INFORMING REINFORCEMENT LEARNING AGENTS BY GROUNDING LANGUAGE TO MARKOV DECISION PRO CESSES

Anonymous authors

006

008 009 010

011

013

014

015

016

017

018

019

021

023

025

026 027 028

029

Paper under double-blind review

ABSTRACT

Natural language advice has the potential to accelerate reinforcement learning, but despite significant efforts to leverage natural language for RL, utilizing *diverse* forms of language efficiently remains unsolved. Existing methods focus on mapping natural language to individual elements of MDPs such as reward functions or policies, but such approaches limit the scope of language they consider to make such mappings possible. We propose to leverage general language advice by translating sentences to a grounded formal language for expressing information about *every* element of an MDP and its solution, including policies, plans, reward functions, and transition functions. We also introduce a new model-based reinforcement learning algorithm, RLang-Dyna-Q, capable of leveraging all such advice, and demonstrate in two sets of experiments that grounding language to every element of an MDP leads to significant performance gains. In additional symbol-grounding demonstrations we show how vision-language models can annotate important structure in the environent in the form of RLang vocabulary files, eliminating the need for human labels.

1 INTRODUCTION

i innobeenien

Language serves as a powerful means for humans to share information about the world, allowing us to learn more quickly or even skip learning altogether by drawing upon the domain expertise of others in the form of advice. An open question in reinforcement learning is how language advice can be leveraged to speed up learning in Markov Decision Processes (MDPs), as learning tasks *tabula rasa* is exceptionally difficult—and often impossible—in the real world. While many methods of leveraging advice for learning have emerged in the literature, a coherent theory of *language grounding* that can comprehensively support the use of language for reinforcement learning has not.

Virtually all research in language and RL grounds language to individual elements of MDPs such as policies (Liang et al., 2023; Vemprala et al., 2024; Wu et al., 2023; Andreas et al., 2017), reward functions (MacGlashan et al., 2015), and goals (Colas et al., 2020). The main drawbacks of these 040 works is that they restrict their approach to narrow fragments of natural language. For example, the 041 statement "if a mug is tipped over, its contents will spill out" clearly refers to a transition function, 042 and mapping this information to a policy is not straightforward. For this reason, works that ground 043 language to policies primarily focus on *imperative* sentences (e.g. "put the pallet on the truck") that 044 naturally correspond to policies, plans, or reward functions. Likewise, works that ground language to transition functions focus mainly on *declarative* sentences that provide information about the dynamics of a domain. This divergence in methodology suggests that not all language should be 046 grounded to the same component of an MDP, and that a general language grounding system for 047 reinforcement learning agents should be capable of grounding language to every element of an MDP, 048 and its solution. 049

We propose a novel approach to grounding natural language for use in reinforcement learning that formulates the language grounding problem as a machine translation task from natural language to RLang (Rodriguez-Sanchez et al., 2023), a formal language designed to express information about MDPs. Our approach is akin to semantic parsing (Mooney, 2007)—a problem in natural language understanding that involves translating natural language into a formal representation—because RLang



Figure 1: Translating natural language advice to RLang. We extend the original RLang pipeline (bottom) to include natural language translation and a Dyna-Q agent capable of leveraging all forms of RLang advice.

069 is a grounded formal language that offers a systematic means of expressing knowledge about an MDP. Such an approach calls for a learning agent capable of leveraging all such MDP components. 071 including a partial policy, reward function, plan, and transition function. We therefore also introduce 072 RLang-Dyna-Q, a model-based tabular RL agent based on Dyna-Q (Sutton et al., 1998), that can effectively leverage such advice. We demonstrate the strength and generality of our approach by 073 grounding a variety of natural language advice to RLang programs, which RLang-Dyna-Q can use to 074 significantly improve performance, sometimes making it possible to solve tasks that vanilla Dyna-Q 075 cannot solve. Our pipeline for these experiments relies on hand-specified RLang groundings that 076 generalize across tasks in the same domain (e.g. one grounding file for all Minigrid tasks), however, 077 we perform demonstrations showing how these groundings can instead be partially specified by a vision-language model. 079

2 BACKGROUND

054

056

058

060

065

066

067 068

081

083

084

085

087

090

091

098

103 104

105

Reinforcement learning tasks are typically modeled as Markov decision processes (MDPs), which can be represented by a tuple $\langle S, A, R, T, \gamma \rangle$, where S is the set of states, A is the set of actions, R is the reward function, T is the transition function, and γ is the discount factor. The goal of an agent is to find a policy, $\pi(a|s)$ —a function that selects an action for each state—which maximizes the expected sum of discounted rewards:

$$\mathbb{E}_{\pi}\left[\sum_{t=0}^{\infty}\gamma^{t}R(s_{t},a_{t},s_{t+1})\right].$$

Value-based reinforcement learning algorithms rely on estimating the optimal action-value function q_* , defined as 092

$$q_*(s,a) = \max_{\pi} q_{\pi}(s,a),$$

094 providing the expected return for taking action a in state s and subsequently following an optimal 095 policy (Sutton et al., 1998). Q-learning (Watkins, 1989) works to approximate q_* by applying the 096 following update rule after taking roll-outs in the environment:

$$Q(S_t, A_t) \leftarrow Q(S_t, A_t) + \alpha [R_{t+1} + \gamma \max_{a} Q(S_{t+1}, a) - Q(S_t, A_t)].$$

099 Building on Q-learning, Dyna-Q (Sutton et al., 1998) introduces an additional component: a model 100 of the environment. While Q-learning learns from direct interaction with the environment alone, 101 Dyna-Q builds an internal model of the environment and updates the action-value function using both 102 real and simulated roll-outs, enabling faster convergence to the optimal action-value function.

2.1 LEVERAGING FORMAL SPECIFICATION LANGUAGES FOR DECISION-MAKING

Formal specification languages have long been a useful tool to inform decision-making agents. In 106 classical planning, for example, it is standard to use the Planning Domain Description Language 107 (PDDL; Ghallab et al. 1998) and its probabilistic extension PPDDL (probabilistic PDDL; Younes & 108 Table 1: Selected MDP elements, corresponding RLang groundings, and natural language interpre-109 tations. The first column shows a component of the MDP, the second shows an RLang expression 110 about such component, and the last column contains a description of the expression.

MDP Component	RLang Declaration	Natural Language Interpre- tation
Policy $\pi: \mathcal{S} \times \mathcal{A} \to [0, 1]$	<pre>Policy build_bridge: if at_workbench: Execute use else: Execute go_to(workbench)</pre>	If you are at a workbench, use it. Otherwise, go to it.
Plan $\{A_0, A_1,, A_n\}$	<pre>Plan gather_materials: Execute go_to(wood) Execute pickup Execute go_to(string) Execute pickup</pre>	Go to the wood and pick it up, then go to the string and pick it up.
Reward, Transition Func. $R_e: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to \mathbb{R}$ $T_e: \mathcal{S} \times \mathcal{A} \times \mathcal{S} \to [0, 1]$	<pre>Effect common_sense: if at(Wall) and A == walk: Reward 0, S' -> S if at(Lava) and A == walk: Reward -1, S' -> S*0</pre>	Walking into walls will get you nowhere. Walking into lava will kill you.

127 Littman, 2004) to specify the complete dynamics of an environment. Other languages like Linear 128 Temporal Logic (LTL; Littman et al., 2017; Jothimurugan et al., 2019) and Policy Sketches (Andreas 129 et al., 2017) are sufficient for describing goals and hierarchical policies, respectively, for instruction-130 following agents. While effective, one limiting factor of these formal languages is their narrow scope. 131 Natural language, by contrast, can be used to express rich and varied information about nearly all 132 formal elements of decision-making.

133 RLang (Rodriguez-Sanchez et al., 2023) is a recent formal language to emerge from the literature. 134 While previous languages for decision-making narrowly focus on individual components of an 135 MDP such as a policy or reward function, RLang was designed to provide information about every 136 component of a structured MDP and its solution. Formally, an RLang specification is a set of RLang 137 groundings \mathcal{G} given by an RLang program \mathcal{P} and an RLang vocabulary \mathcal{V} , a file containing a set of 138 primitives that ground to important structured abstractions in the agent (e.g. as options, lifted skills, 139 etc.) and the environment (e.g. as objects, factors, etc.). Importantly, vocabulary files can be designed 140 once and reused across many tasks that occur in the same domain with little to no modifications. Some example RLang programs and their natural language interpretations can be seen in Table 1. 141 Crucially, advice specified by RLang can be compiled directly into many components of an MDP 142 including policies, transition functions, reward functions, and plans. Leveraging such components in a 143 learning algorithm is not always straightforward, however, and integrating more than one component 144 into an agent is a non-trivial problem that has not been addressed. 145

146

111

125 126

2.2 LEVERAGING NATURAL LANGUAGE FOR DECISION-MAKING 147

148 **Language in RL** Luketina et al. (2019) identify two variations of language usage in the reinforce-149 ment learning literature. The first, language-conditional RL, is one in which language use is a 150 necessary component of the task. This includes environments where agents must execute commands 151 in natural language (Mirchandani et al., 2021), or otherwise deal with language that is part of the 152 MDP, e.g., in the observation or action space (Fulda et al., 2017; Kostka et al., 2017). The second 153 variation is language-assisted RL, in which natural language is used to communicate task-related 154 information to an agent that is *not necessary* for solving the task. In these settings, language can be 155 used to inform policy structure (Watkins et al., 2021), reward functions (Goyal et al., 2019), transition dynamics (Narasimhan et al., 2018), or Q-functions (Branavan et al., 2012). 156

157

158 Grounding Natural to Formal Languages for Planning and Learning The notion of grounding 159 natural language to a formal language for use in learning and planning is not new. Gopalan et al. (2018) and Berg et al. (2020) translate natural language commands into Linear Temporal Logic 160 (LTL), which they use as reward functions for a learning agent or planning objectives, and Silver 161 et al. (2024) and Miglani & Yorke-Smith (2020) ground natural language into PDDL, which is fed to

162 a recurrent neural network to output solution plans. However, the advancement of large language 163 models (LLMs) has led to even more capable agents that for leveraging formal languages. In the 164 planning literature, Ahn et al. (2022); Huang et al. (2022); Song et al. (2023) use primitive formal 165 languages for executing policies on real robots or in embodied environments, Liu et al. (2023a); Xie 166 et al. (2023) translate natural language commands into PDDL plans with the help of LLMs, and Liu et al. (2023b) proposed a modular system to ground natural language into LTL formulas. Code is also 167 a popular choice for formal languages: Liang et al. (2023); Vemprala et al. (2024); Wu et al. (2023) 168 use an LLM to generate Python functions as policies from natural language instructions; Singh et al. (2023) also generates programs by prompting LLMs for code completion. For learning, more recent 170 works focus on reward design with LLMs for RL agents: Yu et al. (2023) specifies reward with LLMs 171 through code generation and Du et al. (2023) leverage commonsense reasoning for designing reward 172 functions. While many methods excel at grounding to formal languages Cohen et al. (2024), no 173 existing method seeks to ground language to every component of the MDP. 174

2.3 LARGE LANGUAGE MODELS FOR MACHINE TRANSLATION

175

176

180 181

191

192

193 194

196

197

199

200 201

202

203

204

177 Large Language Models (LLMs), often based on architectures like the Transformer (Vaswani 178 et al., 2017), are trained to predict the next token x_t in a sequence given the preceding tokens 179 $\{x_1, x_2, ..., x_{t-1}\}$ via the following objective:

$$\mathcal{L} = -\sum_{t} \log P(x_t | x_1, x_2, \dots, x_{t-1}).$$

182 In very large models, this objective results in emergent capabilities such as natural language under-183 standing and generation, making them suitable for a variety of tasks beyond mere text completion 184 including question-answering, summarization, and more (Bubeck et al., 2023). One useful emergent 185 capability of LLMs is the machine translation of text from one language to another. While specialized neural machine translation systems are trained using a parallel corpus to maximize the conditional probability P(y|x), where x is the source sequence and y is the target sequence (Bahdanau et al., 187 2015), LLMs have achieved similar translation capabilities despite not being trained explicitly on this 188 objective (Brown et al., 2020). Furthermore LLMs have been shown to be proficient at generating 189 text in formal languages such as Python given a language prompt (Chen et al., 2021; Li et al., 2023). 190

3 GROUNDING NATURAL LANGUAGE ADVICE TO RLANG PROGRAMS



Figure 2: LavaCrossing Experiment. The agent was given the following advice: "Walking into lava will kill you. Walking into walls will do nothing." The initial state of LavaCrossing is pictured left, reward curves are in the center, and the grounded RLang advice is on the right.

205 One major motivation for leveraging language advice in RL is to supply agents with the kinds of 206 commonsense reasoning that language can easily express. Consider the LavaCrossing environment in 207 Figure 2. Any human interacting with this environment would quickly learn that walking into the lava squares kills you, or likewise that walking into walls will do nothing at all. Communicating 208 this knowledge to others with language is natural for humans, but leveraging such language advice 209 in RL is a major unsolved problem. An alternative approach to supplying commonsense advice to 210 RL agents involves specifying it in a formal language relevant to decision-processes, which can 211 more straightforwardly be used by a learning agent to improve learning. While such an approach is 212 limited by the expressivity of the formal language and how it is used by the learning agent, it is fully 213 grounded, easily interpretable, and can be extremely powerful. 214

As formal languages for decision-making grow more expressive, a natural next step for leveraging language advice in reinforcement learning is to translate pieces of natural language advice into

216 statements in such formal languages. RLang is a highly expressive candidate for language grounding 217 because it is capable of specifying information about every element of a structured MDP and 218 its solution, including plans, policies, transition functions, and reward functions (see Table 2 in 219 Rodriguez-Sanchez et al. (2023)). Furthermore, we hypothesize that different kinds of advice can 220 most naturally be represented by different components of the MDP, and that methods that ground language to a single component are insufficient to capture general language advice. For example, 221 the statement, "stacked dishes can topple if unevenly piled," is precisely a statement about transition 222 dynamics, and while it can ultimately be used to inform a plan or policy, the information contained 223 in the statement would not be retrievable if it were not represented as a partial transition function; 224 the most harmonious representation of the advice is as a partial transition function. Likewise, the 225 sentence, "wear oven mitts whenever handling pots and pans," is a statement about a policy, and 226 representing it as a reward function would only indirectly capture its meaning. 227

We therefore formulate the language grounding problem in RL as a machine translation task from 228 natural language to RLang. Our task is as follows: given an RLang vocabulary \mathcal{V} —a set of task-229 general groundings that act as primitives in an RLang program-for a given MDP and a piece of 230 natural language advice u, we seek a function $\phi : u \times \mathcal{V} \to \mathcal{P}_u$, where \mathcal{P}_u is an executable RLang 231 program capturing the advice in u that can be leveraged by a learning agent. We propose to do 232 this translation entirely in-context using a general-purpose large language model in a two-stage 233 pipeline by 1) identifying which RLang grounding type would best capture the language advice; 234 and 2) translating the advice into an RLang program. Stage 1, the selection stage, instructs the 235 LLM to classify a novel piece of advice u into RLang grounding types such as Effects, Policies, and 236 Plans, consulting a small number of example classifications in the prompt. This ensures that the advice will be represented by an appropriate component of the MDP.¹ Stage 2, the translation stage, 237 instructs the LLM to translate u to an RLang program specifying the grounding type given by Stage 1 238 using roughly 5 example translations in the prompt that were hand-engineered to cover a wide range 239 of RLang's syntax. These programs are compiled using RLang's compiler into Python functions 240 corresponding to transition functions, reward functions, policies, and plans that can be leveraged by a 241 learning agent. In experiments we demonstrate that this pipeline effectively grounds the advice to 242 useful MDP components. Our pipeline is illustrated in Figure 1. 243

244 245

3.1 RLANG-DYNA-Q: A SINGLE AGENT FOR LEVERAGING ALL OF RLANG

246 In the original RLang paper, the authors presented a number of RLang-enabled agents—including 247 ones based on Q-Learning, PPO (Schulman et al., 2017), and DOORmax (Diuk et al., 2008)-each 248 capable of leveraging *individual* RLang groundings to improve learning. However, leveraging general 249 language advice requires integrating potentially all RLang groundings into a single learning agent. 250 We therefore introduce RLang-Dyna-Q, a learning agent based on Dyna-Q (Sutton et al., 1998) that is capable of simultaneously leveraging a partial policy, plan, reward function, and transition function 251 given by an RLang program. Similar to Dyna-Q, RLang-Dyna-Q leverages the Bellman update rule 252 to update Q-values using rollouts collected both from environment interaction and from simulated 253 interaction, which is generated from a partial model of the environment that is learned over time. 254 However, RLang-Dyna-Q also leverages a partial model given by an RLang program to generate 255 simulated rollouts before learning begins (see Algorithm 1, our modifications to Dyna-Q are in blue). 256 Dyna-Q is an appropriate core learning agent because integrating actions and dynamics is most 257 natural in a model-based learning algorithm that explicitly represents a policy, transition function, 258 and reward function.

259 260

261 262

4 EXPERIMENTS

We hypothesize that RLang is an effective grounding for natural language advice in the context of reinforcement learning. However, evaluating whether language advice u and RLang program \mathcal{P}_u have the same semantic content is difficult, so we designed our experiments to test the objective of primary interest: the agent's performance on a learning task. If we provide advice that is helpful to the agent, then grounding it properly should improve performance. We therefore assessed our translation pipeline by evaluating agent performance on multiple custom tasks based on the Minigrid/BabyAI

²⁶⁸ 269

¹We assume that each piece of advice—which may contain multiple sentences—grounds to a single RLang grounding type. This constraint can easily be relaxed in future work.

Algorithm 1 RLang-Dyna-Q Agent
Given: π_{RLang} , T_{RLang} , R_{RLang} from an RLang program
Init $Q(s, a), T(s, a), R(s, a)$ for all $s \in S, a \in \mathcal{A}(s)$
loop
$s \leftarrow \text{current (nonterminal) state}$
$a \leftarrow \epsilon_1, \epsilon_2$ -greedy $(s, \pi_{\text{RLang}}, Q)$ # With prob. ϵ_2 , we execute the RLang plan or policy
Execute action a ; observe next state s' , and reward r
$Q \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$
$T(s, a), R(s, a) \leftarrow s', r$ # Update our model
for $i = 1$ to N_1 do
$s \leftarrow$ random previously observed state
$a \leftarrow$ random action previously taken in s
$s', r \leftarrow T(s, a), R(s, a)$
$Q \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$
end for
for $i = 1$ to N_2 do
$s \leftarrow$ random previously observed state
$a \leftarrow random action not previously taken in s$
$s', r \leftarrow T_{\text{RLang}}(s, a), R_{\text{RLang}}(s, a)$ Predict s', r using dynamics given by RLang
$Q \leftarrow Q(s, a) + \alpha \left[r + \gamma \max_{a'} Q(s', a') - Q(s, a) \right]$
end for
end loop

(Chevalier-Boisvert et al., 2023; Chevalier-Boisvert et al., 2019) and VirtualHome (Puig et al., 2018) environments given expert advice. We include ablations of different RLang components to demonstrate our auxiliary hypothesis, that language is best grounded to *every* element of an MDP.
 Additionally, we run a small user study to assess our pipeline's efficacy on a wide range of non-expert language advice. Finally, we perform a series of symbol-grounding demonstrations showing how vision-language models can eliminate the need for human-readable RLang vocabulary files, which require human labeling.

299 300 301

292

4.1 LEVERAGING EXPERT ADVICE IN MINIGRID

302 We designed custom environments using the Minigrid/BabyAI library, a platform for studying the 303 behavior of language-informed agents. In a typical Minigrid environment, an agent might reason 304 about opening and closing doors using keys which may be hidden in other rooms, managing a small 305 inventory of items, removing obstacles like balls out of the way to reach other rooms or objects, and 306 avoiding lava, all for the ultimate purpose of reaching a goal. Minigrid environments are an ideal 307 setting for our experiments for three reasons: 1) they can be solved using tabular RL algorithms, 308 which our informed, model-based RLang-Dyna-Q agent is based on; 2) there are clear and obvious referents of language in both the state and action spaces of these environments (e.g. keys, doors, and 309 balls are represented neatly in a discrete state space and skills such as walking towards objects are 310 easy to implement); 3) many objects are shared across environments enabling the reuse of a common 311 RLang vocabulary for referencing these objects, which makes it easier for our translation pipeline to 312 ground novel advice. 313

We provide a domain-general RLang vocabulary file \mathcal{V} containing a set of RLang groundings to be used as primitives in a full RLang program. These vocabulary files are generated automatically for each minigrid environment given a single general template, and include perception abstractions such as the objects in the environment (e.g., yellow_key, red_door) and a short list of predicates for reasoning with them (e.g., carrying(), reachable(), at ()), as well as a single abstract action in the form of a lifted skill for walking to any reachable object (go_to()). Importantly, these groundings have semantically-meaningful labels, which enable a simple translation process.² All agents in the experiments, including the Random, Dyna-Q, and RLang-Dyna-Q agents, are

 ²In our final demonstrations, we relax this constraint by determining the semantic label of entity groundings
 and utilizing images of the objects in the environment to ground the referents of ambiguous advice using an off-the-shelf vision-language model.

324 given access to this lifted skill. However, we do not provide the Dyna-Q agent with any perception 325 abstractions, as including them induces an equivalent state space in the tabular RL setting. Likewise, 326 the Random agent does not consider state when selecting an action. In Stage 2 of the translation 327 pipeline, we provide the list of available RLang groundings that can be referenced in an RLang 328 program along with the language advice. This prevents the LLM from hallucinating imaginary skills, objects, and predicates when translating the advice into an RLang program. The LLM never interacts 329 with the MDP directly. The translation examples used in the prompts in both stages of translation 330 did not change across experiments, though these translations are vocabulary-specific and grounding 331 advice to environments outside of Minigrid will require domain-compatible example translations. 332

333 We evaluated our grounding pipeline on four diverse Minigrid environments: LavaCrossing, Multi-334 Room, MidMazeLava, and HardMaze. For each environment, we collected multiple pieces of natural language advice from human experts—people familiar with both the environment and how the agent 335 interacts with it via perception and action, i.e. the skills the agent has access to and the fact that 336 its perception consists of objects and simple predicates—and translated them into RLang programs 337 using our two-stage pipeline. Each piece of advice was translated to a single RLang grounding type, 338 and each piece of advice contained multiple sentences. We then evaluated our RLang-Dyna-Q agent 339 on each environment with the translated RLang programs. In the MidMazeLave environment, advice 340 was grounded to multiple RLang types. We include additional results for RLang-Dyna-Q utilizing 341 only one type of advice at a time-Effects, Plans, or Policies-to isolate their impact on performance. 342

RLang-Dyna-Q significantly outperformed vanilla Dyna-Q in all of the experiments. In LavaCrossing 343 (see Figure 2), the agent is tasked with reaching a goal while avoiding lava, and merely advising the 344 agent about the dangers of lava and futility of walking into walls greatly increases performance. In 345 the MultiRoom environment (see Figure 9), in which the agent must open a series of doors to reach 346 a goal, providing a plan in natural language significantly increased performance. In MidMazeLava 347 (see Figure 3) and HardMaze (see Figure 4), the agent is faced with significantly more difficult 348 tasks. In the former, the agent must unblock doors and open them with keys to reach a goal while 349 avoiding lava, and in the latter the agent must traverse through many rooms, bringing keys across 350 rooms to doors which must be unblocked to reach a goal. We collected paragraphs-worth of advice 351 for these environments, which we translated into RLang plans, policies, and effects. In HardMaze, this language advice made it possible to solve the task, as the vanilla Dyna-Q agent did no better than 352 random. For each experiment, 10 instances of each agent were run to generate a 95% confidence 353 interval on their cumulative reward over 50 episodes (LavaCrossing was run for 25 episodes only). 354 The number of timesteps per episode varied across environments. 355

- 356
- 357 358

359

4.2 LEVERAGING EXPERT ADVICE IN VIRTUALHOME

360 We ran additional experiments on custom environments based on the VirtualHome library, a platform 361 for simulating complex household activities. VirtualHome has an object-oriented state space, which 362 can be referenced natively in RLang. We engineered 2 tasks in a kitchen environment to assess our 363 language grounding pipeline: FoodSafety (see Figure 5), where the agent is tasked with putting a 364 pie into the fridge and salmon into the microwave, and CouchPotato, where the agent is tasked with bringing a remote control to a sofa and putting cereal into a kitchen cabinet, while avoiding picking 366 up toothpaste. In these environments, agents are given an RLang vocabulary file with groundings 367 for object-oriented perception and action abstractions such as the objects in the environment (e.g. 368 salmon_327, fridge_305), a short list of predicates (e.g. inside (), holding (), near ()), and a set of lifted skills (e.g. walk_to(), open, grab). These groundings have semantically 369 meaningful labels, which make them easy targets for grounding natural language. 370

Our experimental design here is identical to the Minigrid experiments: for each environment we collected multiple pieces of language advice from human experts and translated them into RLang programs via our two-stage pipeline. We then evaluated the performance of an RLang-Dyna-Q agent on our environments in comparison to a vanilla Dyna-Q agent. In all of our experiments, the RLang-informed agents significantly outperformed Dyna-Q. For each experiment, 10 instances of each agent were run to generate a 95% confidence interval on their cumulative reward over 50 and 70 episodes for FoodSafety and CouchPotato, respectively. The maximum number of timesteps per episode varied. The agent parameters are listed in the appendix.



Figure 3: **MidMazeLava Experiment.** Language advice given to the agent was grounded to RLang effects, plans, and policies. The full translated RLang program is available in the appendix. All RLang-Dyna-Q agents outperformed Dyna-Q.



Figure 4: **HardMaze Experiment.** Language advice given to the agent was grounded to RLang effects, plans, and policies. The full translated RLang program is available in the appendix. Vanilla Dyna-Q was not able to complete this task.



Figure 5: **FoodSafety Experiment.** Language advice given to the agent was grounded to RLang effects, plans, and policies. The full translated RLang program is in the appendix. All RLang-Dyna-Q agents outperformed Dyna-Q.



Figure 6: **CouchPotato Experiment.** Language advice given to the agent was grounded to RLang effects, plans, and policies. All the RLang agents outperformed Dyna-Q with the exception of the Effect-enabled agent. We note that bugs in the simulator non-deterministically prevent certain actions from executing, so the advice specified only applies part of the time, leading to decreased performance.



432 The impact of each kind of advice (e.g. plans, policies, transitions, and rewards) varied across tasks 433 in the VirtualHome and Minigrid experiments, with some environments benefiting primarily from 434 plan-centric advice and others benefiting most from policy advice. In virtually all cases, model-centric 435 advice-about transitions and rewards-was less valuable than other forms of advice. We suggest 436 that this discrepancy is due to how useful model-based advice is in comparison to explicit policy and planning advice. While policy and planning advice describe which actions to take in a given context, 437 model-based advice was often used to suggest which actions not to take, relying on the underlying 438 learning agent to find the best action. Furthermore, model-based advice was useful less of the time, 439 i.e. in fewer states. This is best illustrated by comparing the relative performance of effect-enabled 440 RLang-Dyna-Q agents with policy and plan-enabled agents in the MidMazeLava Experiment in 441 Figure 3 and the FoodSafety Experiment in Figure 5. The model-based advice in the first experiment 442 is to avoid lava, which there are many opportunities to walk into, resulting in the performance of the 443 effect-enabled agent closer to the plan and policy-enabled agents. By comparison, the model-based 444 advice in the second experiment is more niche, accounting only for a handful of transitions, and 445 the effect-enabled agent correspondingly performs closer to baseline Dyna-Q than to the plan and 446 policy-enabled agents.

447 448

449

4.3 GROUNDING SYMBOLS WITH A VISION-LANGUAGE MODEL

450 A crucial assumption made by our pipeline is that we are given semantically-meaningful labels for the 451 groundings we have, including labels for objects (e.g. salmon), skills (e.g. qo_to(kitchen)), and truth-valued predicates (e.g. is_open (fridge)). Assigning relevant labels for these ground-452 ings enables a relatively simple translation from natural language into RLang. In a real-world setting, 453 we imagine that the labels for these groundings can be generated in two ways: 1) prescriptively, in the case of skill engineering by humans, and 2) via a pre-trained foundation model for identifying 455 predicates and objects in the environment. Implementing a full symbol-grounding system³ is outside 456 the scope of this work, however, we performed an additional demonstration showing how the labels 457 for object groundings could be easily extracted from images of the VirtualHome environments using 458 a vision-language model. 459

In addition to performing the initial semantic labeling of objects, we demonstrated how incorporating 460 a vision-language model in-the-loop could expand the variety of advice our system is capable of 461 grounding. By asking a VLM to disambiguate referents using images of entities in the environment, 462 we are able to successfully ground entities for semantically-ambiguous advice. Given 11 images 463 of the entities in the VirtualHome environment, we asked GPT-40 to ground the referents of 17 464 ambiguous commands that would require visual and operational knowledge of the entities in the 465 scene. For example, we can ground noun phrases like "the white box you might put food in" to a 466 white microwave in the scene or "the tall red box" to a tall red refrigerator (see Figure 7). These 467 additional experiments-the automatic semantic labeling of entities in the environment and the in-theloop semantic grounding of ambiguous referents in advice-ameliorate the need for expert-crafted, 468 semantically perfect RLang grounding files. 469

- 471 4.4 EVALUATING TRANSLATION AND RLANG EFFICACY
- 472

470

To assess RLang's ability to capture the breadth of general language advice, we ran a small user 473 study. We asked 10 undergraduate students to solve the LockedRoom MiniGrid task (pictured in 474 Figure 8) and then asked them to describe in one or two sentences any advice they would give to 475 an agent completing the task for the first time. We collected their responses and ran them through 476 our translation pipeline to arrive at the RL ang groundings in Table 3 of the Appendix. Of 10 pieces 477 of advice collected, 9 were translated into valid RLang programs, while 1 referenced groundings 478 that did not exist (e.g. second_left_door). We used the remaining valid RLang programs to 479 inform 9 separate RLang-Dyna-Q agents that we compared against a baseline Dyna-Q agent given no 480 advice. With a few exceptions, providing advice either did not meaningfully impact performance 481 over the baseline or led to dramatic improvements in performance (see Table 2). In the cases where 482 advice did not impact performance, it was translated into a parsable RLang program that referenced 483 groundings that were not in the RLang vocabulary file (e.g. "the second left door" was translated to second_left_door, but a proper reference would be yellow_door). We address this symbol-484

³Learning a mapping from symbolic labels to groundings is explored in Steels & Hild (2012).



Figure 7: By capturing images of unnamed objects in the environment specified only by numeric object ids we can prompt a VLM to provide us with semantic labels. The labels output by VLM can be 502 provided to the translation pipeline as semantic primitives (see right). Additionally, we demonstrate how the VLM can be re-prompted after an initial labeling to resolve semantic ambiguities that require 504 visual knowledge of the entities in question. After disambiguation, the advice can be re-written to 505 incorporate the true grounding labels (e.g. cereal_235 instead of "blue box"), and then processed by the remaining grounding pipeline. Additional grounding examples are reported in the appendix. 506

'red box' refers to fridae'

grounding failure in the VirtualHome environment by using a VLM to ground semantically-ambiguous 509 referents (see section 4.3). Failures also occurred when users specified plans whose pre-conditions 510 were not met at the start state of the environment and failed to execute (e.g. the last piece of advice suggests to go to the room with the red key, but the agent cannot visit the room without first opening the grey door). 513

514 515

516

511

512

499 500

501

507 508

DISCUSSION AND CONCLUSION 5

Natural language grounding (Steels & Hild, 2012) has critical implications for all of AI. Just as 517 RL is intended as a model of intelligent decision making, we propose that its core formalisms offer 518 a natural target for language grounding. If MDPs model human decision-making, and humans 519 invented language to share information that aids their decision-making, then the appropriate target 520 for language grounding should be an MDP, or a richer and perhaps more structured decision process 521 reflecting the complexity of human decision-making. One line of evidence for this claim is the direct 522 correspondence between parts of speech and elements of structured decision-processes (Rodriguez-523 Sanchez et al., 2020). For example, the object classes in Object Oriented MDPs (Diuk et al., 524 2008) naturally correspond to the concept of **common nouns** requiring **determiners** to single out 525 class instances, and the parameters in Parameterized Action MDPs Masson et al. (2016) naturally 526 correspond to **adverbs** for modifying the execution of discrete macro-actions (**verbs**).

527 More practically, knowledge expressed in natural language has immense potential to inform rein-528 forcement learning agents, and thereby alleviate the high sample complexity of having to learn *tabula* 529 rasa. We present a novel method for leveraging general natural language advice to expedite learning 530 in Markov Decision Processes by translating it into RLang, a formal language designed to specify 531 information about every element of an MDP and its solution. Our method can ground advice to 532 reward functions, transition functions, plans, and policies. We also introduce a modified Dyna-Q 533 agent capable of leveraging all of the types of information present in the partial MDP specification represented by RLang. Our findings show that our approach can leverage a wide variety of language 534 advice to accelerate learning. 535

536

REFERENCES 538

Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Keerthana Gopalakrishnan, Karol Hausman, Alexander Herzog, Daniel Ho, Jasmine Hsu,

540 541 542 543 544 545 546	Julian Ibarz, Brian Ichter, Alex Irpan, Eric Jang, Rosario M Jauregui Ruano, Kyle Jeffrey, Sally Jesmonth, Nikhil Jayant Joshi, Ryan C. Julian, Dmitry Kalashnikov, Yuheng Kuang, Kuang- Huei Lee, Sergey Levine, Yao Lu, Linda Luu, Carolina Parada, Peter Pastor, Jornell Quiambao, Kanishka Rao, Jarek Rettinghouse, Diego M Reyes, Pierre Sermanet, Nicolas Sievers, Clayton Tan, Alexander Toshev, Vincent Vanhoucke, F. Xia, Ted Xiao, Peng Xu, Sichun Xu, and Mengyuan Yan. Do As I Can, Not As I Say: Grounding Language in Robotic Affordances. In <i>Conference on Robot Learning</i> , 2022.
547 548 549 550	Jacob Andreas, Dan Klein, and Sergey Levine. Modular Multitask Reinforcement Learning with Policy Sketches. In <i>Proceedings of the 34th International Conference on Machine Learning-Volume</i> 70, pp. 166–175. JMLR. org, 2017.
551 552 553 554	Dzmitry Bahdanau, Kyunghyun Cho, and Yoshua Bengio. Neural Machine Translation by Jointly Learning to Align and Translate. In Yoshua Bengio and Yann LeCun (eds.), <i>3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings</i> , 2015.
555 556 557	Matthew Berg, Deniz Bayazit, Rebecca Mathew, Ariel Rotter-Aboyoun, Ellie Pavlick, and Ste- fanie Tellex. Grounding Language to Landmarks in Arbitrary Outdoor Environments. In <i>IEEE</i> <i>International Conference on Robotics and Automation (ICRA)</i> , 2020.
558 559 560	S. R. K. Branavan, David Silver, and Regina Barzilay. Learning to Win by Reading Manuals in a Monte-Carlo Framework. <i>J. Artif. Intell. Res.</i> , 43:661–704, 2012.
561 562 563	Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. <i>Advances in neural information processing systems</i> , 33:1877–1901, 2020.
564 565 566 567 568	Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. Sparks of Artificial General Intelligence: Early experiments with GPT-4, 2023.
569 570 571 572 573 574 575 576 577 578	Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Henrique Pondé de Oliveira Pinto, Jared Kaplan, Harri Edwards, Yuri Burda, Nicholas Joseph, Greg Brockman, Alex Ray, Raul Puri, Gretchen Krueger, Michael Petrov, Heidy Khlaaf, Girish Sastry, Pamela Mishkin, Brooke Chan, Scott Gray, Nick Ryder, Mikhail Pavlov, Alethea Power, Lukasz Kaiser, Mohammad Bavarian, Clemens Winter, Philippe Tillet, Felipe Petroski Such, Dave Cummings, Matthias Plappert, Fotios Chantzis, Elizabeth Barnes, Ariel Herbert-Voss, William Hebgen Guss, Alex Nichol, Alex Paino, Nikolas Tezak, Jie Tang, Igor Babuschkin, Suchir Balaji, Shantanu Jain, William Saunders, Christopher Hesse, Andrew N. Carr, Jan Leike, Joshua Achiam, Vedant Misra, Evan Morikawa, Alec Radford, Matthew Knight, Miles Brundage, Mira Murati, Katie Mayer, Peter Welinder, Bob McGrew, Dario Amodei, Sam McCandlish, Ilya Sutskever, and Wojciech Zaremba. Evaluating Large Language Models Trained on Code. <i>CoRR</i> , abs/2107.03374, 2021.
579 580 581 582 583	Maxime Chevalier-Boisvert, Dzmitry Bahdanau, Salem Lahlou, Lucas Willems, Chitwan Saharia, Thien Huu Nguyen, and Yoshua Bengio. BabyAI: A Platform to Study the Sample Efficiency of Grounded Language Learning. In 7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019. OpenReview.net, 2019.
584 585 586 587	Maxime Chevalier-Boisvert, Bolun Dai, Mark Towers, Rodrigo de Lazcano, Lucas Willems, Salem Lahlou, Suman Pal, Pablo Samuel Castro, and Jordan Terry. Minigrid & Miniworld: Modular & Customizable Reinforcement Learning Environments for Goal-Oriented Tasks. <i>CoRR</i> , abs/2306.13831, 2023.
588 589 590 591	Vanya Cohen, Jason Xinyu Liu, Raymond Mooney, Stefanie Tellex, and David Watkins. A Survey of Robotic Language Grounding: Tradeoffs Between Symbols and Embeddings. <i>arXiv preprint arXiv:2405.13245</i> , 2024.
592 593	Cédric Colas, Ahmed Akakzia, Pierre-Yves Oudeyer, Mohamed Chetouani, and Olivier Sigaud. Language-Conditioned Goal Generation: a New Approach to Language Grounding for RL. <i>CoRR</i> , abs/2006.07043, 2020.

- Carlos Diuk, Andre Cohen, and Michael L Littman. An Object-Oriented Representation for Efficient Reinforcement Learning. In *Proceedings of the 25th international conference on Machine learning*, pp. 240–247, 2008.
- Yuqing Du, Olivia Watkins, Zihan Wang, Cédric Colas, Trevor Darrell, Pieter Abbeel, Abhishek Gupta, and Jacob Andreas. Guiding Pretraining in Reinforcement Learning with Large Language Models. In Andreas Krause, Emma Brunskill, Kyunghyun Cho, Barbara Engelhardt, Sivan Sabato, and Jonathan Scarlett (eds.), *International Conference on Machine Learning, ICML 2023, 23-29 July 2023, Honolulu, Hawaii, USA*, volume 202 of *Proceedings of Machine Learning Research*, pp. 8657–8677. PMLR, 2023.
- Nancy Fulda, Daniel Ricks, Ben Murdoch, and David Wingate. What Can You Do with a Rock?
 Affordance Extraction via Word Embeddings. In Carles Sierra (ed.), *Proceedings of the Twenty- Sixth International Joint Conference on Artificial Intelligence, IJCAI 2017, Melbourne, Australia, August 19-25, 2017*, pp. 1039–1045. ijcai.org, 2017.
- M. Ghallab, A. Howe, C. Knoblock, D. Mcdermott, A. Ram, M. Veloso, D. Weld, and D. Wilkins.
 PDDL—The Planning Domain Definition Language, 1998.
- ⁶¹¹ Nakul Gopalan, Dilip Arumugam, Lawson Wong, and Stefanie Tellex. Sequence-to-Sequence
 ⁶¹² Language Grounding of Non-Markovian Task Specifications. In *Proceedings of Robotics: Science* ⁶¹³ and Systems, Pittsburgh, Pennsylvania, 2018.
- Prasoon Goyal, Scott Niekum, and Raymond J. Mooney. Using Natural Language for Reward Shaping in Reinforcement Learning. In Sarit Kraus (ed.), *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI 2019, Macao, China, August 10-16, 2019*, pp. 2385–2391. ijcai.org, 2019.
- Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language Models as Zero-Shot
 Planners: Extracting Actionable Knowledge for Embodied Agents. In Kamalika Chaudhuri,
 Stefanie Jegelka, Le Song, Csaba Szepesvári, Gang Niu, and Sivan Sabato (eds.), *International Conference on Machine Learning, ICML 2022, 17-23 July 2022, Baltimore, Maryland, USA*,
 volume 162 of *Proceedings of Machine Learning Research*, pp. 9118–9147. PMLR, 2022.
- Kishor Jothimurugan, Rajeev Alur, and Osbert Bastani. A Composable Specification Language for Reinforcement Learning Tasks. In H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (eds.), *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019.
- Bartosz Kostka, Jaroslaw Kwiecieli, Jakub Kowalski, and Pawel Rychlikowski. Text-based adventures of the golovin AI agent. In 2017 IEEE Conference on Computational Intelligence and Games (CIG). IEEE, aug 2017.
- 632 Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muennighoff, Denis Kocetkov, Chenghao 633 Mou, Marc Marone, Christopher Akiki, Jia Li, Jenny Chim, Qian Liu, Evgenii Zheltonozhskii, 634 Terry Yue Zhuo, Thomas Wang, Olivier Dehaene, Mishig Davaadorj, Joel Lamy-Poirier, João 635 Monteiro, Oleh Shliazhko, Nicolas Gontier, Nicholas Meade, Armel Zebaze, Ming-Ho Yee, Logesh Kumar Umapathi, Jian Zhu, Benjamin Lipkin, Muhtasham Oblokulov, Zhiruo Wang, 636 Rudra Murthy V, Jason Stillerman, Siva Sankalp Patel, Dmitry Abulkhanov, Marco Zocca, Manan 637 Dey, Zhihan Zhang, Nour Moustafa-Fahmy, Urvashi Bhattacharyya, Wenhao Yu, Swayam Singh, 638 Sasha Luccioni, Paulo Villegas, Maxim Kunakov, Fedor Zhdanov, Manuel Romero, Tony Lee, 639 Nadav Timor, Jennifer Ding, Claire Schlesinger, Hailey Schoelkopf, Jan Ebert, Tri Dao, Mayank 640 Mishra, Alex Gu, Jennifer Robinson, Carolyn Jane Anderson, Brendan Dolan-Gavitt, Danish 641 Contractor, Siva Reddy, Daniel Fried, Dzmitry Bahdanau, Yacine Jernite, Carlos Muñoz Ferrandis, 642 Sean Hughes, Thomas Wolf, Arjun Guha, Leandro von Werra, and Harm de Vries. StarCoder: 643 may the source be with you! CoRR, abs/2305.06161, 2023. 644
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and Andy Zeng. Code as Policies: Language Model Programs for Embodied Control. In *IEEE International Conference on Robotics and Automation, ICRA 2023, London, UK, May 29 June 2, 2023*, pp. 9493–9500. IEEE, 2023.

- Michael L. Littman, Ufuk Topcu, Jie Fu, Charles Lee Isbell Jr., Min Wen, and James MacGlashan.
 Environment-Independent Task Specifications via GLTL. *CoRR*, abs/1704.04341, 2017.
 - Bo Liu, Yuqian Jiang, Xiaohan Zhang, Qiang Liu, Shiqi Zhang, Joydeep Biswas, and Peter Stone. LLM+P: Empowering Large Language Models with Optimal Planning Proficiency. *CoRR*, abs/2304.11477, 2023a.
- Jason Xinyu Liu, Ziyi Yang, Ifrah Idrees, Sam Liang, Benjamin Schornstein, Stefanie Tellex, and Ankit Shah. Lang2LTL: Translating Natural Language Commands to Temporal Robot Task Specification. *CoRR*, abs/2302.11649, 2023b.
- Jelena Luketina, Nantas Nardelli, Gregory Farquhar, Jakob Foerster, Jacob Andreas, Edward Grefenstette, Shimon Whiteson, and Tim Rocktäschel. A Survey of Reinforcement Learning Informed by Natural Language. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pp. 6309–6317. International Joint Conferences on Artificial Intelligence Organization, 7 2019.
- James MacGlashan, Monica Babes-Vroman, Marie desJardins, Michael L. Littman, Smaranda
 Muresan, Shawn Squire, Stefanie Tellex, Dilip Arumugam, and Lei Yang. Grounding English
 Commands to Reward Functions. In Lydia E. Kavraki, David Hsu, and Jonas Buchli (eds.), *Robotics: Science and Systems XI, Sapienza University of Rome, Rome, Italy, July 13-17, 2015*,
 2015. doi: 10.15607/RSS.2015.XI.018. URL http://www.roboticsproceedings.org/
 rss11/p18.html.
- Warwick Masson, Pravesh Ranchod, and George Dimitri Konidaris. Reinforcement Learning with Parameterized Actions. In Dale Schuurmans and Michael P. Wellman (eds.), *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, February 12-17, 2016, Phoenix, Arizona, USA*, pp. 1934–1940. AAAI Press, 2016.
- Shivam Miglani and Neil Yorke-Smith. Nltopddl: One-shot learning of pddl models from natural
 language process manuals. In *ICAPS'20 Workshop on Knowledge Engineering for Planning and Scheduling (KEPS'20)*. ICAPS, 2020.
- Suvir Mirchandani, Siddharth Karamcheti, and Dorsa Sadigh. ELLA: Exploration through Learned Language Abstraction. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), *Advances in Neural Information Processing Systems* 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 29529–29540, 2021.
 - Raymond J Mooney. Learning for semantic parsing. In *International Conference on Intelligent Text Processing and Computational Linguistics*, pp. 311–324. Springer, 2007.
- Karthik Narasimhan, Regina Barzilay, and Tommi S. Jaakkola. Grounding Language for Transfer in
 Deep Reinforcement Learning. *J. Artif. Intell. Res.*, 63:849–874, 2018.
- Kavier Puig, Kevin Ra, Marko Boben, Jiaman Li, Tingwu Wang, Sanja Fidler, and Antonio Torralba.
 Virtualhome: Simulating household activities via programs. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pp. 8494–8502, 2018.
 - R. Rodriguez-Sanchez, B.A. Spiegel, J. Wang, R. Patel, G.D. Konidaris, and S. Tellex. RLang: A Declarative Language for Describing Partial World Knowledge to Reinforcement Learning Agents. In *Proceedings of the Fortieth International Conference on Machine Learning*, July 2023.
- Rafael Rodriguez-Sanchez, Roma Patel, and George Konidaris. On the Relationship Between
 Structure in Natural Language and Models of Sequential Decision Processes. In *Language in Reinforcement Learning Workshop at ICML 2020*, 2020.
- 700

652

653 654

663

674

678

684

685

686

693

694

696

701 John Schulman, Filip Wolski, Prafulla Dhariwal, Alec Radford, and Oleg Klimov. Proximal Policy Optimization Algorithms. CoRR, abs/1707.06347, 2017.

702 703 704 705 706 707	Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B. Tenenbaum, Leslie Pack Kaelbling, and Michael Katz. Generalized Planning in PDDL Domains with Pretrained Large Language Models. In Michael J. Wooldridge, Jennifer G. Dy, and Sriraam Natarajan (eds.), <i>Thirty-Eighth AAAI Conference on Artificial Intelligence, AAAI 2024, Thirty-Sixth Conference on Innovative Applications of Artificial Intelligence, IAAI 2024, Fourteenth Symposium on Educational Advances in Artificial Intelligence, EAAI 2014, February 20-27, 2024, Vancouver, Canada</i> , pp. 20256–20264. AAAI Dress, 2004
708 709 710 711 712	 Fress, 2024. Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter Fox, Jesse Thomason, and Animesh Garg. ProgPrompt: Generating Situated Robot Task Plans using Large Language Models. In 2023 IEEE International Conference on Robotics and Automation (ICRA), pp. 11523–11530, 2023.
713 714 715 716 717	Chan Hee Song, Brian M. Sadler, Jiaman Wu, Wei-Lun Chao, Clayton Washington, and Yu Su. LLM-Planner: Few-Shot Grounded Planning for Embodied Agents with Large Language Models. In <i>IEEE/CVF International Conference on Computer Vision, ICCV 2023, Paris, France, October 1-6, 2023</i> , pp. 2986–2997. IEEE, 2023.
718 719	Luc Steels and Manfred Hild. <i>Language grounding in robots</i> . Springer Science & Business Media, 2012.
720 721 722	Richard S Sutton, Andrew G Barto, et al. <i>Introduction to Reinforcement Learning</i> , volume 135. MIT press Cambridge, 1998.
723 724 725 726 727	 Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is All you Need. In Isabelle Guyon, Ulrike von Luxburg, Samy Bengio, Hanna M. Wallach, Rob Fergus, S. V. N. Vishwanathan, and Roman Garnett (eds.), Advances in Neural Information Processing Systems 30: Annual Conference on Neural Information Processing Systems 2017, December 4-9, 2017, Long Beach, CA, USA, pp. 5998–6008, 2017.
728 729	Sai Vemprala, Rogerio Bonatti, Arthur Bucker, and Ashish Kapoor. ChatGPT for Robotics: Design Principles and Model Abilities. <i>IEEE Access</i> , 12:55682–55696, 2024.
730	Christopher John Cornish Hellaby Watkins. Learning from delayed rewards. 1989.
732 733 734 735 736	Olivia Watkins, Abhishek Gupta, Trevor Darrell, Pieter Abbeel, and Jacob Andreas. Teachable Reinforcement Learning via Advice Distillation. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan (eds.), Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual, pp. 6920–6933, 2021.
737 738 739 740	Jimmy Wu, Rika Antonova, Adam Kan, Marion Lepert, Andy Zeng, Shuran Song, Jeannette Bohg, Szymon Rusinkiewicz, and Thomas Funkhouser. TidyBot: Personalized Robot Assistance with Large Language Models. <i>Autonomous Robots</i> , 2023.
741 742	Yaqi Xie, Chen Yu, Tongyao Zhu, Jinbin Bai, Ze Gong, and Harold Soh. Translating Natural Language to Planning Goals with Large-Language Models. <i>CoRR</i> , abs/2302.05128, 2023.
743 744 745 746	Håkan L. S. Younes and Michael L. Littman. PPDDL 1.0: An Extension to PDDL for Expressing Planning Domains with Probabilistic Effects. In <i>PPDDL 1.0: An Extension to PDDL for Expressing Planning Domains with Probabilistic Effects</i> , 2004.
747 748 749 750 751 752 753 754	Wenhao Yu, Nimrod Gileadi, Chuyuan Fu, Sean Kirmani, Kuang-Huei Lee, Montserrat Gonzalez Arenas, Hao-Tien Lewis Chiang, Tom Erez, Leonard Hasenclever, Jan Humplik, Brian Ichter, Ted Xiao, Peng Xu, Andy Zeng, Tingnan Zhang, Nicolas Heess, Dorsa Sadigh, Jie Tan, Yuval Tassa, and Fei Xia. Language to Rewards for Robotic Skill Synthesis. In Jie Tan, Marc Toussaint, and Kourosh Darvish (eds.), <i>Conference on Robot Learning, CoRL 2023, 6-9 November 2023, Atlanta, GA, USA</i> , volume 229 of <i>Proceedings of Machine Learning Research</i> , pp. 374–404. PMLR, 2023.

A APPENDIX

A.1 USER STUDY

Table 2: User Study. We collected 10 pieces of advice from 10 undergraduate students for the LockedRoom environment. For each piece of advice, 5 agent instances were run for 25 episodes on the LockedRoom environment for 500 steps. The cumulative discounted reward for the 25 episodes is in the first column along with a 95% confidence interval. The average percent increase in cumulative discounted reward over the baseline is present in the second column. The second-to-last piece of advice did not ground to a valid RLang program, so no experiment was run.

Avg Cumula- tive Return	% improve- ment	Natural Language Advice
17.86 ± 2.36		No advice
22.01 ± 0.71	+23.24	"Remember to toggle to open doors."
16.79 ± 1.75	-5.99	"You don't need to carry keys to open the grey door."
17.22 ± 1.75	-3.60	"Identify the room with the red key, move to that room by opening th
		door. Pick up the key. Identify the room with the red door, procee
		there. Open the red door. Find the green square and go there to finis
		the game."
23.55 ± 0.69	+31.86	"Move to the grey door, open it and enter the room until you get to the
		red key, pick it up. Exit the room and move towards the red door, ope
		it and get into that room. Move to the green block and enter it."
23.93 ± 0.37	+33.99	"Go to the grey door. open the grey door. go to the red key. pick up the
		red key. go to the red door. open the red door. go to the green square
24.07 ± 0.22	+34.77	"Pick up the red key after opening the grey door. Then walk to the red
	1 (1)	door, open it, and go to the goal."
17.57 ± 0.78	-1.63	"You cannot open the red door without a red key."
17.77 ± 0.54	-0.50	"Walking towards the red door is not very useful if it is closed."
_		"Go down until the second door on the left and pick up the key. In
		exit the room and go down until the next door on the left and use it
10.95 1.09	10.70	open the door and get to the green box.
18.30 ± 1.93	+2.70	Go to the foom that has the fed key, pick it up, and then go to the foo
		with a red door. Enter the room, and go to the green goal object.

	_	_	_	_	_	_	_
				_			_
							_
	-						_
							_
							_
							_
				_	_	_	_
							_
							_
							_
							_
				_	-	-	_
							- 1
							- 1
							_
الكال الديد							
و و و از و و و							- 1

Figure 8: The initial state of the LockedRoom environment.



Figure 9: **MultiRoom Experiment.** The agent was given the following advice: "First go to the blue door, then the green door, then the grey door, then the purple door." The initial state of MultiRoom is pictured on the left, reward curves are center, and the translated advice is on the right.



A.2.1 [MINIGRID] PROMPT USED FOR STAGE 1 OF THE TRANSLATION PIPELINE. GIVEN A NEW PIECE OF ADVICE, WE PROMPT THE LLM TO CLASSIFY IT AS AN EFFECT, PLAN, OR POLICY.

RLang is a formal language for specifying information about every element of a Markov
Decision Process (S,A,R,T). Each RLang object refers to one or more elements of an MDP.
Here is a description of three important RLang groundings:

- Policy: a direct function from states to actions, best used for more general commands.
- Effect: a prediction about the state of the world or the reward function. Plan: a sequence of specific steps to take.
- Your task is to decide which RLang grounding most naturally corresponds to a given piece of advice:
- Advice = "Don't touch any mice unless you have gloves on."
- 854 Grounding: Effect
- 855 Advice = "Walking into lava will kill you."
- 856 Grounding: Effect Advice = "First get the money, then go to the green square."
- 857 Grounding: Plan

819

820

836

837

847

848

849

850 851

852

- Advice = "Go through the door to the goal."
- 859 Grounding: Plan
- Advice = "If you have the key, go to the door, otherwise you need to get the key."
- 861 Grounding: Policy
- 862 Advice = "If there are any closed doors, open them."
- 863 Grounding: Policy

Π.2.2	[MINIORID] I ROMFT USED FOR STADE 2 OF THE FIFELINE TO TRANSLATE A FIECE O
	ADVICE INTO AN KL ANG PLAN.
You	Ir task is to translate natural language advice to RLang plan, which is a sequence of
spee	cific steps to take. For each instance, we provide a piece of advice in natural language, a
list	of allowed primitives, and you should complete the instance by filling the missing plan
fun	ction. Don't use any primitive outside the provided primitive list corresponding to each
inst	ance, e.g., if there is no 'green_door' in the primitive list you must not use 'green_door'
Ior	the plan function.
Adv	vice – "Open the door with the key and go through it to the goal"
Prir	nitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right'.
'for	ward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left',
'poi	nting_up', 'go_to', 'step_towards', 'vellow_kev', 'vellow_door', 'agent', 'goal', 'at',
'in_i	inventory']
Pla	n main.
E	<pre>xecute go_to(yellow_key)</pre>
E:	xecute pickup
E	xecute toggle
E:	<pre>xecute go_to(goal)</pre>
Adv	vice = "Get the key behind the red door to open the grey door. Then drop the key to the
left	".
Prin	nitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left'. 'right'.
'for	ward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left',
'poi	nting_up', 'go_to', 'step_towards', 'yellow_key', 'yellow_door', 'agent', 'goal', 'at',
'in_i	inventory']
Pla	n main:
E	xecute go_to(red_door)
E:	xecute toggle xecute go to(grey key)
E	xecute pickup
	xecute go to (grev door)
E	
E: E:	xecute left

A.2.2 [MINIGRID] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A PIECE OF

919

971

A.2.3 [MINIGRID] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A PIECE OF ADVICE INTO AN RLANG POLICY.

924 925 926 Your task is to translate natural language advice to RLang policy, which is a direct function 927 from states to actions. For each instance, we provide a piece of advice in natural language, a 928 list of allowed primitives, and you should complete the instance by filling the missing policy 929 function. Don't use any primitive outside the provided primitive list corresponding to each 930 instance, e.g., if there is no 'green_door' in the primitive list you must not use "green_door' 931 for the policy function. 932 Advice = "If the yellow door is open, go through it and walk to the goal. Otherwise open the 933 yellow door if you have the key.' 934 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 935 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 936 'pointing_up', 'go_to', 'step_towards', 'yellow_key', 'yellow_door', 'agent', 'goal', 'at', 'car-937 rying'] 938 **Policy** main: 939 if yellow_door.is_open: 940 Execute go_to(goal) elif carrying (yellow_key) and at (yellow_door) and not yellow_door.is_open: 941 Execute toggle 942 Advice = "If you don't have the key, go get it." 943 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 944 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 945 'pointing_up', 'go_to', 'step_towards', 'grey_key', 'red_door', 'grey_door', 'agent', 'pur-946 ple_ball', 'at', 'carrying'] 947 Policy main: 948 if at(grey_key): 949 Execute pickup elif not carrying(grey_key): 950 Execute go_to(grey_key) 951 Advice = "If you are carrying a ball and its corresponding box is closed, open the box if you 952 are at it, otherwise go to the box if you can reach it.' 953 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 954 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 955 'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-956 ple_ball', 'at', 'reachable', 'carrying'] 957 Policy main: 958 if carrying(green_ball) and not green_box.is_open: 959 if at(green box): Execute toggle 960 elif reachable(green box): 961 Execute go_to (green_box) 962 Advice = "Drop any balls for boxes you can't reach" 963 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 964 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 965 'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-966 ple_ball', 'at', 'reachable', 'carrying'] 967 Policy main: 968 if carrying(green_ball) and not reachable(green_box): Execute drop 969 if carrying(purple_ball) and not reachable(purple_box): 970 Execute drop

	Advise "if you have one have for a data that you connect much you should drag it?
	Advice = 11 you have any key for a door that you cannot reach, you should drop it Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right'
	'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing right', 'pointing down', 'pointing left',
l	'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-
	ple_ball', 'at', 'reachable', 'carrying']
	Policy main
	<pre>if carrying(green_key) and not reachable(green_door):</pre>
	Execute drop if carrying(nurnle key) and not reachable(nurnle door).
	Execute drop
	<pre>if carrying(red_key) and not reachable(red_door): Execute drop</pre>
	Advice = "Hey listen, you can open the door if you have the key and at the door when the
	door is closed"
	Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right',
	forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left',
	pointing_up, 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-
	pie_bail, at, leachable, carrying j
	Policy main: if corruing (purple key) and not purple deer is open and at (purple deer):
	Execute toggle

```
A.2.4
                 [MINIGRID] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A PIECE OF
1027
                 ADVICE INTO AN RLANG EFFECT.
1028
1029
1030
1031
1032
1033
1034
1035
1036
1038
1039
1040
1041
1042
1043
1044
            Your task is to translate natural language advice to RLang effect, which is a prediction about
1045
            the state of the world or the reward function. For each instance, we provide a piece of advice
1046
            in natural language, a list of allowed primitives, and you should complete the instance by
1047
            filling the missing effect function. Don't use any primitive outside the provided primitive list
1048
            corresponding to each instance, e.g., if there is no 'green_door' in the primitive list you must
1049
            not use 'green_door' for the effect function.
1050
1051
            Advice = "Don't go to the door without the key"
1052
            Primitives = ['yellow_door', 'goal', 'pickup', 'yellow_key', 'toggle', 'go_to', 'carrying', 'at']
1053
            Effect main:
1054
              if at(yellow_door) and not carrying(yellow_key):
                Reward -1
1055
1056
            Advice = "Don't walk into closed doors. If you're tired, don't go forward."
1057
            Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right',
1058
            'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left',
            'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-
1059
            ple_ball', 'at', 'reachable', 'carrying']
            Effect main:
              if at(yellow_door) and yellow_door.is_closed and A == forward:
1062
                Reward -1
1063
                s′
                   -> S
              elif tired() and A == forward:
1064
                Reward -1
1065
            Advice = "Walking into balls is pointless. You will die if you walk into keys. Trying to open
1066
            a box when you aren't near it will do nothing."
1067
            Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right',
1068
            'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left',
1069
            'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-
1070
            ple_ball', 'at', 'reachable', 'carrying']
1071
1072
            Effect main:
1073
              if at (Ball) and A == forward:
                Reward 0
1074
                S' -> S
1075
              elif at(Key) and A == forward:
                Reward -1
1076
                S' -> S*0
1077
              elif at (Box) and A == toggle:
                Reward 0
                s'
                   -> S
1079
```

1080	A.3	PROMPTS USED FOR TRANSLATION PIPELINE FOR VIRTUALHOME EXPERIMENTS
1081		
1082		
1083		
1084		
1085		
1086		
1087		
1088		
1089		
1090		
1091		
1092		
1093		
1094	A.3.1	[VIRTUALHOME] PROMPT USED FOR STAGE 1 OF THE TRANSLATION PIPELINE. GIVEN
1095		A NEW PIECE OF ADVICE, WE PROMPT THE LLM TO CLASSIFY IT AS AN EFFECT, PLAN,
1095		or Policy.
1097		
1090		
1100		
1100		
1101		
1102		
1103		
1104		
1100		
1107		
1107		
1100		
1110		
1111	RI	ang is a formal language for specifying information about every element of a Markov
1112	De	ecision Process (S.A.R.T). Each RLang object refers to one or more elements of an MDP.
1113	He	ere is a description of three important RLang groundings:
1114		
1115	Pc	licy: a direct function from states to actions, best used for more general com-
1116	m	ands.
1117	Ef	fect: a prediction about the state of the world or the reward function.
1118		an: a sequence of specific steps to take.
1119	Vo	ur teale is to decide which PI and grounding most neturally corresponds to a given piece of
1120	01 ad	vice.
1121	A	lyice = "Don't touch any mice unless you have gloves on."
1122	G	ounding: Effect
1123	A	lvice = "Walking into lava will kill you."
1124	Gı	rounding: Effect
1125	A	lvice = "First get the money, then go to the green square."
1126	G	ounding: Plan
1127	A	lvice = "Go through the door to the goal."
1128	G	ounding: Plan
1129		livice = "If you have the key, go to the door, otherwise you need to get the key."
1130		ounding: Policy
1131		ivice = 11 inere are any closed doors, open them."
1132		ounding. Foncy hvice = "Open any doors if they are closed."
1133		rounding: Policy
1100		

1134 A.3.2 [VIRTUALHOME] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A 1135 PIECE OF ADVICE INTO AN RLANG PLAN. 1136 1137 1138 1139 1140 1141 1142 1143 1144 1145 1146 1147 1148 1149 Your task is to translate natural language advice to RLang plan, which is a sequence of 1150 specific steps to take. For each instance, we provide a piece of advice in natural language, a 1151 list of allowed primitives, and you should complete the instance by filling the missing plan 1152 function. Don't use any primitive outside the provided primitive list corresponding to each 1153 instance, e.g., if there is no 'green_door' in the primitive list you must not use 'green_door' 1154 for the plan function. 1155 1156 Advice = "Open the door with the key and go through it to the goal" 1157 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 1158 'pointing_up', 'go_to', 'step_towards', 'yellow_key', 'yellow_door', 'agent', 'goal', 'at', 1159 'in_inventory'] 1160 1161 **Plan** main: Execute go_to(yellow_key) 1162 Execute pickup 1163 Execute go_to(yellow_door) Execute toggle 1164 Execute go_to(goal) 1165 Advice = "Get the key behind the red door to open the grey door. Then drop the key to the 1166 left." 1167 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 1168 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 1169 'pointing_up', 'go_to', 'step_towards', 'yellow_key', 'yellow_door', 'agent', 'goal', 'at', 1170 'in_inventory'] 1171 **Plan** main: 1172 Execute go to (red door) 1173 Execute toggle Execute go_to(grey_key) 1174 Execute pickup 1175 Execute go_to(grey_door) Execute toggle 1176 Execute left 1177 Execute drop 1178 Advice = "Get the key behind the red door to open the grey door." Primitives = ['Agent', 1179 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 'forward', 'walk_to', 1180 'open', 'close', 'putin', 'grab', 'inside', 'grey_key_11', 'red_door', 'grey_door_127', 'agent', 1181 'purple_ball', 'is_on_a', 'at', 'at_any', 'in_inventory'] 1182 Plan main: 1183 Execute walk to (red door) Execute open(red door) 1184 Execute walk_to(grey_key_11) 1185 Execute grab(grey_key_11) Execute walk_to(grey_door_127) 1186 Execute open(grey_door_127) 1187

1188 A.3.3 [VIRTUALHOME] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A 1189 PIECE OF ADVICE INTO AN RLANG POLICY. 1190 1191 1192 1193 1194 1195 1196 1197 1198 1199 1200 1201 1202 1203 1204 1205 1206 1207 Your task is to translate natural language advice to RLang policy, which is a direct function from states to actions. For each instance, we provide a piece of advice in natural language, a 1208 list of allowed primitives, and you should complete the instance by filling the missing policy 1209 function. Don't use any primitive outside the provided primitive list corresponding to each 1210 instance, e.g., if there is no 'green_door' in the primitive list you must not use "green_door' 1211 for the policy function. 1212 1213 Advice = "If the yellow door is open, go through it and walk to the goal. Otherwise open the 1214 yellow door if you have the key.' 1215 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 1216 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 'pointing_up', 'go_to', 'step_towards', 'yellow_key', 'yellow_door', 'agent', 'goal', 'at', 'car-1217 rying'] 1218 1219 Policy main: if yellow_door.is_open: 1220 Execute go to (goal) elif carrying (yellow_key) and at (yellow_door) and not yellow_door.is_open: Execute toggle 1222 1223 Advice = "If you don't have the key, go get it" 1224 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 1225 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 1226 'pointing_up', 'go_to', 'step_towards', 'grey_key_11', 'red_door', 'grey_door', 'agent', 'purple_ball', 'is_on_a', 'at', 'at_any', 'in_inventory'] 1227 1228 Policy main: if at(grey_key_11): 1229 Execute pickup 1230 elif not carrying(grey_key_11): Execute go_to(grey_key_11) 1231 1232 Advice = "If you're at the fridge, close it." 1233 Primitives = ['Toothpaste', 'Bedroom', 'Character', 'Cereal', 'Bathroom', 'Sofa', 'Cabinet', 'Salmon', 'Pie', 'Kitchentable', 'Remotecontrol', 'Fridge', 'Microwave', 'Kitchen', 'Bookshelf', 'Livingroom', 'walk_to', 'open', 'close', 'putin', 'puton', 'grab', 'drop', 'can_drop', 'is_drop', 'inside', 'inside_something', 'on', 'at', 'is_closed', 'is_open', 'holding', 'near', 'char-1236 acter_1', 'kitchen_205', 'bookshelf_249', 'fridge_305', 'oven_133', 'pie_319', 'chicken_127', 1237 'cabinet_19'] 1238 1239 Policy main: if at (fridge 305): 1240 Execute close (fridge 305) 1241

1242 1243	A.3.4 [VIRTUALHOME] PROMPT USED FOR STAGE 2 OF THE PIPELINE TO TRANSLATE A PIECE OF ADVICE INTO AN RLANG EFFECT.
1244	
1246	
247	
248	
249	
250	
1251	
1252	
1253	
1254	
1255	
1256	
1237	
1259	Your task is to translate natural language advice to RLang effect, which is a prediction about
1260	the state of the world or the reward function. For each instance, we provide a piece of advice
261	in natural language, a list of allowed primitives, and you should complete the instance by
262	ning the missing effect function. Don't use any primitive outside the provided primitive list corresponding to each instance, e.g. if there is no 'green door' in the primitive list you must
1263	not use 'green door' for the effect function.
1264	
1265	Advice = "Don't go to the door without the key"
1266	Primitives = ['yellow_door', 'goal', 'pickup', 'yellow_key', 'toggle', 'go_to', 'carrying', 'at']
1267	Effect main:
1268	<pre>if at(yellow_door) and not carrying(yellow_key): Reward -1</pre>
1269	
1270	Advice = "Don't walk into closed doors, since it takes no effect" Primitives = ['A gent', 'Well', 'GoelTile', 'Leve', 'Key', 'Deet', 'Perl', 'Perl', 'right', 'right',
1271	'forward' 'nickun' 'dron' 'togele' 'done' 'nointing right' 'nointing down' 'nointing left'
1272	'pointing_up', 'go_to', 'step_towards', 'agent', 'goal', 'is_on_a', 'at', 'at_any', 'in_inventory']
1274	Effect main:
275	<pre>if at(yellow_door) and not yellow_door.is_open and A == forward:</pre>
276	S' -> S
277	Advice – "Welking to a broken object won't do envthing. You can't grab the bell if it's inside
278	something."
279	Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'is_broken',
1280	'left', 'right', 'forward', 'grab', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down',
281	'pointing_left', 'pointing_up', 'inside_something', 'go_to', 'step_towards', 'agent', 'goal',
1282	'is_on_a', 'at', 'at_any', 'in_inventory', 'gate_12', 'door_16', 'ball_121']
1203	<pre>Effect main: if A == walk to(gate 12) and is broken(gate 12).</pre>
1285	S' -> S
1286	<pre>if A == walk_to(door_16) and is_broken(door_16): S' -> S</pre>
1287	<pre>if A == grab(ball_121) and inside_something(ball_121):</pre>
1288	
1289	Advice = "Don't go to the purple ball"
1290	Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 'forward', 'walk to', 'open', 'close', 'putip', 'graph', 'inside', 'holding', 'graph key 11'
1291	'red door', 'grey door 127', 'agent' 'nurnle hall' 'is on a' 'at' 'at any' 'in inventory'
1292	Effort main.
1293	<pre>if A == walk_to(purple_ball):</pre>
1294	Reward -1
1293	

1296 Advice = "If you put the pie into the microwave and the chicken into the oven, and make sure 1297 that they are both on, you will get reward and the episode will end." 1298 Primitives = ['Toothpaste', 'Bedroom', 'Character', 'Cereal', 'Bathroom', 'Sofa', 'Cabinet', 1299 'Salmon', 'Pie', 'Kitchentable', 'Remotecontrol', 'Fridge', 'Microwave', 'Kitchen', 'Book-1300 shelf', 'Livingroom', 'walk_to', 'open', 'turn_on', 'close', 'putin', 'puton', 'grab', 'drop', 1301 'can_drop', 'is_drop', 'inside', 'inside_something', 'on', 'at', 'is_closed', 'is_open', 'holding', 1302 'near', 'character_1', 'kitchen_205', 'bookshelf_249', 'fridge_305', 'oven_133', 'pie_319', 1303 'chicken_127', 'microwave_19'] 1304 Effect main: 1305 if inside (pie_319, microwave_19) and inside (chicken_127, oven_133): if is_closed(microwave_19) and at(oven_133) and A == turn_on(oven_133): 1306 Reward 5 S' -> S elif is_closed(oven_133) and at(microwave_19) and A == turn_on(microwave_19): Reward 5 1309 S' -> S 1310 Advice = "If you're not trying to pick up the fridge, you will be penalized" Primitives = 1311 ['Sofa', 'Kitchentable', 'Bathroom', 'Salmon', 'Kitchen', 'Bookshelf', 'Cereal', 'Cabinet', 1312 'Livingroom', 'Fridge', 'Bedroom', 'Character', 'Toothpaste', 'Pie', 'Microwave', 'Remote-1313 control', 'walk_to', 'open', 'close', 'putin', 'puton', 'grab', 'drop', 'can_drop', 'is_drop', 1314 'inside', 'inside_something', 'on', 'at', 'fridge_305', 'is_pickup', 'is_closed', 'is_open', 'hold-1315 ing', 'near', 'character_1', 'bathroom_11', 'toothpaste_62', 'bedroom_73', 'kitchen_205', 1316 'kitchentable_231', 'bookshelf_249', 'fridge_305', 'microwave_313', 'pie_319', 'salmon_327', 1317 'cereal_334', 'livingroom_335', 'sofa_368', 'cabinet_415', 'remotecontrol_452'] 1318 Effect main: 1319 if fridge_305(fridge_305) and not is_pickup(A): Reward -1 1320 1321 Advice = "if you have any key for a door that you cannot reach, you should drop it" 1322 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 1323 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-1324 ple_ball', 'at', 'reachable', 'carrying'] 1325 1326 Policy main: if carrying(green_key) and not reachable(green_door): Execute drop 1328 if carrying(purple_key) and not reachable(purple_door): Execute drop 1329 if carrying(red_key) and not reachable(red_door): 1330 Execute drop 1331 Advice = "Hey listen, you can open the door if you have the key and at the door when the 1332 door is closed" 1333 Primitives = ['Agent', 'Wall', 'GoalTile', 'Lava', 'Key', 'Door', 'Box', 'Ball', 'left', 'right', 1334 'forward', 'pickup', 'drop', 'toggle', 'done', 'pointing_right', 'pointing_down', 'pointing_left', 'pointing_up', 'go_to', 'step_towards', 'green_ball', 'green_box', 'purple_box', 'agent', 'pur-1335 1336 ple_ball', 'at', 'reachable', 'carrying'] 1337 Policy main: 1338 if carrying(purple_key) and not purple_door.is_open and at(purple_door): Execute toggle 1339 1340 1341

A.4 USER STUDY - TRANSLATED ADVICE FOR MINIGRID EXPERIMENTS

1343 1344 1345

- 1346
- 1347

1348

1350 1351 1352 1353 1354 1355 1356 1357 Table 3: Advice from the user study translated to RLang. 1358 1359 **RLang Translation** Language Advice 1360 Policy main: "Remember to toggle to open doors." 1361 if at (yellow_door) and not yellow_door.is_open 1362 1363 Execute toggle elif at(red_door) and not red_door.is_open: 1364 Execute toggle elif at (purple_door) and not purple_door. 1365 is_open: 1366 Execute toggle 1367 elif at(blue_door) and not blue_door.is_open: Execute toggle 1368 elif at (green_door) and not green_door.is_open 1369 Execute toggle 1370 elif at (grey_door) and not grey_door.is_open: Execute toggle 1371 1372 "You don't need to carry keys to open the grey Effect main: 1373 if at(grey_door) and carrying(red_key): door." $S' \rightarrow S$ 1374 Reward -1 1375 elif at(grey_door) and carrying_something(): S' -> S 1376 Reward -1 1377 1378 "Identify the room with the red key, move Plan main: Execute go_to(red_key) to that room by opening the door. Pick up 1379 **Execute** pickup the key. Identify the room with the red door, 1380 Execute go_to(red_door) proceed there. Open the red door. Find the **Execute** toggle 1381 green square and go there to finish the game." Execute go_to(goal) 1382 "Move to the grey door, open it and enter the Plan main: Execute go_to(grey_door) room until you get to the red key, pick it up. 1384 **Execute** toggle Exit the room and move towards the red door, Execute go_to(red_key) 1385 open it and get into that room. Move to the Execute pickup 1386 green block and enter it." Execute go_to(grey_door) Execute toggle 1387 Execute go_to(red_door) 1388 Execute toggle Execute go_to(goal) 1389 1390 "Go to the grey door. open the grey door. go Plan main: 1391 Execute go_to(grey_door) to the red key. pick up the red key. go to the **Execute** toggle 1392 red door. open the red door. go to the green Execute go_to(red_key) square." 1393 **Execute** pickup Execute go_to(red_door) 1394 Execute toggle 1395 Execute go_to(goal) 1396 1397 1398 1399 1400 1401 1402

Table 4: Advice from the user study translated to RLang (continued).

Language Advice	RLang Translation			
"Pick up the red key after opening the grey	Plan main:			
door. Then walk to the red door, open it, and	Execute go_to(grey_door) Execute toggle			
go to the goal."	Execute go_to(red_key)			
	Execute pickup Execute go to(red door)			
	Execute toggle			
	Execute go_to(goal)			
"You cannot open the red door without a red	<pre>Effect main: if at(red_door) and not carrying(red_key): S' -> S Reward -1</pre>			
key."				
-				
"Walking towards the red door is not very	<pre>Effect main: if at (red door) and not (red door is open) and</pre>			
useful if it is closed."	<pre>if at (red_door) and not (red_door.is_open) and A == forward: S' -> S Reward -1</pre>			
"Go down until the second door on the left	Plan main:			
and pick up the key. Then exit the room and	Execute go_to(second_left_door) Execute pickup			
go down until the next door on the left and	Execute go_to(exit)			
box "	Execute go_to(next_left_door) Execute toggle			
oox.	Execute go_to(green_box)			
"C	Plan main.			
it up, and then go to the room with a red	Execute go_to(red_key)			
door. Enter the room, and go to the green	Execute pickup			
goal object."	Execute go_to(red_door) Execute toggle			
	Execute go_to(goal)			

1458 A.5 MIDMAZELAVA - TRANSLATED ADVICE 1459

1460 1461

1511

1462 1463 1464 1465 1466 1467 1468 1469 1470 1471 1472 1473 Advice: "Pick up the blue ball and drop it to your right. Then pick up the green key and 1474 unlock the green door. Then drop the key to your right. Some general advice: If you are 1475 carrying a key and its corresponding door is closed, open the door if you are at it, otherwise 1476 go to the door if you can reach it. Otherwise, drop any keys for doors you can't reach. If you 1477 can reach the goal, go to it. Walking into lava will kill you. If you're not at a door, toggling will do nothing. Trying to pick something up while you're carrying something is pointless. 1478 Walking into walls will do nothing." 1479 1480 **Plan** main: Execute go_to(blue_ball) 1481 Execute pickup 1482 Execute right Execute drop 1483 Execute go_to(green_key) 1484 Execute pickup Execute go to (green door) 1485 Execute toggle 1486 Execute right Execute drop 1487 1488 Policy main: if carrying(green_key) and not green_door.is_open: 1489 if at(green_door): 1490 Execute toggle elif reachable(green_door): 1491 Execute go_to(green_door) 1492 elif carrying(grey_key) and not grey_door.is_open: 1493 if at(grey_door): 1494 Execute toggle elif reachable(grey_door): 1495 Execute go_to(grey_door) 1496 elif reachable(goal): 1497 Execute go_to(goal) 1498 elif carrying(green_key) and not reachable(green_door): 1499 Execute drop 1500 elif carrying(grey_key) and not reachable(grey_door): 1501 Execute drop 1502 Effect main: 1503 if at(Lava) and A == forward: 1504 S' -> S*0 Reward -11505 if not at(Door) and A == toggle: 1506 $S' \rightarrow S$ Reward 0 1507 if carrying_something() and A == pickup: 1508 S' -> S Reward 0 1509 if at(Wall) and A == forward: 1510 S' -> S Reward 0

1512 A.6 HARDMAZELIGHT - TRANSLATED ADVICE

Advice: "Go and pick up the green ball, and drop it on your left, and then go pick up the blue
key, and go to the blue door and open it up and drop the key on your left, and then go pick up the green law, and go to the green door to open it and drop the key on your left, and then go pick up
nick up the purple ball and drop it on your right. Nothing will happen if you walk towards the
wall or try to open a purple door without the purple key if it is locked. The applies for the
vellow door and key as well as the red door and key. If you can reach the grey door and it is
closed but you have the key, open it if you are at it or otherwise go to it. The same applies to
the purple door, yellow door, and red door. Lastly, if you find the goal is reachable just go to
the goal directly."
Plan main:
<pre>Execute go_to(green_ball)</pre>
Execute pickup Execute left
Execute drop
Execute go_to(blue_key)
Execute go_to(blue_door)
Execute toggle
Execute drop
Execute drop Execute go_to(green_key)
Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door)
Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door) Execute toggle
Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door) Execute toggle Execute right Execute drop
Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball)
Execute drop Execute go_to(green_key) Execute pickup Execute toggle Execute toggle Execute drop Execute drop Execute go_to(purple_ball) Execute pickup Execute right
Execute drop Execute go_to(green_key) Execute pickup Execute dogle Execute toggle Execute right Execute drop Execute pickup Execute right Execute drop
Execute drop Execute go_to(green_key) Execute pickup Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute drop Execute drop Effect main:
<pre>Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at(Wall) and A == forward:</pre>
<pre>Execute drop Execute go_to(green_key) Execute pickup Execute go_to(green_door) Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute pickup Execute right Execute drop Effect main: if at(Wall) and A == forward: Reward 0 S' -> S</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at(Wall) and A == forward: Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at(Wall) and A == forward: Reward 0 S' -> S elif at(purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at (Wall) and A == forward: Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute toggle Execute right Execute drop Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at (Wall) and A == forward: Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at (yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yellow_key):</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute drop Execute right Execute drop Execute go_to(purple_ball) Execute jickup Execute right Execute drop Effect main: if at (Wall) and A == forward: Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at (yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yellow_key): Reward 0 S' -> S elif at (yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yellow_key): Reward 0 S' -> S</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute go_to(green_door) Execute go_to(purple_ball) Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at (Wall) and A == forward: Reward 0 S' -> S elif at (purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at (yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yurple_key): Reward 0 S' -> S elif at (yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yellow_key): Reward 0 S' -> S elif at (red_door) and red_door.is_locked and A == toggle and not carrying(red_key):</pre>
<pre>Execute drop Execute go_to(green_key) Execute go_to(green_door) Execute go_to(green_door) Execute drogle Execute right Execute go_to(purple_ball) Execute go_to(purple_ball) Execute pickup Execute right Execute drop Effect main: if at(Wall) and A == forward: Reward 0 S' -> S elif at(purple_door) and purple_door.is_locked and A == toggle and not carrying(purple_key): Reward 0 S' -> S elif at(yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yurple_key): Reward 0 S' -> S elif at(yellow_door) and yellow_door.is_locked and A == toggle and not carrying(yellow_key): Reward 0 S' -> S elif at(red_door) and red_door.is_locked and A == toggle and not carrying(red_key): Reward 0 S' -> S elif at(red_door) and red_door.is_locked and A == toggle and not carrying(red_key): Reward 0 S' -> S</pre>

Policy main:

else:

else:

if at(grey_door):

Execute toggle

if at(purple_door):

Execute toggle

Execute go_to (grey_door)

Execute go_to (purple_door)

```
1566
1567
1568
1569
1570
1571
1572
1573
1574
1575
1576
1578
1579
1580
1581
1582
1583
1584
1585
1586
1587
1588
```

1591 1592 1593 1594 1595

1596

1597

1598

1599

1590

if at(yellow_door):
 Execute toggle
else:
 Execute go_to(yellow_door)
elif reachable(red_door) and carrying(red_key) and red_door.is_locked:
 if at(red_door):
 Execute toggle
else:
 Execute go_to(red_door)
elif reachable(goal):
 Execute go_to(goal)

if reachable(grey_door) and carrying(grey_key) and grey_door.is_locked:

elif reachable (purple_door) and carrying (purple_key) and purple_door.is_locked:

elif reachable(yellow_door) and carrying(yellow_key) and yellow_door.is_locked:

A.7 FOODSAFETY - TRANSLATED ADVICE

Advice: "Go to fridge and open it, and then go find the pie and pick it up, walk back to the fridge and put the pie in the fridge. You have to close the fridge too", "If the salmon is in the microwave, and you are at the microwave and it's open, close it. Otherwise if you are holding salmon, do the following: open the microwave if you are near it but it's closed, put the salmon into the microwave if it's open and you're near it, else walk to the microwave.", "If the pie is in the fridge, and the salmon is in the microwave, then closing the fridge if the microwave is closed or closing the microwave if the fridge is closed will give you reward and end the episode."

```
Plan main:
              Execute walk_to(fridge_305)
              Execute open(fridge_305)
1602
              Execute walk_to(pie_319)
1603
              Execute grab(pie_319)
              Execute walk_to(fridge_305)
1604
              Execute putin(fridge_305)
1605
              Execute close (fridge_305)
1606
1607
            Policy main:
              if inside (salmon_327, microwave_313) and at (microwave_313) and is_open (microwave_313):
1608
                Execute close (microwave_313)
1609
              elif holding(salmon_327):
                if at (microwave_313) and is_closed(microwave_313):
1610
                  Execute open(microwave_313)
1611
                elif at (microwave_313) and is_open (microwave_313):
                  Execute putin (microwave_313)
1612
                else:
1613
                  Execute walk_to (microwave_313)
1614
            Effect main:
1615
              if inside (pie_319, fridge_305) and inside (salmon_327, microwave_313):
                if is_closed(fridge_305) and at(microwave_313) and A == close(microwave_313):
1616
                  Reward 5
1617
                  s′
                      -> S
                elif is_closed(microwave_313) and at(fridge_305) and A == close(fridge_305):
1618
                  Reward 5
1619
                  S' -> S
```

1620	A.8	CouchPotato - Translated Advice
1621		
1622		
1623		
1624		
1625		
1626		
1627		
1628		
1629		
1630		
1631		
1632		
1633		
1634		
1626		
1637		
1638		
1639		
1640		
1641		
1642		
1643		
1644		
1645		
1646		
1647		
1648		
1649		
1650		
1651		
1652		
1653		
1654		
1655	A	dvice: "If you're holding the toothpaste and can drop it, drop it.", "Go grab the remote
1656	cc	ntrol and put it on the sofa.", "If you're holding the toothpaste and not trying to drop it, you
1650	W	ill be penalized. Also, nothing will happen if you try to walk to the remote control, cereal,
1000	1 10	ounpaste, or salmon, if you try to walk to them and they are contained inside anything.
1009	Ef	<pre>fect main: if holding(toothpaste 62) and not is drop(A);</pre>
1661		Reward -1
1662		<pre>if inside_something(remotecontrol_452) and A == walk_to(remotecontrol_452): S' -> S</pre>
1663		<pre>if inside_something(cereal_334) and A == walk_to(cereal_334):</pre>
1664		<pre>S' -> S if inside_something(toothpaste_62) and A == walk_to(toothpaste_62):</pre>
1665	1	S' -> S
1666	1	S' -> S
1667	1	
1668	Po	blicy main:
1669	1	<pre>if holding(toothpaste_62) and can_drop(toothpaste_62): Execute drop(toothpaste_62)</pre>
1670	1	
1671	P1	an main: Execute walk_to(remotecontrol_452)
1672	1	Execute grab(remotecontrol_452)
1673		Execute puton (remotecontrol_452, sofa_368)

1674 A.9 GROUNDING SEMANTICALLY-AMBIGUOUS ADVICE

	Advice: Put the baked good into the white box.
	'the baked good' refers to Pie, and 'the white box' refers to Microwave.
	Advice: Put the baked dessert in the red food-box.
	'the baked dessert' refers to Pie, and 'the red food-box' refers to Fridge.
	Advice: Store the pastry in the tall red box.
	'the pastry' refers to Pie, and 'the tall red box' refers to Fridge.
	Advice: Open the appliance for heating.
	'the appliance for heating' refers to Microwave.
	Advice: Open the cooking device.
	'the cooking device' refers to Microwave.
	Advice: Grab the blue box and put it on the wooden white surface.
	In this context, 'the blue box' refers to Cereal, and 'the wooden white surface' refers
	Kitchentable.
	Advice: Place the seafood on the dining surface.
	the seafood refers to Salmon and the dining surface refers to Kitchentable.
	Advice: Store the breakfast grains on the black shelving.
	the breakfast grains refers to Cereal, and the black shelving refers to Bookshelf.
,	Advice: Set the channel changer on the seating furniture.
	Advice: Open the cooling appliance and grab the seafood
	Advice. Open the cooling appliance and grab the scarood.
	Advice: Put the small tube on the wooden hox
	'the small tube' refers to Toothnaste, and 'the wooden box' likely refers to Cabinet
	Advice: Put the food with the white packaging into the big red box
	'the food with the white packaging' refers to Cereal, and 'the big red box' refers to Fridge.
	Advice: Place the baked dessert on the white hard surface.
	'the baked dessert' refers to Pie, and 'the white hard surface' refers to Kitchentable.
	Advice: Place the seafood on the white wooden surface.
	'the seafood' refers to Salmon, and 'the white wooden surface' refers to Kitchentable.
	Advice: Put the orange food in the red box.
	'the orange food' refers to Salmon, and 'the red box' refers to Fridge.
	Advice: Open the red food-box.
	'the red food-box' refers to Fridge.
	Advice: Open the white box you might put food in.
	'the white box you might put food in' refers to Microwave.

B ADDITIONAL EXPERIMENT: GROUNDING COMMANDS TO RLANG PLANS -COMPARISON TO SAYCAN

1713 1714

1712

We compare our method to SayCan (Ahn et al., 2022), which uses the commonsense reasoning
capacity of LLMs to satisfy a natural language request by generating a simple plan consisting of a
series of pre-engineered high-level robot skills. Adopting the same

In this experiment we demonstrate that, in a simulated 3-dimensional physical environment, RLang can express the full range of natural language instructions necessary for a robot to complete various tasks. By grounding natural language instructions to RLang policies over this environment, we achieve performance on par with the results from the open-source tasks that the original SayCan paper evaluated on, showing that RLang can be easily substituted for the formal language that the SayCan authors developed for this specific task, allowing for generalization without sacrificing performance.

1724 Similar to the SayCan work, we assume that we are given a grounding tuple $\langle \Pi, S, A, \rangle$, and a set 1725 of skills Π , where each skill $\pi \in \Pi$ performs an action with the robot arm to manipulate a block or 1726 a bowl. We evaluate on the 8 unique tasks made available in the open-source version of SayCan, 1727 running each task across 10 different randomly selected initial states, using both the native SayCan 1829 language and RLang as the DSL for grounding natural language instructions to robot behavior.



Figure 10: Top: One configuration of the initial and completed states in the SayCan environment.
Bottom: the action sequence to execute on the instruction: "put all the blocks in the green bowl.",
from the robot arm's perspective

Table 5: Success rates of SayCan and RLang-based instruction grounding rate on each task, out of 10 random initial states.
 1762

Instruction	SayCan	NL2RLang
put all the blocks in different corners.	10	10
move the block to the bowl.	6	6
put any blocks on their matched colored bowls.	7	7
put all the blocks in the green bowl.	7	7
stack all the blocks.	8	8
make the highest block stack.	7	7
put the block in all the corners.	10	10
clockwise, move the block through all the corners.	10	10

Each task configuration that the original SayCan agent completes, is also completed by the RLang agent. While their behavior on failure cases occasionally varied, these were generally caused by errors in the vision model's processing of shadows in the simulated environment. These generally caused the textual scene description fed into GPT-3 to include a block where a bowl should be, and occasionally incorrect color labels, which often provided the text-only planner with a nonsensical task that was impossible to complete. Similarly, in cases where multiple action orders could satisfy the request, the RLang and SayCan pipelines occasionally diverged in the order of actions. Nonetheless, neither language grounding pipeline completed a task configuration that the other one did not.

1782 C EXPERIMENT PARAMETERS

In all experiments, for both Dyna-Q and RLang-Dyna-Q, we set the learning rate α to 0.1, the discount factor γ to 0.99, $\epsilon_1 = \epsilon_2 = 0.1$, except when there is policy or plan advice, uniform exploration, and 16 hallucinatory updates with the learned dynamics model. For the translation step, we use the gpt-3.5-turbo-instruct model with a temperature of 0.