`

Moves to the position and places the object 'object\_id', at the location given by 'position' with the same orientation the object is currently in. The robot will open the gripper and drop the item at the position.

Arguments:

object\_id: int  
Id of the object to place  
position: 3D array  
the place location

Returns: None

`learn_skill(skill_name)`

Adds a new category-level skill to the current list of skills.

Arguments:

skill\_name: str  
a short name for the skill to learn, must be a string that can represent a function name (only alphabets and underscore can be used).  
highly-recommended that the object labels are included in the skill name.

Returns:

skill\_function: method  
a function that takes as input an object\_id and performs the skill on the object represented by the object\_id  
another relevant object\_id can be passed optionally. In particular, the returned function takes in arguments: object\_id\_1 and object\_id\_2:  
object\_id\_1: int  
Id of the object to act upon  
object\_id\_2: int (optional)  
Id of the object to place/interact relative to if relevant  
skill\_docstring: str  
a string that describes how to use the new skill in words, including relevant

inputs and outputs,  
along with any information on appropriate situations to use the skill,  
and misleading/confusing scenarios where it might make sense to use  
the skill but where a different skill should actually be used.  
This docstring should be printed out in the console so that the user can  
copy it and paste it into the skill API for use  
on subsequent runs (since the docstring will provide helpful context for  
how to appropriately use the skill in the future).

For example:

```
# Example 1:
drawer_id = 2
open_drawer, open_drawer_doc = learn_skill("open_drawer")
print(open_drawer_doc)
[Out:]
Grasps the handle of the drawer and executes a linear motion
in the direction away from the drawer, so that the drawer opens.
Arguments:
  object_id: int
              ID of the drawer to open
Returns: None
# opens the drawer represented by drawer_id
open_drawer(drawer_id)

# Example 2:
stick_id = 3
ball_id = 4
hit_ball_with_stick, hit_ball_with_stick_doc = learn_skill("hit_ball_with_stick")
print(hit_ball_with_stick_doc)
[Out:]
Hits a ball with a stick that is currently held in the gripper.
Hitting is performed by moving the stick to an offset position away
from the ball (depending on the current position of the ball obtained
by the ball point cloud and the size of the stick, obtained
by the stick point cloud) and then executing a fast gripper +
wrist motion to move the stick and hit the ball. Can only be executed
after executing some kind of "pick" for the stick, so that the stick
is in the grasp.
Args:
  object_id_1: int
                ID of the stick to use for hitting
  object_id_2: int
                ID of the ball that should be hit with the stick
Returns: None
# hits the ball represented by ball_id with the stick represented by stick_id
hit_ball_with_stick(stick_id, ball_id)
```

## A7. EXPERIMENTS

### A. Learning a new skill

Additional planner outputs for the “Empty the sink out” tasks in the previously described scene are provided in Figure A6 to demonstrate the alternative nonsensical plans when the LM has to conform to the static skill API.

<p><b>Alternative Plan 1</b></p> <p>LM: Let's proceed to the next step: tilting the mug to drop the banana onto the table. However, we don't have a tilt function available, so we'll have to work around it. We can move the mug to a higher position and then drop it, assuming that the banana will fall out. After this, the mug can be picked up again.</p> <pre>banana_id = find(place_description='banana that is inside the mug') # Place the mug at a higher position to drop the banana higher_position = get_location(mug_id) higher_position[2] += 0.5 # Increase the z-coordinate place(mug_id, higher_position)  # banana that is inside the mug was found. placing mug was successful.</pre> <p><b>Task was terminated early: the banana is still in the mug</b></p>	<p><b>Alternative Plan 2</b></p> <p>LM: The next task is to tilt the mug to drop the banana from it, onto the table. However, it appears from your functions that the robot cannot specifically perform a tilt action. So, for this step, we will first have to identify the banana inside the mug and pick it up. After picking the banana, the robot will place it on the table next to the sink.</p> <pre>banana_id = find(place_description='banana that is inside the mug') pick(banana_id) # Pick up the banana banana_place_position = get_place_position(banana_id, sink_id, 'to the right') place(banana_id, banana_place_position)  # banana that is inside the mug was found. grasping banana failed. It seems the pick failed, so aborting the place</pre> <p><b>Task was terminated early: unable to pick up the banana</b></p>	<p><b>Alternative Plan 3</b></p> <p>LM: The next task is to tilt the mug to drop the banana from it, onto the table. However, it appears from your functions that the robot cannot specifically perform a tilt action. Given the limitations of the API, we should approach this differently. We need to first place the mug down, then pick it up in such a way that the banana falls out. However, without a function to control how the mug is picked up or tilted, this might be challenging. Let's try to place the mug on the table and then pick it up again to see if that will allow the banana to fall out.</p> <pre>place(mug_id, get_location(mug_id)) pick(mug_id)  # placing mug was successful. the mug was picked successfully</pre> <p><b>Task was terminated early: the banana is still in the mug</b></p>
---	--	--

Fig. A6: Alternate suggestions from the LM to remove a banana from a mug, when learning a new skill is not an option.



<p><b>(A) Task: Pick up the mug by the handle</b></p> <pre># A table has the following objects: a mug.  start_task() mug_id = find(object_label="mug") # mug found. pick(mug_id) # the mug was picked successfully.</pre> <p><b>Output without Learning</b></p> <p><b>Task was terminated early: the mug was picked up by the rim instead of the handle.</b></p>	<p><b>Output with Learning</b></p> <pre>the mug was picked up by the rim. pick_mug_by_handle =   learn_skill("pick_mug_by_handle") pick_mug_by_handle(mug_id)  # pick_mug_by_handle has been learned and is available to use. picks can only be executed if the robot is not already holding something mug_position = get_location(mug_id) place(mug_id, mug_position)  # placing mug was successful. pick_mug_by_handle(mug_id) # pick_mug_by_handle has been executed. end_task()</pre>
<p><b>(B) Task: Place the bottle into the container on its side</b></p> <pre># A table has the following objects: a container and a bottle. At the diagonal left of the container lies the bottle.  start_task() container_id = find(object_label="container") bottle_id = find(object_label="bottle") # container found. bottle found. pick(bottle_id) # the bottle was picked successfully.</pre> <p><b>Output without Learning</b></p> <pre>place_position = get_place_position(bottle_id, container_id, "inside") place(bottle_id, place_position)  # placing bottle was successful. <b>Task was terminated early: the bottle was placed upright.</b></pre> <p><b>Output with Learning</b></p> <pre>place_bottle_on_side = learn_skill("place_bottle_on_side") place_bottle_on_side(bottle_id, container_id)  # place_bottle_on_side has been learned and is available to use. place_bottle_on_side has been executed end_task()</pre>	

Fig. A7: Planner outputs for additional tasks using baseline compared to our system.

With learning capabilities, it generates programs that accomplish the more difficult tasks by acquiring new policies to execute a handle grasp and a sideways reorientation placement seen in Figure A7.

### B. Scene Description Ablations

Plans for experiments in the scene description category are also shown in Figure A8a and Figure A8b.

 <p><b>Without Scene Description</b></p> <pre>A table has the following objects: a tray, an apple and a banana.  start_task() apple_id = find(object_label="apple") pick(apple_id) first_tray_id = \   find(object_label="the first tray") place_position = \   get_place_position(apple_id, first_tray_id, "inside") place(apple_id, place_position) end_task()</pre> <p><b>Output</b></p>	 <p><b>With Scene Description</b></p> <p><b>Task: Place the apple in the tray without mugs</b></p> <pre>A table has the following objects: 2 trays, 2 mugs, and an apple. At the right of all the objects on the table lies the first tray. Further ahead of the first tray lies the second tray. The second tray contains the first mug and the second mug. The apple lies outside and to the diagonal left of the second tray.  start_task() apple_id = find(object_label="apple") first_tray_id = find(object_label="first tray") second_tray_id = find(object_label="second tray") pick(apple_id) place_position = get_place_position(apple_id, first_tray_id, "inside") place(apple_id, place_position) end_task()</pre> <p><b>Output</b></p>
<p><b>Without Scene Description</b></p> <pre>A table has the following objects: a tray, an apple and a banana.  start_task() tray_id=find(object_label='tray') apple_id=find(object_label='apple') banana_id=find(object_label='banana')  pick(apple_id) place_position = \   get_place_position(apple_id, tray_id, "inside") place(apple_id, place_position) pick(banana_id) place_position = \   get_place_position(banana_id, \ tray_id, "inside") place(banana_id, place_position) end_task()</pre> <p><b>Output</b></p>	<p><b>With Scene Description</b></p> <p><b>Task: Put the fruits in the tray</b></p> <pre>A table has the following objects: a tray, an apple and a banana. At the right of all the objects on the table lies the tray. The tray contains the banana. Further away to the left of the tray lies the apple.  start_task() all_object_ids = get_all_object_ids() first_tray_id = find(object_label="the first tray", object_ids=all_object_ids) objects_in_first_tray = get_objects_contained_and_over(first_tray_id) for obj_id in all_object_ids:   if get_container_id(obj_id) == None:     pick(obj_id)     place_position = get_place_position(second_can_id, second_tray_id, "inside")     place(second_can_id, place_position) end_task()</pre> <p><b>Output</b></p>

(a) Planner output for tasks in scene description category with indirect reference to an object.

(b) Planner output for scene description category involving detecting task progress.

## A8. LIMITATIONS

**Scene Descriptions** The complexity of the textual description is an ongoing research problem, and it is difficult to capture every visual feature needed to plan optimally. Camera sensor noise, clutter, and failure to consolidate multiple views of the same object instance are all factors that may cause perception to fail. The perception module may misidentify some objects, while noise or clutter can cause spatial relationship heuristics to fail, thus generating inaccurate scene description that results in incorrect plans.

**Function Usage Errors** Functions that take in open-vocabulary text as parameters are often used incorrect or insufficiently by the LM. The system may then misidentify placements or objects due to incorrect correspondence of the language embeddings. In addition, occasionally the planner will request a new skill with some assumption about its parameters that are unknown until usage. Therefore the user may misinterpret the request and teach it the wrong skill. Finally the planner may request an entire sequence of actions as a new skill which goes against the purpose of the function. An improved API with better function descriptions and more examples may resolve the issue.

**Skill Primitives Failure** While executing pick and place primitives, the generated robotic arm motion plans may be infeasible due to joint limits, or failure in inverse kinematics. This lowers the success rate of the primitive skill functions. Another failure case in placements is when an object has to be placed inside another object, but due to noise in execution the object topples out of the container, thus failing the task.

**Feedback Automation** Our generated textual scene description and execution feedback aims to mitigate the amount of user feedback needed to form a closed-loop system. We fallback on user feedback to describe object states of actions that go beyond successful control execution, which is difficult to automate (see III-C.0.c). To build a self-sufficient embodied agent, a more sophisticated verification pipeline with automatic environment and object state feedback is required. Some systems use multimodal inputs similar to [9] as feedback.

**GPT-3.5 vs GPT-4** The code generation capabilities of GPT-4 far exceeds that of GPT-3.5, but we have a limited number of messages we are able to send, so it is difficult to collect a large amount of experiment data for that reason.

## A9. TASKS FOR LLM-ONLY `LEARN_SKILL` AND SKILL RE-USE EVALUATION

In the subsections below, we include a brief description of each task used to evaluate the LLM on its own in its ability to request new skills, avoid requesting unnecessary skills, and perform tasks with an expanded skill library after a newly-learned skill is added. Each of these tasks was manually curated by hand and did not reflect any real-world environment, nor were the plans meant to be executable on the robot. For brevity, we include the full prompt (including the manually designed scene description) of one of the tasks for each section, and only include a concise description of the task to be solved for the rest.

*a) Ask for a new skill when needed –:* We considered the tasks listed below for evaluating how frequently GPT-4 correctly requests a new skill to be learned using the `learn_skill` function:

- mix the ingredients (a bowl and a spoon, bowl is filled)

A table has the following objects: a bowl and a spoon.  
At the right of all the objects on the table lies the bowl.  
The bowl is filled. To the left of the bowl lies the spoon.  
If you are commanding a robot, tell me in words the steps  
to mix the ingredients?

- empty the mug (mug contains water)
- hang the mug on the mug rack
- place the mug upside down
- bottle in the container horizontally
- put the bowl on its side on the drying rack
- scrub the bowl
- open the drawer
- slide the tray into the oven
- push the button
- turn on the light

*b) Don't ask for a new skill when not needed –:* Contrasting the above scenario, the tasks listed below are used to evaluate how frequently the planner *ignores* the `learn_skill` function and directly attempts to complete the task with the base set of primitive skills that are available. The tasks are constructed such that they can be achieved using the base set of skills:

- build a block tower in rainbow color order

A table has the following objects: a red block, blue block,



and green block. At the right of all the objects on the table lies the blue block. To the left of the blue block lies the red block. To the left of the second bowl lies the green block. If you are commanding a robot, tell me in words the steps to build a block tower in rainbow color order?

- stack the bowls
- put the dishes away
- get me a cup I can pour water into
- empty the bowl
- put the dishes in the sink

c) *When new skills have been added, how often the agent succeeds at a new task which may or may not require the new skill* –: Besides evaluating the LLM’s ability to correctly use the base skill set and request new abilities with `learn_skill`, we also evaluated the planner’s success when prompted with new tasks after newly-learned skills have been added to the API. These tasks have been constructed to sometimes require the use of the newly-learned skill and other times specifically require *not* using the newly learned skill. The tasks we use for this evaluation are listed below:

**After learning and adding `pick_mug_by_handle` to API:**

- pick up the mug (handle of the mug is broken)

A table has the following objects: a mug. At the right of all the objects on the table lies the mug. The handle of the mug is broken. If you are commanding a robot, tell me in the words the steps to pick up the mug?

- empty the mug
- pick up the mug (box is blocking access to the handle)
- bring the cup of coffee to the living room
- pick up the mug (handle of the mug is covered with a dirty, sticky substance)
- fill the mug with water

**After learning and adding `side_pick_bottle` to API:**

- fetch the bottle and put it on the table (bottle is on a bottom shelf, cannot be approached from the top)

A table has the following objects: a bottle. At the right of all the objects on the table lies the bottle. Due to the position on the table and the low height of the table, the bottle cannot be reached from the side (but could be reached from above). If you are commanding a robot, tell me in the words the steps to fetch the bottle?

- fetch the bottle (bottle cannot be reached from the side, but could be reached from above)
- fetch the bottle and put it on the table (bottle inside a box, box is sideways, top of the box is open)
- fetch the bottle and put it on the table (bottle is on a bottom shelf)

**After learning and adding `place_book_horizontally` and `place_book_vertically` to API:**

- put away the book (bookshelf with single book, book doesn’t fit vertically)

A table has the following objects: a book and a bookshelf. At the right of all the objects on the table lies the bookshelf. Next to the bookshelf lies the book. The book doesn’t fit in the shelf vertically. If you are commanding a robot, tell me in the words the steps to put away the book?

- put away the book (bookshelf with single book, bookshelf contains many vertically aligned books, with space for one more)
- put away the book (bookshelf with single book, book LxWxH and bookshelf LxWxH dimensions provided)
- put away all of the books (bookshelf and 10 books total, only three can fit when horizontal)