

# Learning Compositional Behaviors from Demonstration and Language

Anonymous Author(s)

1       **Abstract:** We introduce Behavior from Language and Demonstration (BLADE), a  
2       framework for long-horizon robotic manipulation by integrating imitation learning  
3       and model-based planning. BLADE leverages language-annotated demonstrations,  
4       extracts abstract action knowledge from large language models (LLMs), and con-  
5       structs a library of structured, high-level action representations. These represen-  
6       tations include preconditions and effects grounded in visual perception for each  
7       high-level action, along with corresponding controllers implemented as neural  
8       network-based policies. BLADE can recover such structured representations auto-  
9       matically, without manually labeled states or symbolic definitions. BLADE shows  
10      significant capabilities in generalizing to novel situations, including novel initial  
11      states, external state perturbations, and novel goals. We validate the effectiveness  
12      of our approach both in simulation and on a real robot with a diverse set of objects  
13      with articulated parts, partial observability, and geometric constraints.

14      **Keywords:** Manipulation, Planning Abstractions, Learning from Language

## 16   1 Introduction

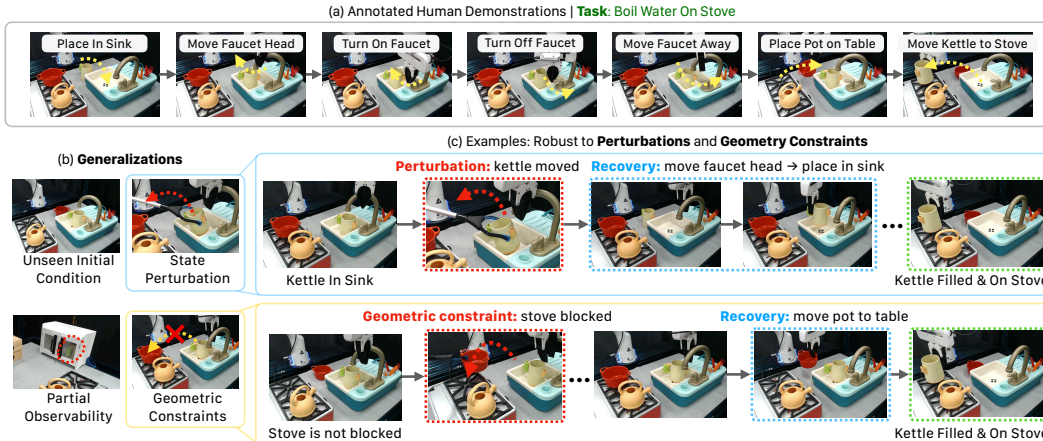
17   Developing autonomous robots capable of completing long-horizon manipulation tasks that involve  
18   interacting with many objects is a significant milestone. We want to build robots that can directly  
19   perceive the world, operate over extended periods, generalize to various states and goals, and are  
20   robust to perturbations. A promising direction is to combine learned policies with model-based  
21   planners, allowing them to operate on different time scales. In particular, imitation learning-based  
22   methods have proven highly successful in learning policies for various “behaviors,” which usually  
23   operate over a short time span [e.g., 1]. To solve more complex and longer-horizon tasks, we can  
24   compose these behaviors by planning in explicit abstract action spaces [2–4], in latent spaces [5], or  
25   via large pre-trained models such as large language models [6].

26   However, one of the key challenges of all high-level planning approaches is the automatic acquisition  
27   of an abstraction for the learned “behaviors” to support long-horizon planning. The goal of this  
28   behavior abstraction learning is to build representations that describe the preconditions and effects of  
29   behaviors, to enable chaining and search. These representations should depend on the environment, the  
30   set of possible goals, and the specifications of individual behaviors. Furthermore, these representations  
31   should be grounded on high-dimensional perception inputs and low-level robot control commands.

32   Our insight into tackling this challenge is to leverage knowledge from two sources: the low-level,  
33   mechanical understanding of robot-object contact, and the high-level, abstract understanding of  
34   object-object interactions described in language that can be extracted from language models as the  
35   knowledge source. We bridge them by learning the grounding of abstract language terms on visual  
36   perception and robot actuation. Our framework, behavior from language and demonstration (BLADE),  
37   takes as input a small number of language-annotated demonstrations (Fig. 1a). It segments each  
38   trajectory based on which object is in contact with the robot. Then, it uses a large language model  
39   (LLM), conditioned on the contact sequences and the language annotations, to propose abstract  
40   behavior descriptions with preconditions and effects that best explain the demonstration trajectories.

---

\* denotes equal contribution.



**Figure 1: BLADE**, a robot manipulation framework combining imitation learning and model-based planning. (a) BLADE takes language-annotated demonstrations as training data. (b) It generalizes to unseen initial conditions, state perturbations, and geometric constraints. (c) In the depicted scenarios, BLADE recovers from perturbations such as moving the kettle out of the sink, and resolves geometric constraints including a blocked stove.

41 During training, we extract the state abstraction terms from the preconditions and effects (e.g.,  
 42 *turned-on*, *aligned-with*), and learn their groundings on perception inputs. We also learn the control  
 43 policies associated with each behavior (e.g., *turn on the faucet*).

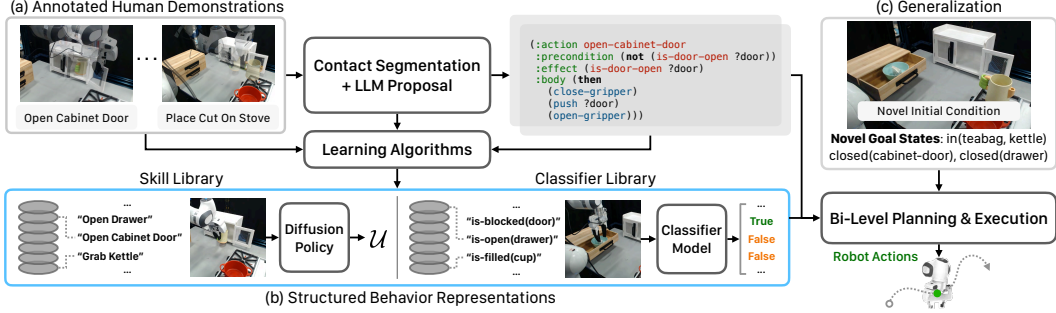
44 Our model offers several advantages. First, unlike prior work that relies on manually defined state  
 45 abstractions or additional state labels, our method automatically generates state abstraction labels  
 46 based on the language annotations and LLM-proposed behavior descriptions. BLADE recovers the  
 47 visual grounding of these abstractions without any additional label. Second, BLADE generalizes to  
 48 novel states and goals by composing learned behaviors using a planner. Shown in Fig. 1b, it can  
 49 handle various novel initial conditions and external perturbations that lead to unseen states. Third,  
 50 our method can handle novel geometric constraints (Fig. 1c), novel goals expressed in learned state  
 51 abstractions, and partial observability from articulated bodies like drawers.

## 52 2 Related Work

53 **Composing skills for long-horizon manipulation.** A large body of model-based planning methods  
 54 use manually-defined transition models [2, 7–9] or models learned from data [10–15] to generate  
 55 long-horizon plans. However, learning dynamics models with accurate long-term predictions and  
 56 strong generalization remains challenging. Another related direction is to introduce hierarchical struc-  
 57 tures into the policy models [16–20], where different methods have been introduced to decompose  
 58 continuous demonstrations into segments for short-horizon skills [20–22]. Unable to model the de-  
 59 pendencies between the skills, these methods are limited to following sequentially specified subgoals  
 60 and struggle to generalize to unseen goals. Researchers have also used learned models to improve  
 61 state estimation [23] and planning efficiency [24]. However, they still require manual definitions of  
 62 planning knowledge. Some work addresses this issue by learning the dependencies between actions  
 63 from data, but they still require large-scale supervised datasets [25–27]. In contrast, BLADE learns  
 64 planning-compatible action representations from only language-annotated demonstrations.

65 **Using LLMs for planning.** Many researchers have explored using LLMs for planning. Methods  
 66 for direct generation of action sequences [28, 29] usually do not produce accurate plans [30, 31].  
 67 Researchers have also leveraged LLMs as translators from natural language instructions to symbolic  
 68 goals [32–35], as generalized solvers [36], as memory modules [37], and as world models [38, 39].  
 69 To improve the planning accuracy of LLMs, prior work has explored techniques including learning  
 70 affordance functions [6, 40], replanning [41], finetuning [42–44], and VLM-based decision-making  
 71 [45, 46]. BLADE shares a similar spirit as methods using LLMs to generate planning-compatible action  
 72 representations [47–49]. However, they all make assumptions on the availability of state abstractions,  
 73 while BLADE automatically grounds LLM-generated action definitions without additional labels.





**Figure 2: Overview of BLADE.** (a) It receives language-annotated human demonstrations, (b) segments demonstrations into contact primitives, and learns a structured behavior representation. (c) It generalizes to novel initial conditions, leveraging bi-level planning and execution to achieve goal states.

### 74 3 Problem Formulation

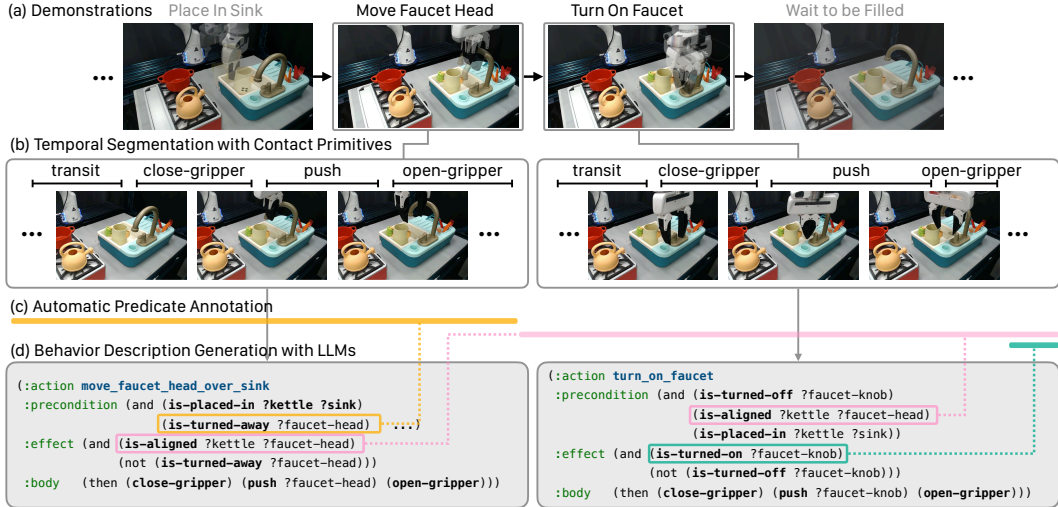
75 We consider the problem of learning a language-conditioned goal-reaching manipulation policy.  
 76 Formally, the environment is modeled as a tuple  $\langle \mathcal{X}, \mathcal{U}, \mathcal{T} \rangle$  where  $\mathcal{X}$  is the raw state space,  $\mathcal{U}$  is the  
 77 low-level action space, and  $\mathcal{T} : \mathcal{X} \times \mathcal{U} \rightarrow \mathcal{X}$  is the transition function (which may be stochastic and  
 78 unknown). Furthermore, the robot will receive observations  $o \in \mathcal{O}$  that may be partially observable  
 79 views of the states. At test time, the robot also receives a natural language instruction  $\ell_t$ , which  
 80 corresponds to a set of goal states. An oracle goal satisfaction function defines whether the language  
 81 goal is reached, i.e.,  $g_{\ell_t} : \mathcal{X} \rightarrow \{T, F\}$ . Given an initial state  $x_0 \in \mathcal{X}$  and the instruction  $\ell_t$ , the  
 82 robot should generate a sequence of low-level actions  $\{u_1, u_2, \dots, u_H\} \in \mathcal{U}^H$ .

83 In the language-annotated learning setting, the robot has a dataset of language-annotated demonstra-  
 84 tions  $\mathcal{D}$ . Each demonstration is a sequence of robot actions  $\{u_1, \dots, u_H\}$  paired with observations  
 85  $\{o_0, \dots, o_H\}$ . Each trajectory is segmented into  $M$  subtrajectories, and natural language descriptions  
 86  $\{\ell_1, \dots, \ell_M\}$  are associated with the segments (e.g., “place the kettle on the stove”). In this paper, we  
 87 assume that there is a finite number of possible  $\ell$ s—each corresponding to a skill to learn.

88 Directly learning a single goal-conditioned policy that can generalize to novel states and goals is  
 89 challenging. Therefore, we recover an *abstract* state and action representation of the environment and  
 90 combine online planning in abstract states and offline policy learning for low-level control to solve  
 91 the task. In BLADE, behaviors are represented as temporally extended actions with preconditions and  
 92 effects characterized by state predicates. Formally, we want to recover a set of predicates  $\mathcal{P}$  that define  
 93 an abstract state space  $\mathcal{S}$ . We focus on a scenario where all predicates are binary. However, they are  
 94 grounded on high-dimensional sensory inputs. Using  $\mathcal{P}$ , a state can be described as a set of grounded  
 95 atoms such as  $\{kettle(A), stove(B), filled(A), on(A, B)\}$  for a two-object scene. BLADE will learn a  
 96 function  $\Phi : \mathcal{O} \rightarrow \mathcal{S}$  that maps observations to abstract states. In its current implementation, BLADE  
 97 requires humans to additionally provide a list of predicate names in natural language, which we  
 98 have found to be helpful for LLMs to generate action definitions. We provide additional ablations  
 99 in the Appendix A.2. Based on  $\mathcal{S}$ , we learn a library of *behaviors* (a.k.a., *abstract actions*). Each  
 100 behavior  $a \in \mathcal{A}$  is a tuple of  $\langle name, args, pre, eff, \pi \rangle$ . *name* is the name of the action. *args* is a list of  
 101 variables related to the action, often denoted by  $?x, ?y$ . *pre* and *eff* are the precondition and effect  
 102 formula defined in terms of the variables *args* and the predicates  $\mathcal{P}$ . A low-level policy  $\pi : \mathcal{O} \rightarrow \mathcal{U}$   
 103 is also associated with  $a$ . The semantics of the preconditions and effects is: for any state  $x$  such that  
 104  $pre(\Phi(x))$  is satisfied, executing  $\pi$  at  $x$  will lead to a state  $x'$  such that  $eff(\Phi(x'))$  [50].

### 105 4 Behavior from Language and Demonstration

106 BLADE is a method for learning abstract state and action representations from language-annotated  
 107 demonstrations. It works in three steps, as illustrated in Fig. 2. First, we generate a symbolic behavior  
 108 definition conditioned on the language annotations and contact sequences in the demonstration using  
 109 a large language model (LLM). Next, we learn the classifiers associated with all state predicates and  
 110 the control policies, all from the demonstration without additional annotations. At test time, we use a  
 111 bi-level planning and execution strategy to generate robot actions.



**Figure 3: Behavior Descriptions Learning.** Starting with (a) human demonstrations with language annotations, BLADE segments (b) the demonstrations into contact primitives such as “close-gripper,” and “push.” Then, BLADE (d) generates operators using an LLM, defining actions with specific preconditions and effects. (c) These operators allow for automatic predicate annotation based on the preconditions and effects.

## 112 4.1 Behavior Description Learning

113 Given a finite set of behaviors with language descriptions  $\{\ell\}$  and corresponding demonstration  
 114 segments, we generate an abstract description for each  $\ell$  by querying large language models. To  
 115 facilitate LLM generation, we provide additional information on the list of objects with which the  
 116 robot has contact. The generated operators are further refined with abstract verification.

117 **Temporal segmentation.** We first segment each demonstration (Fig. 3a) into a sequence of *contact-*  
 118 *based primitives* (Fig. 3b). In this paper we consider seven primitives describing the interactions  
 119 between the robot and other objects: *open/close* grippers without holding objects, *move-to*( $x$ ) which  
 120 moves the gripper to an object, *grasp*( $x, y$ ) and *place*( $x, y$ ) which grasp and place object  $x$  from/onto  
 121 another object  $y$ , *move*( $x$ ) which moves the currently holding object  $x$  and *push*( $x$ ). We leverage  
 122 proprioception, i.e., gripper open state, and object segmentation to automatically segment the con-  
 123 tinuous trajectories into these basis segments. For example, pushing the faucet head away involves  
 124 the sequence of  $\{\textit{close-gripper}, \textit{push}, \textit{open-gripper}\}$ . This segmentation will be used for LLMs to  
 125 generate operator definitions and for constructing training data for control policies.

126 **Behavior description generation with LLMs.** Our behavior description language is based on  
 127 PDDL [51]. We extend the PDDL definition to include a *body* section which is a sequence of contact  
 128 primitives. It will be generated by the LLM based on the demonstration data.

129 Our input to the LLM contains four parts: 1) a general description of the environment, 2) the natural  
 130 language descriptions  $\ell$  associated with the behavior itself and other behaviors that have appeared  
 131 preceding  $\ell$  in the dataset, 3) all possible sequence of contact primitive sequences associated with  
 132  $\ell$  across the dataset, and 4) additional instructions on the PDDL syntax, including a single PDDL  
 133 definition example. We find that the inclusion of previous behaviors and contact primitive sequences  
 134 improves the overall generation quality. As shown in Fig. 3c, in addition to preconditions and effects  
 135 of the operators, we also ask LLMs to predict a *body* of contact primitive sequence associated with  
 136 the behavior, which we call *body*. We assume that each behavior has a single corresponding contact  
 137 primitive sequence, and use this step to account for noises in the segmentation annotations. After  
 138 LLM predicts the definition for all behavior, we will re-segment the demonstrations associated with  
 139 each behavior based on the LLM-predicted body section.

140 **Behavior description refinement with abstract verification.** Besides checking for syntax errors,  
 141 we also verify the generated behavior descriptions by performing *abstract verification* on the demon-  
 142 stration trajectories. In particular, given a segmented sequence of the trajectory where each segment  
 143 is associated with a behavior, we verify whether the preconditions of each behavior can be satisfied

144 by the accumulated effects of the previous behaviors. This verification does not require learning the  
145 grounding of state predicates and can be done at the behavior level to discover incorrect preconditions  
146 and effects, and at the contact primitive level to find missing or incorrect contact primitives (e.g.,  
147 *grasp* cannot be immediately followed by other *grasp*). We resample behavior definitions that do not  
148 pass the verification test.

## 149 4.2 Classifier and Policy Learning

150 Given the dataset of state-action segments associated with each behavior, we train the classifiers for  
151 different state predicates and the low-level controller for each behavior.

152 **Automatic predicate annotation.** We leverage *all* behavior descriptions to automatically label an  
153 observation  $\bar{o} = \{o_1, \dots, o_H\}$  based on its associated segmentation. In particular, at  $o_0$ , we label all  
154 state predicates as “unknown.” Next, we unroll the sequence of behavior executed in  $\bar{o}$ . As illustrated  
155 in Fig. 3c, before applying a behavior  $a$  at step  $o_t$ , we label all predicates in  $pre_a$  true. When  $a$   
156 finishes at step  $o_{t'}$ , we label all predicates in  $eff_a$ . In addition, we will propagate the labels for state  
157 predicates to later time steps until they are explicitly altered by another behavior  $a$ . In contrast to  
158 earlier methods, such as Migimatsu and Bohg [52] and Mao et al. [53], which directly use the first  
159 and last state of state-action segments to train predicate classifiers, our method greatly increases the  
160 diversity of training data. After this step, for each predicate  $p \in \mathcal{P}$ , we obtain a dataset of paired  
161 observations  $o$  and the predicate value of  $p$  at the corresponding time step.

162 **Classifier learning.** Based on the state predicate dataset generated from behavior definitions, we train  
163 a set of state classifiers  $f_\theta(p) : \mathcal{O} \rightarrow \{T, F\}$ , which are implemented as standard neural networks for  
164 classification. We include implementation details in Appendix A.6. In real-world environments with  
165 strong data-efficiency requirements, we additionally use an open vocabulary object detector [54] to  
166 detect relevant objects for the state predicate and crop the observation images. For example, only  
167 pixels associated with the object faucet will be the input to the *turned-on(faucet)* classifier.

168 **Policy learning.** For each behavior, we also train control policies  $\pi_\theta(a) : \mathcal{O} \rightarrow \mathcal{U}$ , implemented as  
169 a diffusion policy [1]. In simulation, we use a combination of frame-mounted and wrist-mounted  
170 RGB-D cameras as the inputs to the diffusion policy, while in the real world, the policy takes raw  
171 camera images as input. The high-level planner orchestrates which of these low-level policies to  
172 deploy based on the scene and states. Once trained on these diverse demonstrations of different skills,  
173 the resulting low-level policies can adapt to local changes, such as variations in object poses.

## 174 4.3 Bi-Level Planning and Execution

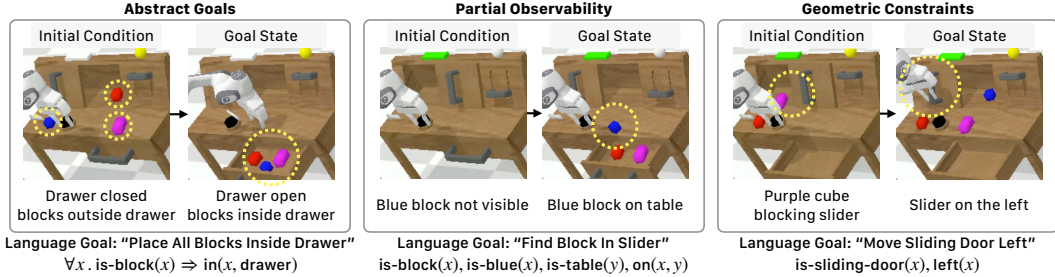
175 At test time, given a novel state and a novel goal, BLADE first uses LLMs to translate the goal into a  
176 first-order logic formula based on the state predicates. Next, it leverages the learned state abstractions  
177 to perform planning in a symbolic space to produce a sequence of behaviors. Then, we execute  
178 the low-level policy associated with the first behavior, and we re-run the planner after the low-level  
179 policy finishes—this enables us to handle various types of uncertainties and perturbations, including  
180 execution failure, partial observability, and human perturbation.

181 Visibility and geometric constraints are also modeled as preconditions, in addition to other object-  
182 state and relational conditions. For example, the behavior “opening the cabinet door” will have  
183 preconditions on the initial state of the door, a visibility constraint that the door is visible, and a  
184 geometric constraint that nothing is blocking the door. When those preconditions are not satisfied,  
185 the planner will automatically generate plans, such as actions that move obstacles away, to achieve  
186 them. Partial observability was handled by using the most-likely state assumption during planning  
187 and performing replanning. We include details in Appendix A.8.

# 188 5 Experiments

## 189 5.1 Simulation Experimental Setup

190 We use the CALVIN benchmark [55] for simulation-based evaluations, which include teleoperated  
191 human-play data. We use the split  $D$  of the dataset, which consists of approximately 6 hours of  
192 interactions. Annotations of the play data are generated by a script that detects goal conditions on



**Figure 4: Generalization Tasks in CALVIN.** Examples from the three generalization tasks in the CALVIN simulation environment. Successfully completing these tasks require planning for and executing 3-7 actions.

**Table 1: Generalization results in CALVIN.** Mean success rates with STD from three seeds are reported. BLADE outperforms latent planning, LLM, and VLM baselines in completing novel long-horizon tasks.

Method	State Classifier	Latent Feasibility	Generalization Task		
			Abstract Goal	Geometric Constraint	Partial Observability
HULC [56]	N/A	N/A	2.78 ± 3.47	11.67 ± 11.55	0.00 ± 0.00
SayCan [6]	N/A	Short	0.00 ± 0.00	0.00 ± 0.00	0.00 ± 0.00
VILA [45]	N/A	N/A	18.38 ± 2.48	0.00 ± 0.00	4.17 ± 5.20
T2M-Shooting [40]	Learned	Long	57.78 ± 12.29	0.00 ± 0.00	13.33 ± 1.44
Ours	Learned	N/A	<b>68.33 ± 10.14</b>	<b>26.67 ± 7.64</b>	<b>75.83 ± 3.82</b>
T2M-Shooting [40]	GT	Long	61.67 ± 5.00	0.00 ± 0.00	0.83 ± 1.44
Ours	GT	N/A	<b>76.11 ± 6.74</b>	<b>56.67 ± 16.07</b>	<b>70.00 ± 5.00</b>

193 simulator states, and there are in total 34 types of behaviors. We use RGB-D images from the mounted  
 194 camera for classifier learning and partial 3D point clouds recovered from the RGB-D cameras for  
 195 policy learning. The original benchmark focuses only on evaluating individual skills. To evaluate the  
 196 ability of different algorithms to compositionally combine previously learned policies to solve novel  
 197 tasks, we design six new generalization tasks, as shown in Fig. 4. Each task has a language instruction,  
 198 a sampler that generates random initial states, and a goal satisfaction function for evaluation. For  
 199 each task, we sample 20 initial states and evaluate all methods with three different random seeds. See  
 200 Appendix B.1 for more details on the benchmark setup.

201 **Baselines.** We compare BLADE with two groups of baselines: hierarchical policies with planning in  
 202 latent spaces and LLM/VLM-based methods for robotic planning. For the former, we use HULC [56],  
 203 the state-of-the-art method in CALVIN, which learns a hierarchical policy from language-annotated  
 204 play data using hindsight labeling. For the latter, we use SayCan [6], Robot-VILA [45], and  
 205 Text2Motion [40]. Note that Text2Motion assumes access to ground-truth symbolic states. Hence we  
 206 compare Text2Motion with BLADE in two settings: one with the ground-truth states and the other  
 207 with the state classifiers learned by BLADE. See Appendix B.2 for more details on these methods.

## 208 5.2 Results in Simulation

209 Table 1 presents the performance of different models in all three types of generalization tasks.

210 **Structured behavior representations improve long-horizon planning.** We first focus on the  
 211 comparison with the hierarchical policy model HULC in Table. 1. BLADE with learned classifiers  
 212 achieves a more than 65% improvement in the success rate for reaching abstract goals while using the  
 213 same language-annotated play data. We attribute this to the particular implementation of hindsight  
 214 labeling in HULC being not sufficient to achieve goals that require chaining together multiple high-  
 215 level actions: for example, the task of placing all blocks in the closed drawer requires chaining  
 216 together a minimum of 7 behaviors.

217 **Structured transition models learned by BLADE facilitate long-horizon planning.** Both SayCan  
 218 and T2M-Shooting learn a long-horizon transition and action feasibility model for planning. Shown  
 219 in Table. 1, learning accurate feasibility models directly from raw demonstration data remains a  
 220 significant challenge. In our experiment, we find that first, when the LLM does not take into account  
 221 state information (SayCan), using the short-horizon feasibility model is not sufficient to produce



222 sound plans. Second, since our model learns a structured transition model, factorized into different  
 223 state predicates, BLADE is capable of producing longer-horizon plans.

224 **Structured scene representations facilitate making feasible plans.** Compared to the Robot-VILA  
 225 method, which directly predicts action sequences based on the image state, BLADE first uses learned  
 226 state classifiers to construct an abstract state representation. This contributes to a 49% improvement  
 227 on the Abstract Goal tasks in Table 1. We observe that the pre-trained VLM used in Robot-VILA  
 228 often predicts actions that are not feasible in the current state. For example, Robot-VILA consistently  
 229 performs better in completing “placing all blocks in a closed drawer” than “placing all blocks in an  
 230 open drawer” since it always predicts opening the drawer as the first step.

231 **Explicit modeling of geometric constraints and object visibility improves performance in these**  
 232 **scenarios.** BLADE can reason about these challenging situations without explicitly being trained  
 233 in those settings. Table. 1 shows that our approach consistently outperforms baselines in these two  
 234 settings. These generalization capabilities are built on the explicit modeling of geometric constraints  
 235 and object visibility in behavior preconditions.

236 **BLADE can automatically propose operators for the specific environment given demonstrations.**  
 237 Our experiment shows that the LLM can automatically propose high-quality behavior descriptions  
 238 that resemble the dependency structures among operators. For example, the LLM discovers from  
 239 the given contact primitive sequences and language-paired demonstration that blocks can only be  
 240 placed after the block is lifted and that a drawer needs to be opened before placing objects inside, etc.  
 241 Some of these dependencies are unique to the CALVIN environment, therefore requiring the LLM to  
 242 generate specifically for this domain. We provide more visualizations in the Appendix A.1.

243 **BLADE’s automatic predicate annotation**  
 244 **enables better classifier learning.** From  
 245 Table 1, we observe that having accurate  
 246 state classifier models is critical for algo-  
 247 rithms’ performance (GT vs. Learned).  
 248 Hence, we perform additional ablation stud-  
 249 ies on classifier learning. Migimatsu and

**Table 2:** Ablation on state classifier learning in CALVIN.

Method	Abstract	Geometric	Partial Obs.
[52]	33.89 ± 5.85	9.17 ± 5.20	3.33 ± 2.89
BLADE	<b>68.33 ± 10.14</b>	<b>26.67 ± 7.64</b>	<b>75.83 ± 3.82</b>

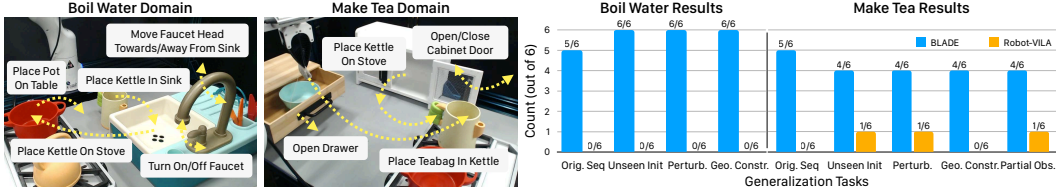
250 Bohg [52] also presented a method for learning the preconditions and effects of actions from seg-  
 251 mented trajectories and symbolic action descriptions. The key difference between BLADE and theirs is  
 252 that they only use the first and last frame of each segment to supervise the learning of state classifiers.  
 253 We compare the two classifier learning algorithms, given the same LLM-generated behavior defini-  
 254 tions, by evaluating the classifier accuracy on held-out states. BLADE shows a 20.7% improvement in  
 255 F1 (16.3% improvement for classifying object states and 38.6% improvement for classifying spatial  
 256 relations) compared to the baseline model. This also translates into significant improvements in the  
 257 planning success rate, as shown in Table 2,

### 258 5.3 Real World Experiments

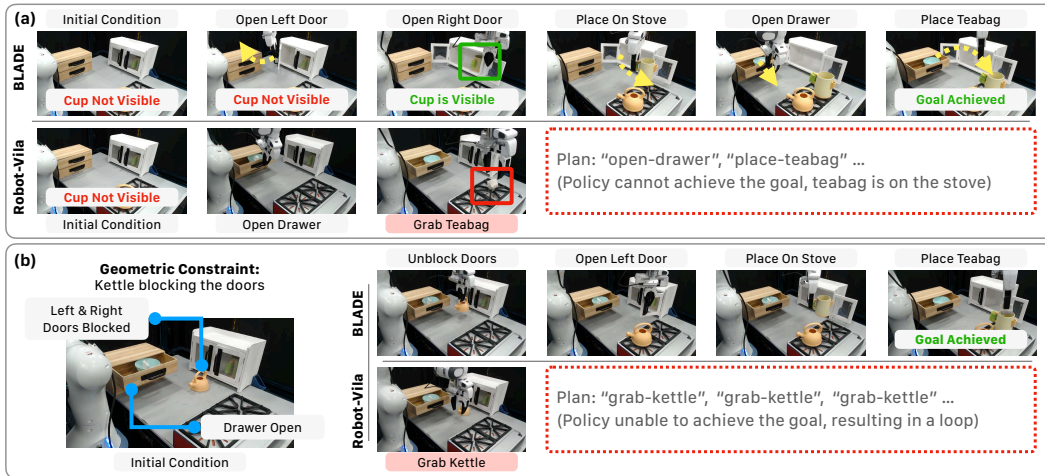
259 **Environments.** We use a Franka Emika robot arm with a parallel jaw gripper. The setup includes  
 260 five RealSense RGB-D cameras, with one being wrist-mounted on the robot and the remaining  
 261 positioned around the workspace. Fig. 5 shows the two domains: Make Tea and Boil Water. For  
 262 each domain, we collect 85 language-annotated demonstrations using teleoperation with a 3D mouse.  
 263 After segmenting the demonstrations using proprioception sensor data, an LLM is used to generate  
 264 behavior descriptions. These descriptions are subsequently used for policy and classifier learning.

265 **Setup.** We compare BLADE against the VLM-based baseline Robot-VILA. We omit SayCan and  
 266 T2M-Shooting since they require additional training data. We first test the original action sequences  
 267 seen in the demonstrations for each domain. We then test on tasks that require novel compositions of  
 268 behaviors for four types of generalizations, i.e., unseen initial condition, state perturbation, geometric  
 269 constraints, and partial observability. For each generalization type, we run six experiments and report  
 270 the number of experiments that have been successfully completed.

271 **Results.** In Fig. 5, we show that our model is able to successfully complete at least 4/6 tasks for all  
 272 generalization types in the two different domains. In comparison, Robot-VILA struggles to generate



**Figure 5: Domains and Results in Real World.** **Make Tea** features a toy kitchen designed to simulate boiling water on a stove. The robot must assess the available space on the stove for the kettle. It also needs to manage the dependencies between actions, such as the faucet must be turned away before the kettle can be placed into the sink to avoid collisions. **Boil Water** involves a tabletop task aimed at preparing tea, incorporating a cabinet, a drawer, and a stove. The robot must locate the kettle, potentially hidden within the cabinet, and a teabag in the drawer. Additionally, it must consider geometric constraints by removing obstacles that block the cabinet doors. In both environments, our model significantly outperforms the VLM-based planner Robot-VILA.



**Figure 6: Real World Planning and Execution.** We show the execution traces from BLADE and Robot-VILA for two generalization tasks: (a) partial observability and (b) geometric constraints.

273 correct plans to complete the tasks. In Fig. 6, we visualize the generated plans and execution traces  
 274 of both models. In example A, we show that BLADE can find the kettle initially hidden in the cabinet  
 275 and then complete the rest of the task. In comparison, Robot-VILA directly predicts placing the  
 276 teabag in the kettle when the kettle is not visible, resulting in a failure.

## 277 6 Conclusion and Discussion

278 BLADE is a novel framework for long-horizon manipulation by integrating model-based planning and  
 279 imitation learning. BLADE uses an LLM to generate behavior descriptions with preconditions and  
 280 effects from language-annotated demonstrations and automatically generates state abstraction labels  
 281 based on behavior descriptions for learning state classifiers. At performance time, BLADE generalizes  
 282 to novel states and goals by composing learned behaviors with a planner. Compared to latent-space  
 283 and LLM/VLM-based planners, BLADE successfully completes significantly more long-horizon tasks  
 284 with various types of generalizations.

285 **Limitations.** One limitation of BLADE is that the automatic segmentation of demonstrations is based  
 286 on gripper states; more advanced contact detection techniques might be required for certain tasks such  
 287 as caging grasps. We also assume the knowledge of a given set of predicate names in natural language  
 288 and focus on learning dependencies between actions using the given predicates. Automatically  
 289 inventing task-specific predicates from demonstrations and language annotations, possibly with the  
 290 integration of vision-language models (VLMs) is an important future direction. In our experiments,  
 291 we also found that noisy state classification led to some planning failures. Therefore, developing  
 292 planners that are more robust to noises in state estimation is necessary. Finally, achieving novel  
 293 compositions of behaviors also requires policies with strong generalization to novel environmental  
 294 states, which remain a challenge for skills learned from a limited amount of demonstration data.

295 **References**

- 296 [1] C. Chi, S. Feng, Y. Du, Z. Xu, E. Cousineau, B. Burchfiel, and S. Song. Diffusion policy:  
297 Visuomotor policy learning via action diffusion. In *RSS*, 2023. 1, 5, 17
- 298 [2] C. R. Garrett, T. Lozano-Pérez, and L. P. Kaelbling. PDDLStream: Integrating Symbolic  
299 Planners and Blackbox Samplers via Optimistic Adaptive Planning. In *ICAPS*, 2020. 1, 2
- 300 [3] D. Xu, A. Mandlekar, R. Martín-Martín, Y. Zhu, S. Savarese, and L. Fei-Fei. Deep affordance  
301 foresight: Planning through what can be done in the future. In *ICRA*, 2021.
- 302 [4] H. Shi, H. Xu, Z. Huang, Y. Li, and J. Wu. RoboCraft: Learning to see, simulate, and shape  
303 elasto-plastic objects in 3d with graph networks. *IJRR*, 43(4):533–549, 2024. 1
- 304 [5] C. Lynch, M. Khansari, T. Xiao, V. Kumar, J. Tompson, S. Levine, and P. Sermanet. Learning  
305 latent plans from play. In *CoRL*, 2020. 1
- 306 [6] A. Brohan, Y. Chebotar, C. Finn, K. Hausman, A. Herzog, D. Ho, J. Ibarz, A. Irpan, E. Jang,  
307 R. Julian, et al. Do as I can, not as I say: Grounding language in robotic affordances. In *CoRL*,  
308 2023. 1, 2, 6
- 309 [7] S. Srivastava, E. Fang, L. Riano, R. Chitnis, S. Russell, and P. Abbeel. Combined Task and  
310 Motion Planning through an Extensible Planner-Independent Interface Layer. In *ICRA*, 2014. 2
- 311 [8] N. T. Dantam, Z. K. Kingston, S. Chaudhuri, and L. E. Kavraki. Incremental task and motion  
312 planning: A constraint-based approach. In *RSS*, 2016.
- 313 [9] M. Toussaint. Logic-Geometric Programming: An optimization-based approach to combined  
314 task and motion planning. In *IJCAI*, 2015. 2
- 315 [10] C. Finn and S. Levine. Deep visual foresight for planning robot motion. In *ICRA*, 2017. 2
- 316 [11] S. Nair and C. Finn. Hierarchical foresight: Self-supervised learning of long-horizon tasks via  
317 visual subgoal generation. In *ICLR*, 2020.
- 318 [12] H. Shi, H. Xu, S. Clarke, Y. Li, and J. Wu. Robocook: Long-horizon elasto-plastic object  
319 manipulation with diverse tools. In *CoRL*, 2023.
- 320 [13] A. Simeonov, Y. Du, B. Kim, F. Hogan, J. Tenenbaum, P. Agrawal, and A. Rodriguez. A long  
321 horizon planning framework for manipulating rigid pointcloud objects. In *CoRL*, 2021.
- 322 [14] X. Lin, C. Qi, Y. Zhang, Z. Huang, K. Fragkiadaki, Y. Li, C. Gan, and D. Held. Planning  
323 with spatial and temporal abstraction from point clouds for deformable object manipulation. In  
324 *CoRL*, 2022.
- 325 [15] Y. Du, M. Yang, P. Florence, F. Xia, A. Wahid, B. Ichter, P. Sermanet, T. Yu, P. Abbeel, J. B.  
326 Tenenbaum, et al. Video language planning. *arXiv:2310.10625*, 2023. 2
- 327 [16] J. Luo, C. Xu, X. Geng, G. Feng, K. Fang, L. Tan, S. Schaal, and S. Levine. Multi-stage cable  
328 routing through hierarchical imitation learning. *IEEE Transactions on Robotics*, 2024. 2
- 329 [17] L. X. Shi, Z. Hu, T. Z. Zhao, A. Sharma, K. Pertsch, J. Luo, S. Levine, and C. Finn. Yell at your  
330 robot: Improving on-the-fly from language corrections. *arXiv:2403.12910*, 2024.
- 331 [18] S. Pirk, K. Hausman, A. Toshev, and M. Khansari. Modeling long-horizon tasks as sequential  
332 interaction landscapes. In *CoRL*, 2020.
- 333 [19] C. Wang, L. Fan, J. Sun, R. Zhang, L. Fei-Fei, D. Xu, Y. Zhu, and A. Anandkumar. Mimicplay:  
334 Long-horizon imitation learning by watching human play. In *CoRL*, 2023.

- 335 [20] C. Lynch and P. Sermanet. Language conditioned imitation learning over unstructured data. In  
336 *RSS*, 2021. 2
- 337 [21] Z. Zhang, Y. Li, O. Bastani, A. Gupta, D. Jayaraman, Y. J. Ma, and L. Weihs. Universal Visual  
338 Decomposer: Long-horizon manipulation made easy. In *ICRA*, 2024.
- 339 [22] Y. Zhu, P. Stone, and Y. Zhu. Bottom-up skill discovery from unsegmented demonstrations  
340 for long-horizon robot manipulation. *IEEE Robotics and Automation Letters*, 7(2):4126–4133,  
341 2022. 2
- 342 [23] A. Curtis, X. Fang, L. P. Kaelbling, T. Lozano-Pérez, and C. R. Garrett. Long-horizon manipu-  
343 lation of unknown objects via task and motion planning with estimated affordances. In *ICRA*,  
344 2022. 2
- 345 [24] D. Driess, O. Oguz, J.-S. Ha, and M. Toussaint. Deep visual heuristics: Learning feasibility of  
346 mixed-integer programs for manipulation planning. In *ICRA*, 2020. 2
- 347 [25] Y. Zhu, J. Tremblay, S. Birchfield, and Y. Zhu. Hierarchical planning for long-horizon manipu-  
348 lation with geometric and symbolic scene graphs. In *ICRA*, 2020. 2
- 349 [26] D.-A. Huang, S. Nair, D. Xu, Y. Zhu, A. Garg, L. Fei-Fei, S. Savarese, and J. C. Niebles. Neural  
350 task graphs: Generalizing to unseen tasks from a single video demonstration. In *CVPR*, 2019.
- 351 [27] D.-A. Huang, D. Xu, Y. Zhu, A. Garg, S. Savarese, F.-F. Li, and J. C. Niebles. Continuous  
352 relaxation of symbolic planner for one-shot imitation learning. In *IROS*, 2019. 2
- 353 [28] W. Huang, P. Abbeel, D. Pathak, and I. Mordatch. Language models as zero-shot planners:  
354 Extracting actionable knowledge for embodied agents. In *ICML*, 2022. 2
- 355 [29] T. Silver, V. Hariprasad, R. S. Shuttlesworth, N. Kumar, T. Lozano-Pérez, and L. P. Kaelbling.  
356 Pddl planning with pretrained large language models. In *NeurIPS 2022 foundation models for*  
357 *decision making workshop*, 2022. 2
- 358 [30] K. Valmeekam, M. Marquez, S. Sreedharan, and S. Kambhampati. On the planning abilities of  
359 large language models—a critical investigation. In *NeurIPS*, 2023. 2
- 360 [31] S. Kambhampati, K. Valmeekam, L. Guan, K. Stechly, M. Verma, S. Bhambri, L. Saldyt, and  
361 A. Murthy. Llms can’t plan, but can help planning in llm-modulo frameworks. *arXiv:2402.01817*,  
362 2024. 2
- 363 [32] Y. Chen, J. Arkin, Y. Zhang, N. Roy, and C. Fan. AutoTAMP: Autoregressive task and motion  
364 planning with llms as translators and checkers. In *ICRA*, 2024. 2
- 365 [33] B. Liu, Y. Jiang, X. Zhang, Q. Liu, S. Zhang, J. Biswas, and P. Stone. LLM+P: Empowering  
366 large language models with optimal planning proficiency. *arXiv:2304.11477*, 2023.
- 367 [34] Y. Xie, C. Yu, T. Zhu, J. Bai, Z. Gong, and H. Soh. Translating natural language to planning  
368 goals with large-language models. *arXiv:2302.05128*, 2023.
- 369 [35] A. Mavrogiannis, C. Mavrogiannis, and Y. Aloimonos. Cook2ltl: Translating cooking recipes  
370 to ltl formulae using large language models. In *ICRA*, 2024. 2
- 371 [36] T. Silver, S. Dan, K. Srinivas, J. B. Tenenbaum, L. Kaelbling, and M. Katz. Generalized  
372 planning in PDDL domains with pretrained large language models. In *AAAI*, 2024. 2
- 373 [37] X. Zhu, Y. Chen, H. Tian, C. Tao, W. Su, C. Yang, G. Huang, B. Li, L. Lu, X. Wang, et al. Ghost  
374 in the Minecraft: Generally capable agents for open-world environments via large language  
375 models with text-based knowledge and memory. *arXiv:2305.17144*, 2023. 2



- 376 [38] K. Nottingham, P. Ammanabrolu, A. Suhr, Y. Choi, H. Hajishirzi, S. Singh, and R. Fox. Do  
377 embodied agents dream of pixelated sheep: Embodied decision making using language guided  
378 world modelling. In *ICML*, 2023. 2
- 379 [39] S. Hao, Y. Gu, H. Ma, J. J. Hong, Z. Wang, D. Z. Wang, and Z. Hu. Reasoning with language  
380 model is planning with world model. In *EMNLP*, 2023. 2
- 381 [40] K. Lin, C. Agia, T. Migimatsu, M. Pavone, and J. Bohg. Text2motion: From natural language  
382 instructions to feasible plans. *Autonomous Robots*, 47(8):1345–1365, 2023. 2, 6
- 383 [41] M. Skreta, Z. Zhou, J. L. Yuan, K. Darvish, A. Aspuru-Guzik, and A. Garg. Replan: Robotic  
384 replanning with perception and language models. *arXiv:2401.04157*, 2024. 2
- 385 [42] D. Driess, F. Xia, M. S. Sajjadi, C. Lynch, A. Chowdhery, B. Ichter, A. Wahid, J. Tompson,  
386 Q. Vuong, T. Yu, et al. Palm-e: An embodied multimodal language model. *arXiv:2303.03378*,  
387 2023. 2
- 388 [43] Z. Wu, Z. Wang, X. Xu, J. Lu, and H. Yan. Embodied task planning with large language models.  
389 *arXiv:2307.01848*, 2023.
- 390 [44] J. Xiang, T. Tao, Y. Gu, T. Shu, Z. Wang, Z. Yang, and Z. Hu. Language models meet world  
391 models: Embodied experiences enhance language models. In *NeurIPS*, 2024. 2
- 392 [45] Y. Hu, F. Lin, T. Zhang, L. Yi, and Y. Gao. Look before you leap: Unveiling the power of  
393 GPT-4v in robotic vision-language planning. *arXiv:2311.17842*, 2023. 2, 6
- 394 [46] N. Wake, A. Kanehira, K. Sasabuchi, J. Takamatsu, and K. Ikeuchi. ChatGPT empowered  
395 long-step robot control in various environments: A case application. *IEEE Access*, 2023. 2
- 396 [47] L. Wong, J. Mao, P. Sharma, Z. S. Siegel, J. Feng, N. Korneev, J. B. Tenenbaum, and J. Andreas.  
397 Learning adaptive planning representations with natural language guidance. In *ICLR*, 2024. 2
- 398 [48] L. Guan, K. Valmeekam, S. Sreedharan, and S. Kambhampati. Leveraging pre-trained large  
399 language models to construct and utilize world models for model-based task planning. In  
400 *NeurIPS*, 2023.
- 401 [49] P. Smirnov, F. Joublin, A. Ceravola, and M. Gienger. Generating consistent PDDL domains  
402 with large language models. *arXiv:2404.07751*, 2024. 2
- 403 [50] V. Lifschitz. On the semantics of STRIPS. In M. Georgeff, Lansky, and Amy, editors, *Reasoning*  
404 *about Actions and Plans*, pages 1–9. Morgan Kaufmann, San Mateo, CA, 1987. 3
- 405 [51] C. Aeronautiques, A. Howe, C. Knoblock, I. D. McDermott, A. Ram, M. Veloso, D. Weld, D. W.  
406 SRI, A. Barrett, and D. Christianson. PDDL: The Planning Domain Definition Language, 1998.  
407 4
- 408 [52] T. Migimatsu and J. Bohg. Grounding predicates through actions. In *ICRA*, 2022. 5, 7
- 409 [53] J. Mao, T. Lozano-Pérez, J. Tenenbaum, and L. Kaelbling. PDSketch: Integrated domain  
410 programming, learning, and planning. In *NeurIPS*, 2022. 5
- 411 [54] S. Liu, Z. Zeng, T. Ren, F. Li, H. Zhang, J. Yang, C. Li, J. Yang, H. Su, J. Zhu, et al. Grounding  
412 dino: Marrying dino with grounded pre-training for open-set object detection. *arXiv:2303.05499*,  
413 2023. 5, 16
- 414 [55] O. Mees, L. Hermann, E. Rosete-Beas, and W. Burgard. Calvin: A benchmark for language-  
415 conditioned policy learning for long-horizon robot manipulation tasks. *RA-L*, 7:7327–7334,  
416 2021. 5
- 417 [56] O. Mees, L. Hermann, and W. Burgard. What matters in language conditioned robotic imitation  
418 learning over unstructured data. *RA-L*, 7:11205–11212, 2022. 6, 19

- 419 [57] Z. Zhang, Y. Li, O. Bastani, A. Gupta, D. Jayaraman, Y. J. Ma, and L. Weihs. Universal visual  
420 decomposer: Long-horizon manipulation made easy. In *ICRA*, 2024. 15
- 421 [58] W. Wan, Y. Zhu, R. Shah, and Y. Zhu. Lotus: Continual imitation learning for robot manipulation  
422 through unsupervised skill discovery. In *ICRA*, 2024. 15
- 423 [59] M.-H. Guo, J.-X. Cai, Z.-N. Liu, T.-J. Mu, R. R. Martin, and S.-M. Hu. Pct: Point cloud  
424 transformer. *Computational Visual Media*, 7:187–199, 2021. 16
- 425 [60] L. P. Kaelbling and T. Lozano-Pérez. Hierarchical task and motion planning in the now. In  
426 *ICRA*, 2011. 17, 18
- 427 [61] C. R. Garrett, C. Paxton, T. Lozano-Pérez, L. P. Kaelbling, and D. Fox. Online replanning in  
428 belief space for partially observable task and motion problems. In *ICRA*. IEEE, 2020. 18
- 429 [62] O. Mees, J. Borja-Diaz, and W. Burgard. Grounding language with visual affordances over  
430 unstructured data. In *ICRA*, 2023. 19

# Supplementary Material for Learning Compositional Behaviors from Demonstration and Language

This supplementary material provides additional details on the BLADE model, the simulation experiments, and qualitative examples. Section A provides a detailed description of the BLADE model, including the behavior description generation, predicate generation, abstract verification, automatic predicate annotation, classifier implementation, and policy implementation. Section B provides details on the simulation experiments, including the task design and baseline implementations. Section C provides qualitative examples of our method and baselines. Section D provides details of our setup of the real-robot experiment. Finally, Section E includes a full list of the prompts for the baselines used in the simulation experiments.

## A BLADE Details

### A.1 Behavior Description Generation with LLMs

In Listing 2, we show the behavior descriptions automatically generated by the LLM for the CALVIN domain. We also show the detailed prompt to the LLM for generating the behavior description. We break down the system prompt into four parts: definitions of primitive actions (Listings 3), definitions of predicates and environment context (Listings 4), an in-context example (Listings 5), and additional instructions (Listings 6). In Listings 7, we show one example of the specific user prompt that is used to generate the behavior description for *place-in-drawer*.

### A.2 Predicate Generation with LLMs

In our main experiment, we assume that the predicates for each task domain are provided in natural language. Here, we show that given the task definition and the environment context, a LLM can automatically generate the relevant predicates for the domain. Listing 8 shows the prompt for generating predicates; Listing 1 shows the generated predicates and the relevant actions for them. In Table A1, we compare the generated predicates with the predicates defined by the domain expert. We observe that the LLM is able to generate 13 out of 15 predicates that match closely with the expert-designed predicates. These predicates provide abstract representations for object states (e.g., *is-open*, *light-on*), relations between objects (e.g., *in-slider*, *in*), and robot-centric states (e.g., *holding*). The LLM incorrectly generates the predicate *next-to(?x, ?y)* to characterize the effects of the *push-left* and *push-right* actions, possibly due to ambiguities in the definition of the actions. Future work can explore interactive generation and refinement of predicates. Our algorithms are agnostic to the source of predicates and can flexibly generate action descriptions based on the given predicates and learn their groundings.

Listing 1: LLM Response for Predicate Generation.

```
### Unary Predicates
1. **(on-table ?x)** - Indicates that object ?x is on the table.
   - Relevant Actions: lift_red_block_table, lift_blue_block_table, lift_pink_block_table,
   place_on_table
2. **(in-slider ?x)** - Indicates that object ?x is inside the slider.
   - Relevant Actions: lift_red_block_slider, lift_blue_block_slider, lift_pink_block_slider,
   place_in_slider
3. **(in-drawer ?x)** - Indicates that object ?x is inside the drawer.
   - Relevant Actions: lift_red_block_drawer, lift_blue_block_drawer, lift_pink_block_drawer,
   place_in_drawer
4. **(is-open ?x)** - Indicates that container ?x (drawer or slider) is open.
   - Relevant Actions: open_drawer, close_drawer, place_in_drawer
5. **(is-closed ?x)** - Indicates that container ?x (drawer or slider) is closed.
   - Relevant Actions: close_drawer
6. **(light-on ?x)** - Indicates that the light (lightbulb or led) is turned on.
   - Relevant Actions: turn_off_lightbulb, turn_off_led
7. **(light-off ?x)** - Indicates that the light (lightbulb or led) is turned off.
   - Relevant Actions: turn_on_lightbulb, turn_on_led
8. **(holding ?x)** - Indicates that the robot is currently holding object ?x.
   - Relevant Actions: grasp, place, move, stack_block, unstack_block
9. **(clear ?x)** - Indicates that there is no object on top of ?x.
```

```

485 - Relevant Actions: stack_block, unstack_block, place_on_table
486 10. (rotated-right ?x) - Indicates that block ?x is rotated to the right.
487 - Relevant Actions: rotate_red_block_left, rotate_blue_block_left, rotate_pink_block_left
488 11. (rotated-left ?x) - Indicates that block ?x is rotated to the left.
489 - Relevant Actions: rotate_red_block_right, rotate_blue_block_right,
490 rotate_pink_block_right
491 12. (slider-left ?x) - Indicates that the slider ?x is in the leftmost position.
492 13. (slider-right ?x) - Indicates that the slider ?x is in the rightmost position.
493
494 ### Binary Relations
495 1. (on ?x ?y) - Indicates that object ?x is on top of object ?y.
496 - Relevant Actions: stack_block, unstack_block
497 2. (next-to ?x ?y) - Indicates that object ?x is next to object ?y.
498 - Relevant Actions: push_red_block_right, push_red_block_left, push_blue_block_right,
499 push_blue_block_left, push_pink_block_right, push_pink_block_left
500

```

### 501 A.3 Temporal Segmentation

502 Before the generation of behavior description, we segment each demonstration into a sequence of  
503 *contact-based primitives*. We consider seven primitives describing the interactions between the robot  
504 and other objects: *open/close* grippers without holding objects, *move-to(x)* which moves the gripper  
505 to an object, *grasp(x, y)* and *place(x, y)* which grasp and place object *x* from/onto another object *y*,  
506 *move(x)* which moves the currently holding object *x* and *push(x)*.

507 We use a set of heuristics to automatically segment the continuous trajectories using proprioception,  
508 i.e., gripper open state, and object segmentation. Specifically, *open* and *close* are directly detected by  
509 checking whether the gripper width is at the maximum or minimum value. *grasp(x, y)* and *place(x, y)*  
510 correspond to the other closing and opening gripper actions. *move(x)*, *push(x)* and *move-to(x)* are  
511 matched to temporal segments between pairs of gripper actions. Their type can be inferred based on  
512 the preceding and following gripper actions. We make a simplifying assumption that the robot moves  
513 freely in space only when the gripper is fully open and pushes objects only when the gripper is fully  
514 closed. These are given as instructions to the human demonstrators. In the simulator, the arguments  
515 of the primitives are obtained from the contact state. In the real world, they are inferred from the  
516 language annotations of the actions (e.g., “place the kettle on the stove” corresponds to *place(kettle,*  
517 *stove)*) procedurally or by the LLMs.

518 In Section 4.1, we discuss that we use LLMs to predict a *body* of contact primitive sequence associated  
519 with each behavior description. This additional step helps account for noises in the segmentation anno-  
520 tations, which are prevalent in CALVIN’s language-annotated demonstrations. For example, the lan-  
521 guage annotation “lift-block-table” correspond to the contact sequence  $\{move-to, grasp, move, place\}$ .  
522 Based on the generated *body*, the behavior can be correctly mapped to  $\{grasp, move\}$  and the demon-

**Table A1: Comparison of Predicates Defined by Domain Expert and Predicates Generated by an LLM.**

Manually Defined	Automatically Generated
<i>rotated-left(?x)</i>	<i>rotated-left(?x)</i>
<i>rotated-right(?x)</i>	<i>rotated-right(?x)</i>
<i>lifted(?x)</i>	<i>holding(?x)</i>
<i>is-open(?x)</i>	<i>is-open(?x)</i>
<i>is-close(?x)</i>	<i>is-closed(?x)</i>
<i>is-turned-on(?x)</i>	<i>light-on(?x)</i>
<i>is-turned-off(?x)</i>	<i>light-off(?x)</i>
<i>is-slider-left(?x)</i>	<i>slider-left(?x)</i>
<i>is-slider-right(?x)</i>	<i>slider-right(?x)</i>
<i>is-on(?x, ?y)</i>	<i>on-table(?x)</i>
<i>is-in(?x, ?y)</i>	<i>in-slider(?x), in-drawer(?x)</i>
<i>stacked(?x, ?y)</i>	<i>on(?x, ?y)</i>
<i>unstacked(?x, ?y)</i>	<i>clear(?x)</i>
<i>pushed-left(?x)</i>	-
<i>pushed-right(?x)</i>	-
-	<i>next-to(?x, ?y)</i>



523 stration trajectories can then be re-segmented. This additional step is crucial for learning accurate  
 524 groundings of the states and actions.

525 In our preliminary studies, we also experiment with other vision-based temporal segmentation  
 526 methods including UVD [57] and Lotus [58]. A main issue for incorporating these methods is that  
 527 they provide less consistent segmentations for different occurrences of the same behavior. As we  
 528 discussed in Section 6, more advanced contact detection techniques will be an important future  
 529 direction for using contact primitives as a meaningful interface between actions and language.

#### 530 A.4 Abstract Verification

531 After the generation of the behavior descriptions, we verify the generated behavior descriptions by  
 532 performing abstract verification on the demonstration trajectories. Given a segmented sequence of  
 533 the trajectory where each segment is associated with a behavior, we verify whether the preconditions  
 534 of each behavior can be satisfied by the accumulated effects of the previous behaviors. Pseudocode  
 535 for this algorithm is shown in Algorithm 1.

---

#### Algorithm 1 Abstract Verification

---

**Input:** Dataset  $\mathcal{D}$ , Behavior descriptions  $\mathcal{A}$

```

1: error_counter  $\leftarrow$  a counter for sequencing errors related to each behavior
2: counter  $\leftarrow$  a counter for storing the occurrences of each behavior
3: for  $i \leftarrow 1$  to  $K$  do
4:   obtain a behavior sequence  $\mathcal{D}_i \leftarrow \{a_1^i, \dots, a_N^i\}$ 
5:   initialize a dictionary for predicate state  $pred \leftarrow \{\}$ 
6:   for  $t \leftarrow 1$  to  $N$  do
7:     for each  $exp$  in  $pre_{a_t^i}$  do
8:        $(p, v) \leftarrow \text{EXTRACTPREDICATEANDBOOL}(exp)$ 
9:       if  $p$  not in  $pred$  then
10:         $pred[p] \leftarrow v$ 
11:       else
12:        if  $pred[p] \neq v$  then
13:         increment  $error\_counter[a_t^i]$ 
14:       for each  $exp$  in  $eff_{a_t^i}$  do
15:         $(p, v) \leftarrow \text{EXTRACTPREDICATEANDBOOL}(exp)$ 
16:         $pred[p] \leftarrow v$ 
17:        increment  $counter[a_t^i]$ 
18:   for each  $a$  in  $error\_counter$  do
19:     if  $error\_counter[a]/counter[a] > threshold$  then
20:       regenerate the behavior description for  $a$ 

```

---

#### 536 A.5 Automatic Predicate Annotation

537 We leverage *all* behavior descriptions to automatically label an observation  $\bar{o} = \{o_1, \dots, o_H\}$  based  
 538 on its associated segmentation. In particular, at  $o_0$ , we label all state predicates as “unknown.” Next,  
 539 we unroll the sequence of behavior executed in  $\bar{o}$ . As illustrated in Fig. 3c, before applying a behavior  
 540  $a$  at step  $o_t$ , we label all predicates in  $pre_a$  true. When  $a$  finishes at step  $o_{t'}$ , we label all predicates in  
 541  $eff_a$ . In addition, we will propagate the labels for state predicates to later time steps until they are  
 542 explicitly altered by another behavior  $a$ . Pseudocode for this algorithm is shown in Algorithm 2.

#### 543 A.6 Classifier Implementation

544 Based on the state predicate dataset generated from behavior definitions, we train a set of state  
 545 classifiers  $f_\theta(p) : \mathcal{O} \rightarrow \{T, F\}$ , which are implemented as standard neural networks for classification.

546 In the simulation experiment, the classifier model is based on a pre-trained CLIP model (ViT-B/32).  
 547 We use the image pre-processing pipeline from the CLIP model to process the input images. We

---

**Algorithm 2** Predicate Annotation

---

**Input:** Behavior sequence  $\{a_1, \dots, a_N\}$ , Observation sequence  $\{o_1, \dots, o_H\}$ , Descriptions  $\mathcal{A}$

- 1:  $propagated \leftarrow$  an empty list of propagated predicates
- 2:  $prev\_effs \leftarrow$  a list for storing effects from previous step
- 3:  $timed\_preds \leftarrow$  an empty list of predicates associated with time steps
- 4:  $pred\_obs \leftarrow$  an empty list for storing predicates paired with observations
- 5: **for**  $t \leftarrow 1$  to  $N$  **do**
- 6:     *// Precondition*
- 7:      $timed\_preds \leftarrow timed\_preds \cup \text{GETTIMEDPREDICATES}(pre_{a_t}, t)$
- 8:      $timed\_preds \leftarrow timed\_preds \cup \text{GETTIMEDPREDICATES}(\neg eff_{a_t}, t)$
- 9:     *// Propagated*
- 10:    **for** each  $p$  in  $propagated$  **do**
- 11:      **if** not  $\text{ALTERED}(p, a_t)$  **then**
- 12:         $\text{UPDATE TIME}(p, t)$
- 13:      **else**
- 14:         $propagated.remove(p)$
- 15:         $timed\_preds.add(p)$
- 16:    *// Previous effects*
- 17:    **for** each  $p$  in  $prev\_effs$  **do**
- 18:      **if** not  $\text{ALTERED}(p, a_t)$  **then**
- 19:         $propagated.add(p)$
- 20:      **else**
- 21:         $timed\_preds.add(p)$
- 22:    *// Store effects for next step*
- 23:     $prev\_effs \leftarrow \text{GETTIMEDPREDICATES}(eff_{a_t}, t)$
- 24:  $timed\_preds.update(propagated)$
- 25:  $timed\_preds.update(prev\_effs)$
- 26: **for** each  $p$  in  $timed\_preds$  **do**
- 27:     $pred\_obs.update(\text{MATCHTIMEDPREDICATEWITHOBSERVATION}(p, \{o_1, \dots, o_H\}))$
- 28: **return**  $pred\_obs$

---

548 use images from the static camera in the simulation. We perform one additional step of image  
549 processing to mask out the robot arm, which we find in our preliminary experiment to help avoid  
550 overfitting. We do not use the global image embedding from the CLIP model, instead we extract the  
551 patch tokens from the output of the vision transformer. We downsize the concatenated patch tokens  
552 with a multilayer perceptron (MLP) and then concatenate with word embeddings of the predicate  
553 arguments (e.g., *red-block, table*). The final embedding is then passed through a predicate-specific  
554 MLP to output the logit for binary classification. The CLIP model is frozen, while all other learnable  
555 parameters are trained.

556 In the real-world experiment, we find that, with more limited data than simulation, the pre-trained  
557 CLIP model often overfits to spurious relations in the training images (e.g., the state of the faucet  
558 is entangled with the location of the kettle). We also experiment with a ResNet-50 model pre-  
559 trained on ImageNet and find similar behavior. To improve generalization, we choose to focus on  
560 relevant objects and regions. We achieve this by using segmented object point clouds. We use open  
561 vocabulary object detector Grounding-Dino [54] to detect objects given object names. The predicted  
562 2D bounding boxes are projected into 3D and used to extract regions of the point cloud surrounding  
563 each object. The point-cloud-based classifier is based on the shape classification model from the  
564 Point Cloud Transformer (PCT) [59]. We concatenate the segmented object point clouds and include  
565 one additional channel to indicate the identity of each point. The PCT is used to encode the combined  
566 point cloud and output the final logit. The PCT model is trained from scratch.

## 567 A.7 Policy Implementation

568 For each behavior, we train control policies  $\pi_\theta(a) : \mathcal{O} \rightarrow \mathcal{U}$ , implemented as a diffusion policy [1].  
569 We make three changes to the original implementation to facilitate chaining the learned behaviors.  
570 First, when training the model to predict the first raw action for each skill, we replace the history  
571 observations with observations sampled randomly from a temporal window prior to when the skill is  
572 executed, to avoid bias in the starting positions of the robot arm. Second, we perform biased sampling  
573 of the training sequences to ensure that the policy is trained on a diverse set of starting positions.  
574 Third, at the end of each training sequence, we append a sequence of zeros actions so the learned  
575 policy can learned to predicate termination. These strategies are implemented for both the simulation  
576 and the real world.

577 In simulation, we construct the point cloud of the scene using the RGB-D image from the frame-  
578 mounted camera. We then obtain segmented object point clouds for the relevant objects of each  
579 behavior (e.g., *table* and *block* for *pick-block-table*) with groundtruth segmentation masks from the  
580 PyBullet simulator. The segmented point clouds of the objects are concatenated to form the input  
581 point cloud observation. The model uses the PCT to encode a sequence of point clouds as history  
582 observations and uses another time-series transformer encoder to reason over the history observations  
583 and predict the next actions. The time-series transformer is similar in design to the transformer-based  
584 diffusion policy [1].

585 In the real world, we use RGB images from four stationary cameras mounted around the workspace  
586 and a wrist-mounted camera as input to an image-based diffusion policy model. The input is processed  
587 using five separate ResNet-34 encoder heads. The policy directly predicts the gripper pose in the  
588 world frame. We found the wrist-mounted camera to be particularly helpful in the real-world setup.

## 589 A.8 Planner Implementation

590 **Planning over geometric constraints.** Geometric constraints, specifically the collision-free con-  
591 straints for each action, are handled “in the now,” right before an action is executed. This is because  
592 in order to classify the geometric constraints, we would need to know the exact pose of all objects in  
593 the environments. However, we do not explicitly learn models for predicting the exact location of  
594 objects after executing certain behaviors.

595 Our approach to handle this is to process them in the now. It follows the hierarchical planning  
596 strategy [60]. In particular, the precondition for actions is an ordered list. In our case, there are two  
597 levels: the second level contains the geometric constraint preconditions and the first level contains the  
598 rest of the semantic preconditions. During planning, only the first set of preconditions will be added  
599 to the subgoal list. After we have finished planning for the first-level preconditions, we consider  
600 the second-level precondition for the first behavior in the resulting plan, by possibly moving other  
601 obstacles away.

602 As an example, let us consider the skill of opening the cabinet door. Its first-level precondition  
603 only considers the initial state of the cabinet door (i.e., it should be initially closed). It also has a  
604 second-level precondition stating that nothing else should be blocking the door. In the beginning, the  
605 planner only considers the first-level preconditions. When this behavior is selected to be executed  
606 next, the planner checks for the second-level precondition. If this non-blocking precondition is not  
607 satisfied in the current state, we will recursively call the planner to achieve it (which will generate  
608 actions that move the blocking obstacles away). If this precondition has already been satisfied, we  
609 will proceed to execute the policy associated with this *opening-cabinet-door* skill.

610 This strategy will work for scenarios where there is enough space for moving obstacles around and  
611 the robot does not need to make dedicated plans for arranging objects. In scenarios where space is  
612 tight and dedicated object placement planning is required, we can extend our framework to include  
613 the prediction of object poses after each skill execution.

614 **Planning over partial observability.** Partial observability is handled assuming the most likely state.  
615 In particular, the effect definitions for all behaviors are deterministic. It denotes the most likely

616 state that it will result in. For example, in the definition of behaviors for finding objects (e.g., the  
617 *find-object-in-left-cabinet*), we have a deterministic and “optimistic” effect statement that the object  
618 will be visible after executing this action.

619 At performance time, since we will replan after executing each behavior, if the object is not visible  
620 after we have opened the left cabinet, the planner will automatically plan for other actions to achieve  
621 this visibility subgoal.

622 This strategy works for simple partially observable Markov decision processes (POMDPs). A  
623 potential extension to it is to model a belief state (e.g., representing a distribution of possible object  
624 poses) and execute belief updates on it. Planners can then use more advanced algorithms such as  
625 observation-based planning to generate plans. Such strategies have been studied in task and motion  
626 planning literature [60, 61].

## 627 B Simulation Experiment Details

### 628 B.1 Task Design

629 To evaluate generalization to new long-horizon manipulation tasks, we designed six tasks that fall  
630 into three categories: Abstract Goal, Geometric Constraint, and Partial Observability. Each task has a  
631 language instruction, a sampler that generates random initial states, and a goal satisfaction function  
632 for evaluation. We provide details for each task below.

#### 633 Task-1

- 634 • **Task Category:** Abstract Goal
- 635 • **Language Instruction:** *turn off all lights.*
- 636 • **Logical Goal:** (and (is-turned-off led) (is-turned-off lightbulb))
- 637 • **Initial State:** Both the led and the lightbulb are initially turned on.
- 638 • **Goal Satisfaction:** The logical states of both the lightbulb and the led are off.
- 639 • **Variation:** The initial states of the led and the lightbulb are both on and the goal is to turn them  
640 off.

#### 641 Task-2

- 642 • **Task Category:** Abstract Goal
- 643 • **Language Instruction:** *move all blocks to the closed drawer.*
- 644 • **Logical Goal:** (and (is-in red-block drawer) (is-in blue-block drawer) (is-in pink-block drawer))
- 645 • **Initial State:** The blocks are visible and not in the drawer. The drawer is closed.
- 646 • **Goal Satisfaction:** The blocks are in the drawer.

#### 647 Task-3

- 648 • **Task Category:** Abstract Goal
- 649 • **Language Instruction:** *move all blocks to the open drawer.*
- 650 • **Logical Goal:** (and (is-in red-block drawer) (is-in blue-block drawer) (is-in pink-block drawer))
- 651 • **Initial State:** The blocks are visible and not in the drawer. The drawer is open.
- 652 • **Goal Satisfaction:** The blocks are in the drawer.

#### 653 Task-4

- 654 • **Task Category:** Partial Observability
- 655 • **Language Instruction:** *place a red block on the table.*
- 656 • **Logical Goal:** (is-on red-block table)
- 657 • **Initial State:** The red block is in the drawer and the drawer is closed.
- 658 • **Goal Satisfaction:** The red block is placed on the table.
- 659 • **Variations:** Find the blue block or the pink block.

#### 660 Task-5

- 661 • **Task Category:** Partial Observability
- 662 • **Language Instruction:** *place a red block on the table.*



- 663 • **Logical Goal:** (is-on red-block table)
- 664 • **Initial State:** The red block is behind the sliding door.
- 665 • **Goal Satisfaction:** The red block is placed on the table.
- 666 • **Variations:** Find the blue block or the pink block.

#### 667 **Task-6**

- 668 • **Task Category:** Geometric Constraint
- 669 • **Language Instruction:** *open the slider.*
- 670 • **Logical Goal:** (is-slider-left slider)
- 671 • **Initial State:** The sliding door is on the right and there is a pink block on the path of the sliding
- 672 door to the left.
- 673 • **Goal Satisfaction:** The sliding door is within 5cm of the left end.
- 674 • **Variations:** Move the slider to the right.

## 675 **B.2 Baseline Implementation**

676 **HULC.** This baseline is a hierarchical policy learning method that learns from language-annotated  
 677 play data using hindsight labeling [56]. It’s one of the best-performing models on the  $D \rightarrow D$  split of  
 678 the CALVIN benchmark. We omit the comparison to the HULC++ method [62], the follow-up work  
 679 of HULC that leverages affordance prediction and motion planning to improve the low-level skills,  
 680 because our evaluation is focused on the task planning ability of the learned hierarchical model.

681 **SayCan.** This baseline combines an LLM-based planner that takes the language instruction and  
 682 learned feasibility functions for skills to perform task planning. We adopt SayCan to our learning-  
 683 from-play-data setting by training our own skill feasibility function by predicting possible next actions  
 684 to be executed at each state. The prompt of the model is listed in Listing 9.

685 **Robot-VILA.** This baseline performs task planning with a VLM. We adopt the prompts pro-  
 686 vided in the original paper to the CALVIN environment. The prompts are divided into the initial  
 687 prompt that is used to generate the task plan given the initial observation (shown in Listing 10)  
 688 and the follow-up prompt that is used for all subsequent steps (shown in Listing 11). We use  
 689 `gpt-4-turbo-2024-04-09` as the VLM. Because the model does not memorize the history. We  
 690 store the history dialogue, including the text input and the image input, and concatenate the history  
 691 dialogue with the current dialogue as the input to the VLM.

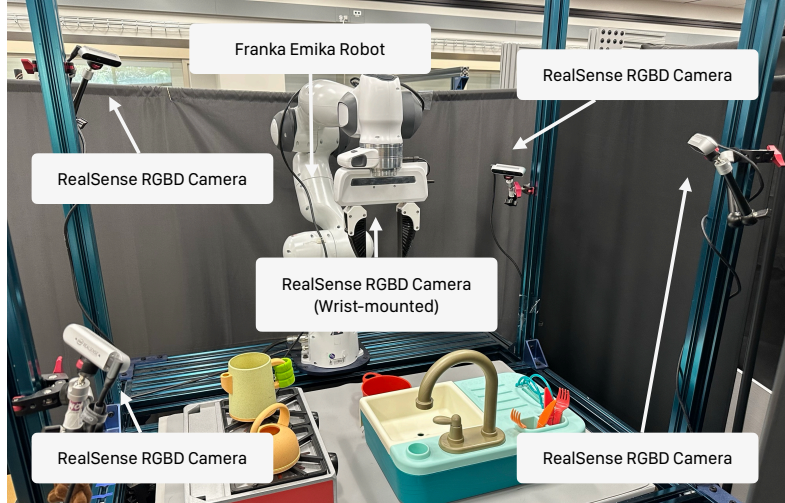
692 **T2M-Shooting.** This baseline (in particular, the shooting-based algorithm) is similar to the SayCan  
 693 algorithm except that: 1) it uses a multi-step feasibility model in contrast to the single-step feasibility  
 694 model used by SayCan; 2) the LLM additionally takes a symbolic state description of object states  
 695 and relationships. The original Text2Motion method assumes access to ground-truth symbolic states.  
 696 For comparison, in this paper, we compare Text2Motion with BLADE in two settings: one with the  
 697 ground-truth states and the other with the state classifiers learned by BLADE. The prompt of the  
 698 model is listed in Listing 12.

## 699 **C Qualitative Examples**

700 In this section, we include three qualitative examples from the CALVIN experiments to compare the  
 701 generalization capabilities of BLADE with baselines. Specifically, Fig. A2 shows generalization to  
 702 abstract goal, Fig. A3 shows generalization to partial observability, and Fig. A4 shows generalization  
 703 to geometric constraint. In summary, BLADE is able to generate accurate long-horizon manipulation  
 704 plans for novel situations while latent planning, LLM, and VLM baselines fail.

## 705 **D Real World Experiment Details**

706 As shown in Fig. A1, we employ a 7-degree of freedom (DOF) Franka Emika robotic arm equipped  
 707 with a parallel jaw gripper. A total of Five Intel RealSense RGB-D cameras are used to provide



**Figure A1:** We use a 7-degree of freedom (DOF) Franka Emika robotic arm with a parallel jaw gripper for our real-world experiment. A total of Five Intel RealSense RGB-D cameras are used to provide observation for our policies and state classifiers. Four cameras are mounted on the frame and an additional one is mounted to the robot’s wrist.

708 observation for our policies and state classifiers. Four cameras are mounted on the frame and one  
 709 additional camera is mounted on the robot’s wrist.

710 We use a teleoperation system with a 3DConnexion SpaceMouse for control. During the collection of  
 711 demonstrations, we record the robot’s joint configurations, the pose of the end effector, the gripper  
 712 width, and the RGB-D images from the five cameras. We collected approximately 80 demonstrations  
 713 for each of the two real-world domains, which provide the training data for the diffusion policy  
 714 models and the state classifiers.

715 Similar to our simulation experiments, our evaluation protocol includes the design of six tasks aimed  
 716 at assessing the model’s generalization capabilities across new long-horizon tasks. These tasks are  
 717 specifically crafted to test the model’s proficiency for four types of generalization: Unseen Initial  
 718 Condition, State Perturbation, Partial Observability, and Geometric Constraint.

719 **Task-1**

- 720 • **Domain:** Boil Water
- 721 • **Task Category:** Unseen Initial Condition
- 722 • **Language Instruction:** *Fill the kettle with water and place it on the stove*
- 723 • **Logical Goal:** (and (is-filled kettle) (is-placed-on kettle stove) (is-turned-off faucet-knob))
- 724 • **Initial State:** The kettle is placed inside the sink, and the stove is not blocked. The faucet is  
 725 turned off with the faucet head turned away.

726 **Task-2**

- 727 • **Domain:** Boil Water
- 728 • **Task Category:** State Perturbation
- 729 • **Language Instruction:** *Fill the kettle with water and place it on the stove*
- 730 • **Logical Goal:** (and (is-filled kettle) (is-placed-on kettle stove) (is-turned-off faucet-knob))
- 731 • **Initial State:** The kettle is placed inside the sink and the stove is blocked.
- 732 • **Perturbation:** The human user moves the kettle from the sink to the table after the robot turns  
 733 the faucet head towards the sink. The robot needs to replan to move the kettle back to the sink.

734 **Task-3**

- 735 • **Domain:** Boil Water
- 736 • **Task Category:** Geometric Constraint
- 737 • **Language Instruction:** *Fill the kettle with water and place it on the stove*

- 738 • **Logical Goal:** (and (is-filled kettle) (is-placed-on kettle stove) (is-turned-off faucet-knob))  
739 • **Initial State:** The kettle is placed inside the sink and the stove is blocked, creating a geometric  
740 constraint.

#### 741 **Task-4**

- 742 • **Domain:** Make Tea  
743 • **Task Category:** Unseen Initial Condition  
744 • **Language Instruction:** *Place the kettle on the stove and place the teabag inside the kettle.*  
745 • **Logical Goal:** (and (is-placed-on kettle stove) (is-placed-inside teabag kettle))  
746 • **Initial State:** The kettle is placed inside a cabinet. The cabinet doors are open. The drawer is  
747 closed.

#### 748 **Task-5**

- 749 • **Domain:** Make Tea  
750 • **Task Category:** State Perturbation  
751 • **Language Instruction:** *Place the kettle on the stove and place the teabag inside the kettle.*  
752 • **Logical Goal:** (and (is-placed-on kettle stove) (is-placed-inside teabag kettle))  
753 • **Initial State:** The kettle is placed inside the cabinet and the cabinet door is open. The drawer is  
754 initially closed.  
755 • **Perturbation:** Once the robot opens the drawer, a human user closes the drawer.

#### 756 **Task-6**

- 757 • **Domain:** Make Tea  
758 • **Task Category:** Geometric Constraint  
759 • **Language Instruction:** *Place the kettle on the stove and place the teabag inside the kettle.*  
760 • **Logical Goal:** (and (is-placed-on kettle stove) (is-placed-inside teabag kettle))  
761 • **Initial State:** There is a teapot blocking the cabinet doors. The kettle is inside the cabinet. The  
762 drawer is open with the teabag visible.

#### 763 **Task-7**

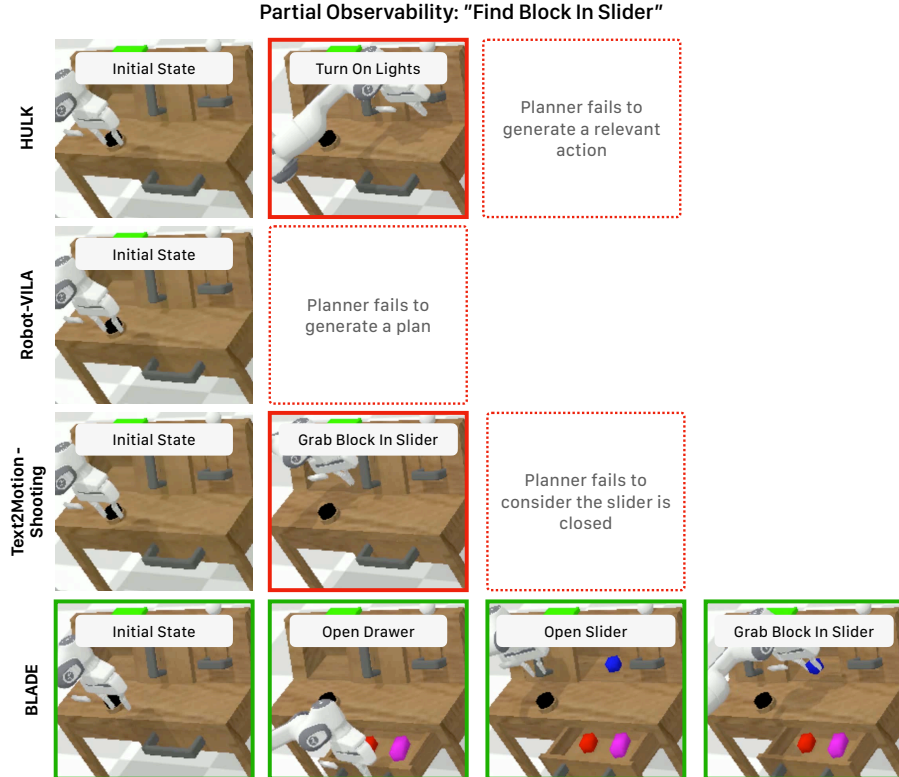
- 764 • **Domain:** Make Tea  
765 • **Task Category:** Partial Observability  
766 • **Language Instruction:** *Place the kettle on the stove and place the teabag inside the kettle.*  
767 • **Logical Goal:** (and (is-placed-on kettle stove) (is-placed-inside teabag kettle))  
768 • **Initial State:** The kettle is placed inside a cabinet and is not visible.

## 769 **E Prompts for Baselines**

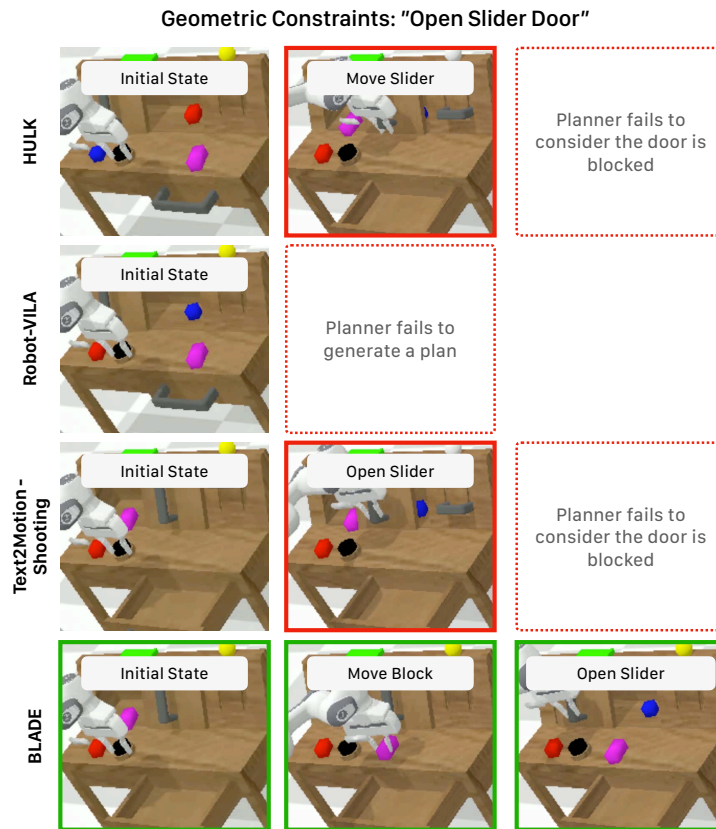
770 In this section, we provide the prompts for the baselines used in the simulation experiments. We  
771 provide the prompts for SayCan in Listing 9, Robot-VILA in Listing 10 and Listing 11, and T2M-  
772 Shooting in Listing 12.



**Figure A2:** BLADE and baseline performance on an Abstract Goal generalization task in the CALVIN environment.



**Figure A3:** BLADE and baseline performance on the Partial Observability generalization task in the CALVIN environment.



**Figure A4:** BLADE and baseline performance on the Geometric Constraint generalization task in the CALVIN environment.

## Listing 2: Behavior descriptions generated by the LLM for the CALVIN domain.

```
773
774 ;; lift_block_table
775 (:action lift-block-table
776 :parameters (?block - item ?table - item)
777 :precondition (and (is-block ?block) (is-table ?table) (is-on ?block ?table) (not (is-lifted
778 ?block)))
779 :effect (and (lifted ?block) (not (is-on ?block ?table)))
780 :body (then
781   (grasp ?block ?table)
782   (move ?block)
783 )
784 )
785
786 ;; lift_block_slider
787 (:action lift_block_slider
788 :parameters (?block - item ?slider - item)
789 :precondition (and (is-block ?block) (is-slider ?slider) (is-in ?block ?slider))
790 :effect (and (lifted ?block) (not (is-in ?block ?slider)))
791 :body (then
792   (grasp ?block ?slider)
793   (move ?block)
794 )
795 )
796
797 ;; lift_block_drawer
798 (:action lift-block-drawer
799 :parameters (?block - item ?drawer - item)
800 :precondition (and (is-block ?block) (is-drawer ?drawer) (is-in ?block ?drawer) (is-open ?
801 drawer))
802 :effect (and (lifted ?block) (not (is-in ?block ?drawer)))
803 :body (then
804   (grasp ?block ?drawer)
805   (move ?block)
806 )
807 )
808
809 ;; place_in_slider
810 (:action place-in-slider
811 :parameters (?block - item ?slider - item)
812 :precondition (and (is-block ?block) (is-slider ?slider) (is-lifted ?block))
813 :effect (and (is-in ?block ?slider) (not (is-lifted ?block)))
814 :body (then
815   (place ?block ?slider)
816 )
817 )
818
819 ;; place_in_drawer
820 (:action place-in-drawer
821 :parameters (?block - item ?drawer - item)
822 :precondition (and (is-block ?block) (is-drawer ?drawer) (is-lifted ?block) (is-open ?drawer)
823 )
824 :effect (and (is-in ?block ?drawer) (not (is-lifted ?block)))
825 :body (then
826   (place ?block ?drawer)
827 )
828 )
829
830 ;; place_on_table
831 (:action place-on-table
832 :parameters (?block - item ?table - item)
833 :precondition (and (is-block ?block) (is-table ?table) (is-lifted ?block))
834 :effect (and (is-on ?block ?table) (not (is-lifted ?block)))
835 :body (then
836   (place ?block ?table)
837 )
838 )
839
840 ;; stack_block
841 (:action stack_block
842 :parameters (?block - item ?target - item)
843 :precondition (and (is-block ?block) (is-block ?target) (is-lifted ?block))
844 :effect (and (stacked ?block ?target) (not (is-lifted ?block)))
845 :body (then
846   (place ?block ?target)
847 )
848 )
849
850
851 ;; unstack_block
852 (:action unstack_block
```



```

853 :parameters (?block1 - item ?block2 - item)
854 :precondition (and (is-block ?block1) (is-block ?block2) (stacked ?block1 ?block2))
855 :effect (and (unstacked ?block1 ?block2) (is-lifted ?block1) (not (stacked ?block1 ?block2)))
856 :body (then
857   (grasp ?block1 ?block2)
858   (move ?block1)
859 )
860 )
861
862 ;; rotate_block_right
863 (:action rotate-block-right
864 :parameters (?block - item ?table - item)
865 :precondition (and (is-block ?block) (is-table ?table) (is-on ?block ?table))
866 :effect (and
867   (rotated-right ?block)
868   (not (rotated-left ?block)))
869 :body (then
870   (grasp ?block ?table)
871   (move ?block)
872   (place ?block ?table)
873 )
874 )
875
876 ;; rotate_block_left
877 (:action rotate_block_left
878 :parameters (?block - item ?table - item)
879 :precondition (and (is-block ?block) (is-table ?table) (is-on ?block ?table))
880 :effect (and (rotated-left ?block))
881 :body (then
882   (grasp ?block)
883   (move ?block)
884   (place ?block)
885 )
886 )
887
888 ;; push_block_right
889 (:action push_block_right
890 :parameters (?block - item ?table - item)
891 :precondition (and (is-block ?block) (is-table ?table) (is-on ?block ?table))
892 :effect (and (pushed-right ?block) (not (pushed-left ?block)))
893 :body (then
894   (close)
895   (push ?block)
896   (open)
897 )
898 )
899
900 ;; push_block_left
901 (:action push-block-left
902 :parameters (?block - item)
903 :precondition (and (is-block ?block))
904 :effect (and (pushed-left ?block))
905 :body (then
906   (close)
907   (push ?block)
908   (open)
909 )
910 )
911
912 ;; move_slider_left
913 (:action move_slider_left
914 :parameters (?slider - item)
915 :precondition (and (is-slider ?slider) (is-slider-right ?slider))
916 :effect (and (is-slider-left ?slider) (not (is-slider-right ?slider)))
917 :body (then
918   (grasp ?slider)
919   (move ?slider)
920   (place ?slider)
921 )
922 )
923
924 ;; move_slider_right
925 (:action move-slider-right
926 :parameters (?slider - item)
927 :precondition (and (is-slider ?slider) (not (is-slider-right ?slider)))
928 :effect (and (is-slider-right ?slider))
929 :body (then
930   (grasp ?slider)
931   (move ?slider)
932   (place ?slider)
933 )

```

```

934 )
935
936 ;; open_drawer
937 (:action open-drawer
938 :parameters (?drawer - item)
939 :precondition (and (is-drawer ?drawer) (is-close ?drawer))
940 :effect (and (is-open ?drawer) (not (is-close ?drawer))))
941 :body (then
942   (close)
943   (push ?drawer)
944   (open)
945 )
946 )
947
948 ;; close_drawer
949 (:action close-drawer
950 :parameters (?drawer - item)
951 :precondition (and (is-drawer ?drawer) (is-open ?drawer))
952 :effect (and (is-close ?drawer) (not (is-open ?drawer))))
953 :body (then
954   (close)
955   (push ?drawer)
956   (open)
957 )
958 )
959
960 ;; turn_on_lightbulb
961 (:action turn-on-lightbulb
962 :parameters (?lightbulb - item)
963 :precondition (and (is-lightbulb ?lightbulb) (is-turned-off ?lightbulb))
964 :effect (and (is-turned-on ?lightbulb) (not (is-turned-off ?lightbulb))))
965 :body (then
966   (close)
967   (push ?lightbulb)
968   (open)
969 )
970 )
971
972 ;; turn_off_lightbulb
973 (:action turn-off-lightbulb
974 :parameters (?lightbulb - item)
975 :precondition (and (is-lightbulb ?lightbulb) (is-turned-on ?lightbulb))
976 :effect (and (is-turned-off ?lightbulb) (not (is-turned-on ?lightbulb))))
977 :body (then
978   (close) (push ?lightbulb) (open)
979 )
980 )
981
982 ;; turn_on_led
983 (:action turn-on-led
984 :parameters (?led - item)
985 :precondition (is-led ?led)
986 :effect (and (is-turned-on ?led) (not (is-turned-off ?led))))
987 :body (then
988   (close)
989   (push ?led)
990   (open)
991 )
992 )
993
994 ;; turn_off_led
995 (:action turn-off-led
996 :parameters (?led - item)
997 :precondition (and (is-led ?led) (is-turned-on ?led))
998 :effect (and (is-turned-off ?led) (not (is-turned-on ?led))))
999 :body (then
1000   (close)
1001   (push ?led)
1002   (open)
1003 )
1004 )
1005
1006 ;; push_into_drawer
1007 (:action push-into-drawer
1008 :parameters (?block - item ?drawer - item)
1009 :precondition (and (is-block ?block) (is-drawer ?drawer) (is-open ?drawer))
1010 :effect (and (is-in ?block ?drawer))
1011 :body (then
1012   (close)
1013   (push ?block)
1014   (open)

```

1015 )  
1019 )

---

### Listing 3: Example Prompt for CALVIN-Contact Primitives.

---

```
1018 **Primitive Actions:**
1019 There are seven primitive actions that the robot can perform. They are:
1020 - (grasp ?x ?y): ?x and ?y are two object variables. ?x is the object that the robot will be
1021 grasping, ?y is the object that ?x is currently on or in.
1022 - (place ?x ?y): ?x and ?y are two object variables. ?x is the object that the robot is
1023 currently holding, ?y is the object that ?x will be placed on or in.
1024 - (move ?x): ?x is the object that the robot is currently holding and will be moved by the
1025 robot.
1026 - (push ?x): ?x is the object that the robot will be pushing.
1027 - (move-to ?x): the robot arm will move without holding any object or pushing any object.
1028 - (open): the robot gripper will open fully.
1029 - (close): the robot gripper will close without grasping any object.
1030
1031 **Combined Primitives:**
1032 The primitive actions can be combined into a high-level routine. For example, (then (grasp ?x
1033 ?y) (move ?x) (place ?x ?y)) means the robot will pick up ?x from ?y, move ?x, and place ?x to
1034 ?z. The possible combination of primitives are:
1035 A. (then (grasp ?x ?y) (move ?x))
1036 B. (then (place ?x ?y))
1037 C. (then (grasp ?x ?y) (move ?x) (place ?x ?z))
1038 D. (then (close) (push ?x) (open))
1039
```

---

### Listing 4: Example Prompt for CALVIN-Environment.

---

```
1041 **Predicates for Preconditions and Effects:**
1042 The list of all possible predicates for defining the preconditions and effects of the high-
1043 level routine are listed below:
1044
1045 For specifying the type of the object:
1046 - (is-table ?x - item): ?x is a table.
1047 - (is-slider ?x - item): ?x is a slider.
1048 - (is-drawer ?x - item): ?x is a drawer.
1049 - (is-lightbulb ?x - item): ?x is a lightbulb.
1050 - (is-led ?x - item): ?x is a led.
1051 - (is-block ?x - item): ?x is a block.
1052
1053 For specifying the attributes of the object:
1054 - (is-red ?x - item): ?x is red. This predicate applies to a block.
1055 - (is-blue ?x - item): ?x is blue. This predicate applies to a block.
1056 - (is-pink ?x - item): ?x is pink. This predicate applies to a block.
1057
1058 For specifying the state of the object:
1059 - (rotated-left ?x - item): ?x is rotated left. This predicate applies to a block.
1060 - (rotated-right ?x - item): ?x is rotated right. This predicate applies to a block.
1061 - (pushed-left ?x - item): ?x is pushed left. This predicate applies to a block.
1062 - (pushed-right ?x - item): ?x is pushed right. This predicate applies to a block.
1063 - (lifted ?x - item): ?x is lifted in the air. This predicate applies to a block.
1064 - (is-open ?x - item): ?x is open. This predicate applies to a drawer.
1065 - (is-close ?x - item): ?x is close. This predicate applies to a drawer.
1066 - (is-turned-on ?x - item): ?x is turned on. This predicate applies to a lightbulb or a led.
1067 - (is-turned-off ?x - item): ?x is turned off. This predicate applies to a lightbulb or a led.
1068 - (is-slider-left ?x - item): the sliding door of the slider ?x is on the left.
1069 - (is-slider-right ?x - item): the sliding door of the slider ?x is on the right.
1070
1071 For specifying the relationship between objects:
1072 - (is-on ?x - item ?y - item): ?x is on top of ?y. This predicate applies when ?x is a block
1073 and ?y is a table.
1074 - (is-in ?x - item ?y - item): ?x is inside of ?y. This predicate applies when ?x is a block
1075 and ?y is a drawer or a slider.
1076 - (stacked ?x - item ?y - item): ?x is stacked on top of ?y. This predicate applies when ?x
1077 and ?y are blocks.
1078 - (unstacked ?x - item ?y - item): ?x is unstacked from ?y. This predicate applies when ?x and
1079 ?y are blocks.
1080
1081 **Task Environment:**
1082 In the environment where the demonstrations are being performed, there are the following
1083 objects:
1084 - A table. Objects can be placed on the table.
1085 - A drawer that can be opened. Objects can be placed into the drawer when it is open.
1086 - A slider which is a cabinet with a sliding door. The sliding door can be moved to the left
1087 or to the right. Objects can be placed into the slider no matter the position of the sliding
1088 door.
1089 - A lightbulb that be can turned on/off with a button.
1090 - A led that can be turned on/off with a button.
1091
```

1093 - Three blocks that can be rotated, pushed, lifted, and placed.

---

### Listing 5: Example Prompt for CALVIN-In-Context Example.

---

```
1094
1095 **Demonstration Parsing:**
1096 Now, you will help to parse several human demonstrations of the robot performing a task and
1097 generate a lifted description of how to accomplish this task.
1098 For each demonstration, a sequence of performed primitives will be given, with actual object
1099 names. Three demonstrations for the task of "place_in_slider" is:
1100
1101 <code name="primitive_sequence">
1102 primitives = [
1103   {"name": "grasp", "arguments": ["red_block", "table"]}
1104   {"name": "move", "arguments": ["red_block"]}
1105   {"name": "place", "arguments": ["red_block", "slider"]}
1106   {"name": "move-to", "arguments": [""]}
1107 ]
1108 </code>
1109
1110 <code name="primitive_sequence">
1111 primitives = [
1112   {"name": "grasp", "arguments": ["blue_block", "table"]}
1113   {"name": "move", "arguments": ["blue_block"]}
1114   {"name": "place", "arguments": ["blue_block", "slider"]}
1115   {"name": "move-to", "arguments": [""]}
1116 ]
1117 </code>
1118
1119 <code name="primitive_sequence">
1120 primitives = [
1121   {"name": "grasp", "arguments": ["pink_block", "table"]}
1122   {"name": "move", "arguments": ["pink_block"]}
1123   {"name": "place", "arguments": ["pink_block", "slider"]}
1124   {"name": "move-to", "arguments": [""]}
1125 ]
1126 </code>
1127
1128 **Previous Tasks:**
1129 A list of tasks that can be performed before the current task will also be provided as context
1130 . For the task of "place_in_slider", the possible previous tasks are:
1131 lift_block_table, lift_block_drawer, move_slider_right
1132
1133 **Example Output:**
1134 You should generate a lifted description, treating all objects as variables. For example, the
1135 lifted description for "place_in_slider" is:
1136 <code name="mechanism">
1137 (:mechanism place-in-slider
1138  :parameters (?block - item ?slider - item)
1139  :precondition (and (is-block ?block) (is-slider ?slider) (is-lifted ?block))
1140  :effect (and (is-in ?block ?slider) (not (is-lifted ?block)))
1141  :body (then
1142    (place ?block ?slider)
1143  )
1144 )
1145 </code>
```

---

### Listing 6: Example Prompt for CALVIN-Instructions.

---

```
1147
1148 **Think Step-by-Step:**
1149 To generate the lifted description, you should think through the task in natural language in
1150 the following steps. Be EXTREMELY CAREFUL to think through step 3a, 3b, and 4a, 4b.
1151 1. Parse the goal. For example "place_in_slider", the goal is to place a block into the slider
1152 .
1153 2. Think about the possible effects achieved by previous tasks and the previous actions that
1154 have been performed. For "lift_block_table", a block is lifted from the table and the effect
1155 is that the block is lifted. For "lift_block_drawer", a block is lifted from the drawer and
1156 the effect is that the block is lifted. For "move_slider_right", the sliding door of the
1157 slider is moved to the right and the effect is that the sliding door is on the right.
1158 3. Parse the demonstrations and choose the combination of primitives for the current task. The
1159 demonstrations are noisy so that the demonstrated primitive sequences may include extra
1160 primitive actions that are not necessary for the current task at the beginning or end. The
1161 extra primitive actions can be for the previous tasks. Combining with the understanding of the
1162 task and previous task to infer the correct combination of primitives for the current task.
1163 3a. In this case, the previous tasks are relevant to the current task. We should think about
1164 how to sequence the previous tasks with the current task. The primitive combination for the
1165 current task should not include primitive actions that have been performed. The above example
1166 for "place_in_slider" is this case. We can infer that "grasp" in the demonstrated sequences is
1167 likely to be for the previous tasks and should not be included in the primitive combination
```

```

1168 for the current task. We therefore choose B. (then (place ?x ?y)). The semantics is that the
1169 robot place the lifted block in the slider.
1170 3b. In this case, the previous tasks are not relevant to the current task.
1171 4. Think about the preconditions. Also specify the types of all relevant objects in the
1172 preconditions.
1173 4a. In this case, previous tasks are relevant to the current task. We should think about the
1174 effects of the previous tasks. For "place_in_slider", the effects of previous tasks include
1175 the block is already lifted. So we should specify that the block is lifted in the
1176 preconditions for the current task.
1177 4b. In this case, previous tasks are not relevant to the current task.
1178 5. Think about the effects. For "place_in_slider", the effects are that the block is in the
1179 slider and the block is not lifted.
1180 6. Write down the mechanism in the format of the example.
1181
1182 **Additional Instructions:**
1183 1. Make sure the generated lifted description starts with <code name="mechanism"> and ends
1184 with </code>.
1185 2. Please do not invent any new predicates for the precondition and effect. You can only use
1186 the predicates listed above.
1187 3. Consider the physical constraints of the objects. For example, a robot arm can not go
1188 through a closed door.
1189 4. For each parameter in :parameters, you should use one of the predicates for specifying the
1190 type of the object to indicate its type (e.g., is-drawer, is-block, and etc).

```

---

### Listing 7: Example Prompt for CALVIN-Task Input.

```

1192
1193 **Current Task:** place_in_drawer
1194
1195 **Example Sequences:**
1196 <code name="primitive_sequence">
1197 primitives = [
1198   {"name": "grasp", "arguments": ["blue_block", "table"]}
1199   {"name": "move", "arguments": ["blue_block"]}
1200   {"name": "place", "arguments": ["blue_block", "drawer"]}
1201   {"name": "move-to", "arguments": [""]}
1202 ]
1203 </code>
1204
1205 <code name="primitive_sequence">
1206 primitives = [
1207   {"name": "grasp", "arguments": ["red_block", "table"]}
1208   {"name": "move", "arguments": ["red_block"]}
1209   {"name": "place", "arguments": ["red_block", "drawer"]}
1210   {"name": "move-to", "arguments": [""]}
1211 ]
1212 </code>
1213
1214 <code name="primitive_sequence">
1215 primitives = [
1216   {"name": "grasp", "arguments": ["pink_block", "table"]}
1217   {"name": "move", "arguments": ["pink_block"]}
1218   {"name": "place", "arguments": ["pink_block", "drawer"]}
1219   {"name": "move-to", "arguments": [""]}
1220 ]
1221 </code>
1222
1223 **Previous Tasks:** push_into_drawer, lift_block_table, lift_block_slider

```

---

### Listing 8: Example Prompt for Predicate Generation.

```

1225
1226 You are a helpful agent in helping a robot interpret human demonstrations and discover a
1227 generalized high-level routine to accomplish a given task.
1228 **Primitive Actions:**
1229 There are seven primitive actions that the robot can perform. They are:
1230 - (grasp ?x ?y): ?x and ?y are two object variables. ?x is the object that the robot will be
1231 grasping, ?y is the object that ?x is currently on or in.
1232 - (place ?x ?y): ?x and ?y are two object variables. ?x is the object that the robot is
1233 currently holding, ?y is the object that ?x will be placed on or in.
1234 - (move ?x): ?x is the object that the robot is currently holding and will be moved by the
1235 robot.
1236 - (push ?x): ?x is the object that the robot will be pushing.
1237 - (move-to ?x): the robot arm will move without holding any object or pushing any object.
1238 - (open): the robot gripper will open fully.
1239 - (close): the robot gripper will close without grasping any object.
1240
1241 **Task Environment:**
1242 In the environment where the demonstrations are being performed, there are the following
1243 objects:
1244 - A table. Objects can be placed on the table.

```

```

1245 - A drawer that can be opened. Objects can be placed into the drawer when it is open.
1246 - A slider which is a cabinet with a sliding door. The sliding door can be moved to the left
1247 or to the right. Objects can be placed into the slider no matter the position of the sliding
1248 door.
1249 - A lightbulb that be can turned on/off with a button.
1250 - A led that can be turned on/off with a button.
1251 - Three blocks that can be rotated, pushed, lifted, and placed.
1252
1253 **Task**
1254 You will help the robot to write PDDL definitions for the following actions:
1255 1. lift_red_block_table
1256 2. lift_red_block_slider
1257 3. lift_red_block_drawer
1258 4. lift_blue_block_table
1259 5. lift_blue_block_slider
1260 6. lift_blue_block_drawer
1261 7. lift_pink_block_table
1262 8. lift_pink_block_slider
1263 9. lift_pink_block_drawer
1264 10. stack_block
1265 11. unstack_block
1266 12. place_in_slider
1267 13. place_in_drawer
1268 14. place_on_table
1269 15. rotate_red_block_right
1270 16. rotate_red_block_left
1271 17. rotate_blue_block_right
1272 18. rotate_blue_block_left
1273 19. rotate_pink_block_right
1274 20. rotate_pink_block_left
1275 21. push_red_block_right
1276 22. push_red_block_left
1277 23. push_blue_block_right
1278 24. push_blue_block_left
1279 25. push_pink_block_right
1280 26. push_pink_block_left
1281 27. move_slider_left
1282 28. move_slider_right
1283 29. open_drawer
1284 30. close_drawer
1285 31. turn_on_lightbulb
1286 32. turn_off_lightbulb
1287 33. turn_on_led
1288 34. turn_off_led
1289
1290 Before writing the operators, define the predicates that should be used to write the
1291 preconditions and effects of the operators. Group the predicates into unary predicates that
1292 define the states of objects and binary relations that specify relations between two objects.
1293 For each predicate, list actions that are relevant.

```

---

### Listing 9: Prompt for SayCan.

---

```

1295
1296 **Objective:**
1297 You are a helpful agent in helping a robot plan a sequence of actions to accomplish a given
1298 task.
1299 I will first provide context and then provide an example of how to perform the task.
1300
1301 **Task Environment:**
1302 In the robot's environment, there are the following objects:
1303 - A table. Objects can be placed on the table.
1304 - A drawer that can be opened. Objects can be placed into the drawer when it is open.
1305 - A slider which is a cabinet with a sliding door. The sliding door can be moved to the left
1306 or to the right. Objects can be placed into the slider no matter the position of the sliding
1307 door.
1308 - A lightbulb that be can turned on/off with a button.
1309 - A led that can be turned on/off with a button.
1310 - Three blocks that can be rotated, pushed, lifted, and placed.
1311
1312 **Actions:**
1313 There are the following actions that the robot can perform. They are:
1314 - lift_red_block_table: lift the red block from the table.
1315 - lift_red_block_slider: lift the red block from the slider.
1316 - lift_red_block_drawer: lift the red block from the drawer.
1317 - lift_blue_block_table: lift the blue block from the table.
1318 - lift_blue_block_slider: lift the blue block from the slider.
1319 - lift_blue_block_drawer: lift the blue block from the drawer.
1320 - lift_pink_block_table: lift the pink block from the table.
1321 - lift_pink_block_slider: lift the pink block from the slider.
1322 - lift_pink_block_drawer: lift the pink block from the drawer.
1323 - stack_block: stack the blocks.

```



```

1324 - place_in_slider: place the block in the slider.
1325 - place_in_drawer: place the block in the drawer.
1326 - place_on_table: place the block on the table.
1327 - rotate_red_block_right: rotate the red block to the right.
1328 - rotate_red_block_left: rotate the red block to the left.
1329 - rotate_blue_block_right: rotate the blue block to the right.
1330 - rotate_blue_block_left: rotate the blue block to the left.
1331 - rotate_pink_block_right: rotate the pink block to the right.
1332 - rotate_pink_block_left: rotate the pink block to the left.
1333 - push_red_block_right: push the red block to the right.
1334 - push_red_block_left: push the red block to the left.
1335 - push_blue_block_right: push the blue block to the right.
1336 - push_blue_block_left: push the blue block to the left.
1337 - push_pink_block_right: push the pink block to the right.
1338 - push_pink_block_left: push the pink block to the left.
1339 - move_slider_left: move the slider to the left.
1340 - move_slider_right: move the slider to the right.
1341 - open_drawer: open the drawer.
1342 - close_drawer: close the drawer.
1343 - turn_on_lightbulb: turn on the lightbulb.
1344 - turn_off_lightbulb: turn off the lightbulb.
1345 - turn_on_led: turn on the led.
1346 - turn_off_led: turn off the led.
1347 - do_nothing: do nothing.
1348
1349 **Example Task:**
1350 Now, you will help to parse the goal predicate and generate a list of candidate actions the
1351 robot can potentially take to accomplish the task. You should rank the actions in terms of how
1352 likely they are to be performed next.
1353 Goal predicate: (is-turned-off led)
1354 Task output:
1355 ```python
1356 ['turn_off_led', 'do_nothing']
1357 ```
1358 In this example above, if the led is on, the robot should turn it off. If the led is already
1359 off, the robot should do nothing.
1360
1361 **Additional Instructions:**
1362 1. Make sure the generated plan is a list of actions. Place the list between ```python and
1363 ends with ```.
1364 2. Think Step-by-Step.

```

---

### Listing 10: Initial Prompt for Robot-VILA.

```

1366 You are highly skilled in robotic task planning, breaking down intricate and long-term tasks
1367 into distinct primitive actions.
1368 If the object is in sight, you need to directly manipulate it. If the object is not in sight,
1369 you need to use primitive skills to find the object
1370 first. If the target object is blocked by other objects, you need to remove all the blocking
1371 objects before picking up the target object. At
1372 the same time, you need to ignore distracters that are not related to the task. And remember
1373 your last step plan needs to be "done".
1374
1375 Consider the following skills a robotic arm can perform.
1376 - lift_red_block_table: lift the red block from the table.
1377 - lift_red_block_slider: lift the red block from the slider.
1378 - lift_red_block_drawer: lift the red block from the drawer.
1379 - lift_blue_block_table: lift the blue block from the table.
1380 - lift_blue_block_slider: lift the blue block from the slider.
1381 - lift_blue_block_drawer: lift the blue block from the drawer.
1382 - lift_pink_block_table: lift the pink block from the table.
1383 - lift_pink_block_slider: lift the pink block from the slider.
1384 - lift_pink_block_drawer: lift the pink block from the drawer.
1385 - stack_block: stack the blocks.
1386 - place_in_slider: place the block in the slider.
1387 - place_in_drawer: place the block in the drawer.
1388 - place_on_table: place the block on the table.
1389 - rotate_red_block_right: rotate the red block to the right.
1390 - rotate_red_block_left: rotate the red block to the left.
1391 - rotate_blue_block_right: rotate the blue block to the right.
1392 - rotate_blue_block_left: rotate the blue block to the left.
1393 - rotate_pink_block_right: rotate the pink block to the right.
1394 - rotate_pink_block_left: rotate the pink block to the left.
1395 - push_red_block_right: push the red block to the right.
1396 - push_red_block_left: push the red block to the left.
1397 - push_blue_block_right: push the blue block to the right.
1398 - push_blue_block_left: push the blue block to the left.
1399 - push_pink_block_right: push the pink block to the right.
1400 - push_pink_block_left: push the pink block to the left.
1401 - move_slider_left: move the slider to the left.
1402

```

```

1403 - move_slider_right: move the slider to the right.
1404 - open_drawer: open the drawer.
1405 - close_drawer: close the drawer.
1406 - turn_on_lightbulb: turn on the lightbulb.
1407 - turn_off_lightbulb: turn off the lightbulb.
1408 - turn_on_led: turn on the led.
1409 - turn_off_led: turn off the led.
1410 - done: the goal has reached.
1411
1412 You are only allowed to use the provided skills. You can first itemize the task-related
1413 objects to help you plan.
1414 For the actions you choose, list them as a list in the following format.
1415
1416 <code>
1417 ['turn_off_led', 'open_drawer', 'done']
1418 </code>

```

---

### Listing 11: Follow-Up Prompt for Robot-VILA.

```

1420 This image displays a scenario after you have executed some steps from the plan generated
1421 earlier. When interacting with people,
1422 sometimes the robotic arm needs to wait for the person's action. If you do not find the target
1423 object in the current image, you need to
1424 continue searching elsewhere. Continue to generate the plan given the updated environment
1425 state.
1426

```

---

### Listing 12: Prompt for Text2Motion.

```

1428
1429 **Objective:**
1430 You are a helpful agent in helping a robot plan a sequence of actions to accomplish a given
1431 task.
1432 I will first provide context and then provide an example of how to perform the task.
1433
1434 **Task Environment:**
1435 In the robot's environment, there are the following objects:
1436 - A table. Objects can be placed on the table.
1437 - A drawer that can be opened. Objects can be placed into the drawer when it is open.
1438 - A slider which is a cabinet with a sliding door. The sliding door can be moved to the left
1439 or to the right. Objects can be placed into the slider no matter the position of the
1440 sliding door.
1441 - A lightbulb that be can turned on/off with a button.
1442 - A led that can be turned on/off with a button.
1443 - Three blocks that can be rotated, pushed, lifted, and placed.
1444
1445 **Predicates for symbolic state:**
1446 The list of all possible predicates for defining the symbolic state are listed below:
1447 - (rotated-left ?x - item): ?x is rotated left. This predicate applies to a block.
1448 - (rotated-right ?x - item): ?x is rotated right. This predicate applies to a block.
1449 - (pushed-left ?x - item): ?x is pushed left. This predicate applies to a block.
1450 - (pushed-right ?x - item): ?x is pushed right. This predicate applies to a block.
1451 - (lifted ?x - item): ?x is lifted in the air. This predicate applies to a block.
1452 - (is-open ?x - item): ?x is open. This predicate applies to a drawer.
1453 - (is-close ?x - item): ?x is close. This predicate applies to a drawer.
1454 - (is-turned-on ?x - item): ?x is turned on. This predicate applies to a lightbulb or a led.
1455 - (is-turned-off ?x - item): ?x is turned off. This predicate applies to a lightbulb or a
1456 led.
1457 - (is-slider-left ?x - item): the sliding door of the slider ?x is on the left.
1458 - (is-slider-right ?x - item): the sliding door of the slider ?x is on the right.
1459 - (is-on ?x - item ?y - item): ?x is on top of ?y. This predicate applies when ?x is a block
1460 and ?y is a table.
1461 - (is-in ?x - item ?y - item): ?x is inside of ?y. This predicate applies when ?x is a block
1462 and ?y is a drawer or a slider.
1463 - (stacked ?x - item ?y - item): ?x is stacked on top of ?y. This predicate applies when ?x
1464 and ?y are blocks.
1465 - (unstacked ?x - item ?y - item): ?x is unstacked from ?y. This predicate applies when ?x
1466 and ?y are blocks.
1467
1468 **Actions:**
1469 There are the following actions that the robot can perform. They are:
1470 - lift_red_block_table: lift the red block from the table.
1471 - lift_red_block_slider: lift the red block from the slider.
1472 - lift_red_block_drawer: lift the red block from the drawer.
1473 - lift_blue_block_table: lift the blue block from the table.
1474 - lift_blue_block_slider: lift the blue block from the slider.
1475 - lift_blue_block_drawer: lift the blue block from the drawer.
1476 - lift_pink_block_table: lift the pink block from the table.
1477 - lift_pink_block_slider: lift the pink block from the slider.
1478 - lift_pink_block_drawer: lift the pink block from the drawer.
1479 - stack_block: stack the blocks.

```

```

1480 - place_in_slider: place the block in the slider.
1481 - place_in_drawer: place the block in the drawer.
1482 - place_on_table: place the block on the table.
1483 - rotate_red_block_right: rotate the red block to the right.
1484 - rotate_red_block_left: rotate the red block to the left.
1485 - rotate_blue_block_right: rotate the blue block to the right.
1486 - rotate_blue_block_left: rotate the blue block to the left.
1487 - rotate_pink_block_right: rotate the pink block to the right.
1488 - rotate_pink_block_left: rotate the pink block to the left.
1489 - push_red_block_right: push the red block to the right.
1490 - push_red_block_left: push the red block to the left.
1491 - push_blue_block_right: push the blue block to the right.
1492 - push_blue_block_left: push the blue block to the left.
1493 - push_pink_block_right: push the pink block to the right.
1494 - push_pink_block_left: push the pink block to the left.
1495 - move_slider_left: move the slider to the left.
1496 - move_slider_right: move the slider to the right.
1497 - open_drawer: open the drawer.
1498 - close_drawer: close the drawer.
1499 - turn_on_lightbulb: turn on the lightbulb.
1500 - turn_off_lightbulb: turn off the lightbulb.
1501 - turn_on_led: turn on the led.
1502 - turn_off_led: turn off the led.
1503
1504 **Example Task:**
1505 Now, you will help to parse the goal predicate and generate a sequence of actions to
1506 accomplish this task.
1507 Goal predicate: (is-turned-off led)
1508 Symbolic state: is-turned-on(led), is-turned-on(lightbulb), not(is-turned-off(led)), not(is-
1509 turned-off(lightbulb))
1510 Task output:
1511 ```python
1512 ['turn_off_led']
1513 ```
1514
1515 **Example Task:**
1516 Goal predicate: (is-turned-on led)
1517 Symbolic state: is-turned-on(led), is-turned-on(lightbulb), not(is-turned-off(led)), not(is-
1518 turned-off(lightbulb))
1519 Task output:
1520 ```python
1521 []
1522 ```
1523
1524 **Example Task:**
1525 Goal predicate: (is-in red_block drawer)
1526 Symbolic state: not(is-in(red_block, drawer)), not(is-in(red_block, slider)), is-on(
1527 red_block, table), not(is-open(drawer)), is-close(drawer), is-slider-left(slider), not(is-
1528 slider-right(slider)), not(lifted(red_block))
1529 Task output:
1530 ```python
1531 ['open_drawer', 'lift_red_block_table', 'place_in_drawer']
1532 ```
1533
1534 **Example Task:**
1535 Goal predicate: (is-in red_block drawer)
1536 Symbolic state: not(is-in(red_block, drawer)), not(is-in(red_block, slider)), not(is-on(
1537 red_block, table)), is-open(drawer), not(is-close(drawer)), is-slider-left(slider), not(is-
1538 slider-right(slider)), lifted(red_block)
1539 Task output:
1540 ```python
1541 ['place_in_drawer']
1542 ```
1543
1544 **Example Task:**
1545 Goal predicate: (and (is-turned-on lightbulb) (is-slider-right slider))
1546 Symbolic state: is-slider-left(slider), not(is-slider-right(slider)), is-turned-off(
1547 lightbulb), not(is-turned-on(lightbulb))
1548 Task output:
1549 ```python
1550 ['turn_on_lightbulb', 'move_slider_right']
1551 ```
1552
1553 **Additional Instructions:**
1554 1. Make sure the generated plan is a list of actions. Place the list between ```python and
1555 ends with ```.
```

---