# INNATECODER: LEARNING PROGRAMMATIC OPTIONS WITH FOUNDATION MODELS

Anonymous authors

004

006

008 009

010 011

012

013

014

015

016

017

018

019

021

024

025

026

027 028 029 Paper under double-blind review

#### ABSTRACT

Outside of transfer learning settings, reinforcement learning agents start their learning process from a clean slate. As a result, such agents have to go through a slow process to learn even the most obvious skills required to solve a problem. In this paper, we present INNATECODER, a system that leverages human knowledge encoded in foundation models to provide programmatic policies that encode "innate skills" in the form of temporally extended actions, or options. In contrast to existing approaches to learning options, INNATECODER learns them from the general human knowledge encoded in foundation models in a zero-shot setting, and not from the knowledge the agent gains by interacting with the environment. Then, INNATECODER searches for a programmatic policy by combining the programs encoding these options into larger and more complex programs. We hypothesized that INNATECODER's way of learning and using options could improve the sampling efficiency of current methods for learning programmatic policies. We evaluated our hypothesis in MicroRTS and Karel the Robot, two challenging domains. Empirical results support our hypothesis, since they show that INNATECODER is more sample efficient than versions of the system that do not use options or learn the options from experience. The policies INNATECODER learns are competitive and often outperform current state-of-the-art agents in both domains.

#### 1 INTRODUCTION

031 Outside of transfer learning settings, deep reinforcement learning (DRL) agents begin their learning 032 process with randomly initialized neural networks. As a result, DRL agents must learn from scratch 033 even the most basic skills required to solve a problem. In this paper, we harness the general human 034 knowledge encoded in foundation models to endow agents with helpful skills before they even start interacting with the environment. This is achieved by using programmatic representations of policies (Trivedi et al., 2021)-programs written in a domain-specific language encoding policies-037 and the foundation models' ability to write computer programs. Depending on the language used, 038 programmatic policies were shown to generalize better to unseen scenarios (Inala et al., 2020) and to be human-interpretable (Verma et al., 2018; Bastani et al., 2018). In addition to these advantages and loosely inspired by the innate abilities of animals (Tinbergen, 1951), we show that programmatic 040 representations of policies allow us to harness helpful "innate skills" from foundation models. 041

Given a natural-language description of the problem that the agent needs to learn to solve, our system, which we call INNATECODER, queries a foundation model for programs that encode policies to solve
the problem. Although the programs the model generates are unlikely to encode policies that solve
the problem, we hypothesize that the set of sub-programs we obtain from these programs can encode
helpful temporally extended actions, or options (Sutton et al., 1999). We consider options as functions
the agent can call and that will tell it how to act for a number of steps (Precup et al., 1998).

Options can ease the agent's learning process in different ways. For example, options can allow the agent to better explore the problem space (Machado et al., 2017; Bellemare et al., 2020) or can transfer knowledge between different tasks (Konidaris & Barto, 2007). In this paper, we present a novel way of learning options with foundation models. We also present a novel way of using the learned options, which is inspired by recent work on learning semantic spaces of programming languages (Moraes & Lelis, 2024). We leverage the compositional nature of the programmatic options we harness from a foundation model to learn the underlying semantic space of the programming



Figure 1: Left: The context-free grammar specifying a simplified version of the domain-specific language for Karel the Robot, a benchmark we use in our experiments. Right: the abstract syn-062 tax tree for if markersPresent then pickMarker. In the tree, MP and PM stand for markersPresent and pickMarker, respectively. Karel is a robot acting on a grid, where it needs to accomplish tasks such as collecting and placing markers on different locations of the grid. 064 In this program, Karel will pick up a marker if one is present in its current location on the grid.

065 066 067

060

061

063

language that defines the agent's hypothesis space. In the semantic space, neighbor programs encode 068 similar but different agent behavior, which is a desirable property when inducing spaces conducive 069 to algorithms searching for programmatic policies (Trivedi et al., 2021). The semantic space is approximated by ensuring that neighboring policies differ in term of one sub-policy from the set of 071 options. Instead of searching in the space of programs induced by the syntax of the language (Koza, 072 1992), INNATECODER searches in the space of semantically different programs induced by options. 073 In contrast with previous methods that can benefit from up to hundreds of options (Eysenbach et al., 074 2019), INNATECODER's use of programmatic options allows it to benefit from thousands of them.

075 INNATECODER's approach to harnessing options from foundation models contrasts with previous 076 approaches to automatically learning them, e.g., (Tessler et al., 2017; Bacon et al., 2017; Igl et al., 077 2020; Klissarov & Machado, 2023). This is because options are harnessed from the general knowledge 078 encoded in a foundation model, as opposed to the knowledge the agent gains by interacting with the 079 environment. This zero-shot approach to learning options is enabled by the use of a domain-specific 080 language to bridge the gap between the high-level knowledge encoded in foundation models and the 081 low-level knowledge required at the sensorimotor control level of the agent (Klissarov et al., 2024). For example, foundation models trained on Internet data likely encode the knowledge that, to win a 083 match of a real-time strategy game, the player must collect resources and build structures, which will allow for the training of the units needed to win the game. However, the model cannot issue low-level 084 actions in real time to control dozens of units to accomplish this plan. INNATECODER bridges this 085 gap by distilling the knowledge of the model into options that can be executed in real time.

087 We evaluated our hypothesis that foundation models can generate helpful programmatic options in 088 the domains of MicroRTS, a challenging real-time strategy game (Ontañón, 2017), and Karel the Robot (Pattis, 1994), which has been used as a benchmark for program synthesis and reinforcement 089 learning algorithms (Bunel et al., 2018; Chen et al., 2018; Shin et al., 2018; Trivedi et al., 2021). The 090 results in both domains support our hypothesis, since INNATECODER was more sample-efficient than 091 versions of the system that do not use options or learn options from experience. We also show that 092 the policies INNATECODER learns are competitive and often outperform the current state-of-the-art algorithms. INNATECODER is inexpensive because it uses the foundation model a small number of 094 times as a pre-processing step, making it an accessible system to smaller labs and companies. 095

096

098

#### **PROBLEM DEFINITION** 2

099 We consider sequential decision-making problems that can be formulated as Markov decision pro-100 cesses (MDPs)  $(S, A, p, r, \mu, \gamma)$ . Here, S represents the set of states and A is the set of actions. 101 The function  $p(s_{t+1}|s_t, a_t)$  is the transition model, which gives the probability of reaching state 102  $s_{t+1}$  given that the agent is in  $s_t$  and takes action  $a_t$  at time step t. The agent observes a reward 103 value of  $R_{t+1}$  when transitioning from  $s_t$  to  $s_{t+1}$ . The reward value the agent observes is returned by the function r.  $\mu$  is the distribution of the initial states of the MDP; states sampled from  $\mu$  are 104 denoted  $s_0$ .  $\gamma$  in [0,1] is the discount factor. A policy  $\pi$  is a function that receives a state s and 105 returns a probability distribution over actions available at s. The goal is to learn a policy  $\pi$  that 106 maximizes the expected sum of discounted rewards for  $\pi$  starting in an  $s_0$ :  $\mathbb{E}_{\pi,p,\mu}[\sum_{k=0}^{\infty} \gamma^k R_{k+1}]$ . 107  $V^{\pi}(s) = \mathbb{E}_{p,\pi}[\sum_{k=0}^{\infty} \gamma^k R_{k+1} | s_0 = s]$  is the value function, which measures the expected return

108 Learning Options Semantic Space 109 Programmatic Set of Options  $n' \in \mathcal{N}_k^m(n)$ 110 Policies if b1 then: if b1 then: 111 c1 c1 c2 112 MDP b1 then: c2 Foundation c1113 DSLModel b1 if then: c2 c1 C(n')c1 114 c2 c1 c1 115 c2 116 117 Local Search 118 Hill Climbing Step Initialization Stopping Condition 119 120  $\arg \max \hat{V}^{n'}(s_0)$  $n' {\in} \mathcal{N}_k^m(n)$ 121  $@\sim \mathcal{I}$ Local Optimum? Restart? no ⊾n' 122  $rg \max \hat{V}^{n'}(s_0)$ 123  $n' \in \mathcal{N}_k^x(n)$ 125 126

Figure 2: Schematic view of INNATECODER, with three parts. The "Learning Options" component harnesses options from a foundation model from a natural language description of the MDP and a Backus-Naur form description of the DSL. The model generates a set of programmatic policies that are broken down into a set of options (Section 3.1). The "Semantic Space" component uses the options to approximate the semantic space of the DSL (Section 3.2). The "Local Search" component searches in a mixture of the syntax and semantic spaces for a programmatic policy  $n^*$  (Section 3.3).

134

127

when the agent follows the policy  $\pi$  starting from the state s. In this work, we approximate the value function of a policy  $\pi$  and state s with Monte Carlo roll-outs and denote the approximation as  $\hat{V}^{\pi}(s)$ .

138 We consider programmatic representations of policies, which are policies written in a domain-specific language (DSL). The set of programs a DSL accepts is defined through a context-free grammar 139 (M, N, R, I), where M, N, R, and S are the sets of non-terminals, terminals, the production rules, 140 and the grammar's initial symbol, respectively. Figure 1 shows a DSL for a simplified version 141 of the language we use in our experiments for Karel the Robot (the complete DSL is shown in 142 Appendix H.1). In this DSL, the set M is composed of symbols  $\rho$ , h, a, while the set N includes 143 the symbols if, frontIsClear, markersPresent, move, putMarker, pickMarker. R144 are the production rules (e.g.,  $h \to \text{frontIsClear}$ ), and  $\rho$  is the initial symbol. We denote 145 programmatic policies with letters p and n and their variations such as n' and  $n^*$ . 146

We represent programs as abstract syntax trees (AST), where each node n and its children represent a production rule if n represents a non-terminal symbol. For example, the root of the tree in Figure 1, which represents the non-terminal  $\rho$ ; node  $\rho$  and its children represent the production rule  $\rho \rightarrow \text{if } h$  then a. Leaf nodes in the AST represent terminal symbols. Figure 1 shows an example of an AST for the program if markersPresent then pickMarker. A DSL D defines the possibly infinite space of programs  $[\![D]\!]$ , where in our case each program p in  $[\![D]\!]$  represents a policy.

Given a domain-specific language D, our task is to find a programmatic policy  $p \in \llbracket D \rrbracket$  that maximizes the expected sum of discounted rewards for a given MDP.

155 156

#### 3 INNATECODER

157 158

Figure 2 shows a schematic view of INNATECODER, which receives an MDP and a DSL and returns a programmatic policy, denoted  $n^*$ . INNATECODER is composed of three components: one for learning options, another that uses learned options to induce an approximation of the semantic space of the DSL, and a component to search in such a space. In this section, we explain the three components.

#### 162 3.1 LEARNING OPTIONS 163

164 An option is a program encoding a policy that the agent can invoke at specific states. Once invoked, 165 the program tells the agent what to do for a number of steps. Once completed, the option "returns" 166 the control back to the agent. An option  $\omega$  is defined with a tuple  $(I_{\omega}, \pi_{\omega}, T_{\omega})$ , where  $I_{\omega}$  is the set of 167 states in which the option can be initiated;  $\pi_{\omega}$  is the policy the agent follows once  $\omega$  starts;  $T_{\omega}$  is a function that returns the probability in which  $\omega$  terminates at a given state  $s_t$ . INNATECODER uses a 168 foundation model to learn programs, written in a given DSL, encoding options. The programs receive a state of the MDP and return the action the agent should take at that state, thus encoding  $\pi_{\omega}$ . 170

171 We assume that the set  $I_{\omega}$  is the set S of all states of the MDP, which means that the program can be 172 invoked in any state. However, note that the program may not return any actions in a given state s, which is equivalent to s not being in  $I_{\omega}$ . For example, "if  $b_1$  then  $c_1$ " returns the action given in  $c_1$ 173 only if condition  $b_1$  is satisfied in the state in which the option was queried; it returns no action for 174 states that do not satisfy  $b_1$ . The option termination criterion  $T_{\omega}$  is also determined by the program; 175 the option terminates when the program terminates. This termination criterion means that, depending 176 on the DSL, options have an internal state, representing the line in the program in which the execution 177 will continue the next time the agent interacts with the environment. For example, if the option " $c_1 c_2$ " 178 is invoked for state  $s_t$  and  $c_1$  returns an action, then the agent's action in  $s_{t+1}$  is determined by  $c_2$ . 179

Programmatic options are harnessed from a foundation model as follows. We provide a natural language description of the MDP and the Backus-Naur form of the DSL to the model. The model then 181 provides a set of m programs written in the DSL encoding policies for the MDP. While it is unlikely 182 that the model can provide policies that can maximize the expected return for any MDP of interest, 183 we hypothesize that these programmatic policies can be broken up into sub-programs that can encode helpful agent behavior. Each program p is broken up into one sub-program for each sub-tree rooted 185 at a non-terminal symbol in the AST of p. For example, for the program "if  $b_1$  then  $c_1 c_2$ " we obtain the sub-programs "if  $b_1$  then  $c_1 c_2$ ", " $b_1$ ", " $c_1$ ", " $c_2$ ", and " $c_1 c_2$ ". These sub-programs form a set of 186 187 options O, which INNATECODER uses to approximate the semantic space of the DSL. Note that this 188 set of options is generated zero-shot, before the agent starts interacting with the environment.

189 190

191

197

199

201

#### 3.2 APPROXIMATING THE SEMANTIC SPACE WITH OPTIONS

192 Methods for searching for programmatic policies traditionally search in the space of programs defined 193 by the context-free grammar of the DSL (Koza, 1992; Verma et al., 2018; Carvalho et al., 2024). We 194 refer to this type of space as the syntax space, since it is based on the syntax of the language. 195

196 **Definition 1 (Syntax Space)** The syntax space of a DSL D is defined by  $(D, \mathcal{N}_{k}^{*}, \mathcal{I}, \mathcal{E})$ . With  $[\![D]\!]$ defining the set of candidate programs, or solutions,  $\mathcal{N}_k^x$  (x is for "syntax") is the syntax neighborhood function that receives a candidate and returns k candidates from [D].  $\mathcal{I}$  is the distribution of initial candidates. Finally,  $\mathcal{E}$  is the evaluation function, which receives a candidate in [D] and returns a 200 value in  $\mathbb{R}$ .

202 A common way of defining the distribution of initial candidates  $\mathcal{I}$  is through a procedure that starts 203 with a string that is the initial symbol of the grammar and iteratively, and uniformly at random, 204 samples a production rule to replace a non-terminal symbol in the string. In the example of Figure 1, 205 we replace the initial symbol I with "if(B) then C" with probability 1.0, since this is the only rule available; then, B is replaced with either " $b_1$ " or " $b_2$ " with probability 0.5 each. This iterative process 206 207 stops once the string contains only terminal symbols. If a probabilistic context-free grammar is available, then the distribution  $\mathcal{I}$  can be defined through the same process, but using the probabilities 208 from the grammar as opposed to a uniform distribution over production rules (Trivedi et al., 2021). 209

210 The syntax neighborhood function  $\mathcal{N}_k^k$  defines the structure of the search space, as it determines the 211 set of candidate solutions (programs that encode a policy) that the search procedure can evaluate from a given candidate n. Given a candidate n,  $\mathcal{N}_k^x(n)$  returns a set of k neighbors of n. These candidates 212 213 are generated by selecting uniformly at random a node that represents a non-terminal symbol in the AST of n. Then, the sub-tree rooted at the selected node is replaced by another sub-tree generated 214 using the process described for  $\mathcal{I}$ , but starting at the non-terminal symbol the node represents. This 215 process of replacing a sub-tree in n is repeated k times, to generate k possibly different neighbors of 216 217 218 *n*. Finally,  $\mathcal{E}$  is an approximation of the value function of the policy encoded in *n* from a set of initial states  $s_0$ ,  $\hat{V}^n(s_0)$ ; we obtain  $\hat{V}^n(s_0)$  by averaging the returns after rolling *n* out from states  $s_0$ .

219 Moraes & Lelis (2024) showed that searching in the syntax space can be inefficient because often the neighbors n' of a candidate n encode policies that are semantically identical to n; the programs differ 220 in terms of syntax, but encode exactly the same agent behavior. As a result, the search process wastes 221 time evaluating the same agent behavior. Their solution is to approximate the underlying semantic 222 space of the language, where neighbor programs are similar in terms of syntax, but are likely to differ 223 in terms of behavior. In their setting, the agent learns programs for a set of tasks, which are used 224 to induce the semantic space, and the induced space is used in downstream tasks. INNATECODER 225 overcomes the requirement to operate on a stream of problems by using a foundation model to learn 226 the options in a zero-shot setting. A semantic space is defined as follows.

**Definition 2 (Semantic Space)** The semantic space of a DSL D is defined by  $(D, \mathcal{N}_k^m, \mathcal{I}, \mathcal{E})$ , where

 $\mathcal{I}$  and  $\mathcal{E}$  are identical to the syntax space (Definition 1). The function  $\mathcal{N}_k^m$  (m is for "semantics") is

a semantic neighborhood function that also receives a candidate and returns k candidates from [D].

- 227
- 228

229

- 230
- 231

We define the function  $\mathcal{N}_k^m$  with a set of options  $\Omega$ , where each option  $\omega$  in  $\Omega$  represents a different agent behavior. A neighbor of candidate n is then obtained by selecting, uniformly at random, a node c in the AST of n that represents a non-terminal symbol. Then, we replace the sub-tree rooted at cwith the AST of an option  $\omega$  in  $\Omega$ . The option  $\omega$  is selected, uniformly at random, among those in  $\Omega$ whose AST root represents the same non-terminal symbol c represents. By matching the non-terminal symbols when selecting  $\omega$ , we match  $\omega$  with the type of the sub-tree that is removed from n. Similarly to  $\mathcal{N}_k^x$ ,  $\mathcal{N}_k^m$  generates k possibly different neighbors by repeating this process k times.

The intuition for requiring the options in  $\Omega$  to encode different agent behaviors is to increase the chance of seeing neighbors with different behaviors. For example, if many of the options in  $\Omega$ encoded exactly the same behavior, then the chances of all neighbors of a program also encoding the same behavior would be higher, which is wasteful from a sample efficiency perspective. Next, we describe how we obtain  $\Omega$  from the set of options *O* we harnessed from a foundation model.

244 We filter the set of options O harnessed from the foundation model into a set  $\Omega$  of behaviorally differ-245 ent options by having the agent interact with the environment, as described in previous work (Moraes 246 & Lelis, 2024). Assuming an MDP with discrete actions, each option in O is evaluated in an ordered 247 set of states S of the MDP. This set S is obtained by rolling out all options  $o \in O$  once, from an initial state  $s_0$  sampled from  $\mu$ . The states s observed in this process form S. Then, every option is 248 invoked for each state in S, thus forming an *action signature*  $A_o$  for each o. An action signature is 249 a vector with one action for each state in S, where the *i*-th entry of  $A_o$  corresponds to the action o 250 returns to the *i*-th state in S. The set of options  $\Omega$  is given by one option for each observed  $A_{\rho}$ . If 251 multiple options have the same signature, we arbitrarily select one of them. 252

The number of samples required to filter the set of options into a set of options with different behaviors is negligible: it uses less than 1% of the computation in our experiments. Programs that cannot be rolled out (e.g., Boolean expressions and do not issue actions) are not included in the set of options.

256 257

258

#### 3.3 SEARCHING IN SEMANTIC SPACE

259 INNATECODER uses stochastic hill-climbing (SHC) to search in the semantic space of a given DSL 260 D for a policy that maximizes the agent's return. SHC starts its search by sampling a candidate program n from  $\mathcal{I}$ . In every iteration, SHC evaluates all k neighbors of n in terms of their  $\mathcal{E}$ -value. 261 The search then moves on to the best neighbor of n in terms of  $\mathcal{E}$ , and this process is repeated from 262 there. SHC stops if none of the neighbors has an  $\mathcal{E}$ -value that is better than the current candidate, that 263 is, it reaches a local optimum. SHC uses a restarting strategy: once SHC reaches a local optimum, if 264 SHC has not yet exhausted is search budget, it restarts from another initial candidate sampled from  $\mathcal{I}$ . 265 SHC returns the best solution, denoted  $n^*$ , encountered in all restarts of the search. 266

**267** INNATECODER does not search solely in the semantic space, but mixes both syntax and semantic 268 spaces in the search. This is because the set of options might cover only a part of the space of 269 programs the DSL induces. To guarantee that INNATECODER can access all programs in [D], with probability  $\epsilon$ , SHC uses the syntax neighborhood function in the search, and with probability  $1 - \epsilon$ , it uses the semantic one. We use  $\epsilon = 0.4$  in our experiments. We chose this value because it performed better in preliminary experiments than the value of 0.2 used in previous work (Moraes & Lelis, 2024).

Although other local search algorithms could be used with INNATECODER, such as Simulated Annealing (Kirkpatrick et al., 1983; Husien & Schewe, 2016), we use SHC because previous work showed that it performs well in our test domains (Moraes et al., 2023; Carvalho et al., 2024).

275 276 277

278

284

273

274

# 4 EMPIRICAL EVALUATION

Although foundation models are unlikely to generate programs that encode policies for fully solving
MDPs, we hypothesize that the programs they generate can be broken up into smaller programs that
serve as helpful options. We evaluated the usefulness of these programmatic options by measuring
the sampling efficiency of search algorithms searching in the semantic spaces induced by them. We
evaluate INNATECODER on MicroRTS (Ontañón, 2017) and Karel the Robot (Pattis, 1994).

**MicroRTS** MicroRTS is a real-time strategy game that requires the agent to control dozens of 285 units in real-time, thus making it impractical to use foundation models to decide on agent actions 286 directly. We use the following maps from the MicroRTS repository,<sup>1</sup> with the map size in brackets: 287 NoWhereToRun ( $9 \times 8$ ), basesWorkers ( $24 \times 24$ ), and BWDistantResources ( $32 \times 32$ ), and BloodBath 288  $(64 \times 64)$ . We use these maps because they differ in size and structure. Since MicroRTS is a multi-289 agent problem, we use 2L, a self-play algorithm, to learn programmatic policies (Moraes et al., 2023). 290 MicroRTS is not a symmetric game and the outcome of the game depends on the starting location of 291 the players. To ensure fairness, each pair of policies plays two matches on each map, so that each 292 player can start at each of the two initial locations; the results are then averaged out of these two 293 runs. In the context of 2L, INNATECODER is required to solve an MDP in every iteration of self-play (see Appendix M). We use a new version of the MicroLanguage as the DSL (Mariño et al., 2021). 294 The language offers specialized functions and an action-prioritization scheme through for-loops, 295 where nested for-loops allow for higher priority of actions. We provide a detailed explanation of the 296 MicroLanguage, as well as images of the maps used, in Appendices H.2 and K, respectively. 297

298

Karel Karel the Robot is an environment originally created for teaching people how to write 299 computer programs, which has later been used as a benchmark domain for reinforcement learning 300 algorithms (Trivedi et al., 2021). Karel is a robot interacting with a grid-world, where it can collect 301 markers and place markers. We use the following Karel problems, which were designed in previous 302 works (Trivedi et al., 2021; Liu et al., 2023b): StairClimber, FourCorners, TopOff, Maze, CleanHouse, 303 Harvester, DoorKey, OneStroke, Seeder, and Snake. The problems differ in terms of structure of the 304 grids (e.g., where walls are located) and in terms of the task that Karel needs to accomplish. For 305 example, in CleanHouse, Karel needs to collect all markers placed in the grid, while in TopOff it has 306 to place a marker on top of all existing markers. The problem are described in Appendix I. We use 307 the more difficult version of the environment known as "crashable" (Carvalho et al., 2024), where an 308 episode terminates with a negative reward if Karel bumps into a wall. We use the same DSL used in previous work (Trivedi et al., 2021), which we describe in Appendix H.1. 309

310

**Baselines** The current state-of-the-art methods for both MicroRTS and Karel use programmatic rep-311 resentations of policies, where the policies are written in the DSLs we use in our experiments (Moraes 312 et al., 2023; Trivedi et al., 2021). Therefore, we focus on methods that use programmatic representa-313 tions as baselines. However, we provide comparisons of INNATECODER with deep reinforcement 314 learning baselines in Appendices B.1 (MicroRTS) and B.2 (Karel). For both MicroRTS and Karel we 315 use SHC searching in the syntax space as a baseline, as it represents state-of-the-art performance 316 in both domains (SHC). We also use two variants of INNATECODER where the options are learned 317 without the help of a foundation model. These variants can be seen as implementations of the 318 Library-Induced Semantic Spaces (LISS) (Moraes & Lelis, 2024) for a non-transfer learning setting. 319 In the first variant, LISS learns the options as it learns how to solve the problem. In terms of the 320 scheme shown in Figure 2, we skip the "Learning Options" step and build the set of options from the 321 programs returned in every complete search of SHC. That is, when we reach the box "Restart?", we 322 use the sub-programs of the best program encountered in that search to augment the set of options.

<sup>&</sup>lt;sup>1</sup>https://github.com/Farama-Foundation/MicroRTS/

324 We call this baseline LISS-o, where "o" stands for "online". In the second variant, we sample 325 programs from  $\mathcal{I}$  and use their sub-programs to form the set of options. We call this baseline LISS-r, 326 where "r" stands for "random". We also use the best program the foundation model generated out 327 of all programs used to create the set of options as a baseline, which we call FM, which stands for 328 "foundation model". LISS-o and LISS-r allow us to evaluate the effectiveness of learning options from a foundation model, while FM allows us to evaluate the foundation model as an alternative to 329 solve the problem directly. We also use the Cross Entropy Method (CEM) operating in a learned 330 latent space, which was shown to outperform DRL algorithms in all Karel tasks (Trivedi et al., 2021). 331

332 Foundation Models We use OpenAI's API for GPT 40, whose training cut-off date is October 333 2023. We also perform tests, for MicroRTS, using the LLama 3.1 model with 405 billion parameters, 334 whose training cut-off is December 2021. We used the GPT model in both MicroRTS and Karel 335 experiments, while the Llama model was used in MicroRTS experiments. There were no MicroRTS 336 programs available online prior to the Llama cut-off date, so the Llama evaluations on MicroRTS 337 did not suffer from data leakage. The GPT model might have trained on the MicroRTS and Karel 338 programs that were available online prior to its training cut-off date. We attempt to measure how 339 much a possible data leakage can influence our results by using the FM baseline. If the model can 340 simply retrieve the solutions seen in training, one would expect this baseline to perform well.

342 **Other Specifications** All experiments were run on 2.6 GHz CPUs with 12 GB of RAM. We use 343 k = 1,000 in the neighborhood function. In MicroRTS, SHC is run with a restarting time limit of 344 2,000 seconds for each self-play iteration. In Karel, since we are solving a single MDP, SHC restarts as many times as possible within the computational budget. For MicroRTS, we query the foundation 345 models 120 times to generate the same number of programs; for Karel, we use 100 programs. We use 346 the same number of programs as the LISS-r baseline. We perform 30 independent runs (seeds) of 347 each system, including the generation of the programs by the foundation model. 348

349 **Metrics of Performance** For MicroRTS, performance is measured in terms of winning rate. The 350 winning rate of a policy is computed for a set of opponent policies and is computed as follows: 351 we sum the number of victories and half the number of draws and divide this sum by the total 352 number of matches played (Ontañón, 2017). For Karel, performance is measured in terms of episodic 353 return (Trivedi et al., 2021). We use prompts where we briefly describe each problem and provide 354 a formal description of the DSL used. The prompts used in our experiments are in Appendices L 355 (MicroRTS) and J (Karel). Both MicroRTS and Karel are deterministic, so the value of  $\mathcal{E}$  for policies 356 can be computed with a single roll-out. We report average performance and 95% confidence intervals.

357

341

358 **Efficiency Experiment** We verify the sampling efficiency of INNATECODER, LISS-o, LISS-r, and SHC. Similarly to previous work, we present learning curves, where for MicroRTS, we plot winning 359 rate by the number of games played (Figure 3), and for Karel, we plot episodic return by the number 360 of episodes (Figure 4). For MicroRTS, the winning rate is computed for a system by having the policy 361 the system generated, after a given number of games played, play against the policies each of the 362 other systems generated after the maximum number of games played (rightmost point of each plot). 363

364 **Competition Experiment** We evaluate INNATECODER against COAC, Mayari, and RAISocketAI, the winners of the previous three MicroRTS competitions. We randomly select 9 from the 30 programs 366 generated in the "Efficiency Experiment" and evaluate them against the competition winners. We report the average results of the 9 programs against each opponent in the four maps we use. 368

369 **Size and Information Experiments** We also evaluate the effect of the size of the set of options on 370 the sample efficiency of INNATECODER. We evaluated sets  $\Omega$  with 300, 600, 1400, 5000, 7000 and 371 30000 options on the LetMeOut map  $(16 \times 8)$ ; all options were generated with the Llama 3.1 model. 372 In Appendix C, we evaluate INNATECODER using prompts with more or less information and in 373 Appendix D with GPT 3.5, to verify if performance decreases by using a smaller foundation model.

374

376

367

375 4.1 LEARNING CURVE RESULTS

Figures 3 and 4 show the learning curves for MicroRTS and Karel, respectively, where INNATECODER 377 is denoted as IC-GPT or IC-Llama, depending on the model it uses to learn the programmatic options.



Figure 3: Winning rate (maximum is 100) per number of games played. The winning rate of the policies each system generates for a given number of games played is computed considering as opponents the policies all systems generate at the end of the learning process. The plots show the average winning rate of 30 independent runs (seeds) and the 95% confidence interval.



Figure 4: Average episodic return (maximum is 1.0 for all tasks) per number of episodes. The plots show average episodic return of 30 independent runs (seeds) and the 95% confidence interval.

INNATECODER is often much more sample-efficient than all baselines and, in many cases, by a large margin. We did not observe significant differences between IC-GPT and IC-Llama. LISS-o and LISS-r perform worse than INNATECODER and SHC in MicroRTS. However, LISS-o was competitive with SHC in Karel and LISS-r could outperform SHC (DoorKey and Seeder). This result suggests that the semantic space can be less conducive to search than the syntax space, depending on the quality of the options used to induce it. LISS-r performs better in Karel than in MicroRTS, probably because it uses a distribution  $\mathcal{I}$  that uses a handcrafted probability distribution over the production rules of the language (Trivedi et al., 2021). The resulting grammar allows for the generation of helpful options. We do not have such a distribution for the MicroLanguage, which explains the results. 

FM performs poorly in all experiments; the model is unable to generate effective policies in a zero-shot setting. The results of FM and INNATECODER support our hypothesis that INNATECODER can extract helpful options from foundation models even if the programs the model generates do not encode strong policies. In MicroRTS, some of the options allowed the agent to allocate units to collect resources and train other units. Other systems had to learn such skills from scratch, while INNATECODER's agent had them "innately available". The FM results also suggest that data contamination was not an issue in our experiments, as the model performed poorly on all tasks.

INNATECODER	COAC	RAISocketAI	Mayari	Average
GPT-40	53.75	36.25	71.25	53.75
Llama 3.1	43.79	70.00	58.17	57.32
Llama 3.1 + GPT-40	70.14	72.92	46.39	63.15

Table 1: Winning rate of INNATECODER against winners of previous competitions, averaged across all 4 maps used in our experiments.

442

438

#### 4.2 COMPETITION RESULTS

443 Table 1 shows the results of INNATECODER against the winners of previous MicroRTS competitions. 444 The numbers are the average winning rate of INNATECODER against each system in the 4 maps used 445 in our experiments. COAC and Mayari are human-written programmatic policies, and RAISocketAI 446 is a DRL agent (Goodfriend, 2024). We used the RAISocketAI model submitted to the competition, 447 which was trained with a larger computational budget than what we used with INNATECODER, thus 448 giving RAISocketAI an advantage. We evaluated the models GPT-40 while generating 120 programs 449 from which options are extracted, as well as Llama 3.1 while generating the same number of programs. 450 Finally, we also evaluate INNATECODER when we take the union of the programs generated by both GPT-40 and Llama (denoted Llama 3.1 + GPT-40 in the table). The larger number of programs 451 considered in the combination of GPT-40 and Llama 3.1 resulted in the best average winning rate. 452

The combination of programs written by Llama 3.1 and GPT-40 does not lead to "monotonic improvements", as evidenced by the drop in performance against Mayari. This happens because none of the competition winners is constrained by the DSL we use in our experiments. As a result, the optimization done in self-play might not be specific for the opponents evaluated in Table 1, but to policies written in the DSL and encountered during the self-play process.

458 459

460

#### 4.3 EVALUATING NUMBER OF OPTIONS

461 Figure 5 presents the learning curves for 462 different versions INNATECODER, where 463 we vary the size of the set of options  $\Omega$ . The 464 three lines with the highest winning rate 465 are for option sets of sizes 5000, 7000, and 30000. The versions of INNATECODER 466 with option sets of sizes 300, 600, and 467 1400 perform worse. These results demon-468 strate that INNATECODER can benefit from 469 thousands of options. This is possible due 470 to INNATECODER's way of using options 471 through the induction of the language's un-472 derlying semantic space. 473

These results also show that INNATE-CODER's sample efficiency plateaus at 5000 options, since the use of 7000 and 30000 options does not increase performance. Interestingly, performance does not degrade either as we increase the set size.
We conjecture that this occurs because, for large sets, many of the options will encode



Figure 5: Average winning rate of INNATECODER policies for different sizes of the option set over 10 independent runs (seeds) of each version. We also present the 95% confidence intervals.

different and yet similar behaviors that do not affect the agent's winning rate. For example, an agent could use options  $\omega_1$  or  $\omega_2$  to achieve slightly different behaviors that lead to the same winning rate. That is, although the set of distinct behaviors encoded in the set options grows with larger sets, the relative number of options with behaviors that affect the winning rate remains roughly the same. As a result, the neighborhood function  $\mathcal{N}_k^m$  that uses option sets of sizes 5000, 7000, or 30000 induces spaces that are similarly conducive to search algorithms. In Appendix F, we explain that the difference in performance between INNATECODER with 1400 or fewer options and INNATECODER with 5000 or more is due to an increased chance of sampling helpful options with larger sets.

### 5 RELATED WORK

**Programmatic Policies** One of the key challenges in generating programmatic policies is that the search space is discontinuous and gradient-based optimization cannot be used. Some previous work relied on imitation learning to guide the search for policies (Verma et al., 2018; 2019; Bastani et al., 2018; Milani et al., 2022; Liu et al., 2023d). The issue of this imitation learning approach is known as *representation gap* (Qiu & Zhu, 2022; Medeiros et al., 2022), where the space of programmatic policies does not include the oracle policy that the system tries to imitate. As a result, the oracle might guide the search to unpromising parts of the space. Previous work tried to learn latent spaces of programming languages that are conducive to search (Trivedi et al., 2021; Liu et al., 2023b), which was shown to be outperformed by the syntax space with SHC (Carvalho et al., 2024). Semantic spaces were shown to be more conducive to search than syntax spaces, but required a sequence of tasks, where the agent learns the space in one task and reuses it in others (Moraes & Lelis, 2024). Our work does not require an oracle agent nor a sequence of tasks to learn the semantic space.

503 504

486

487

488 489 490

491 492

493

494

495

496

497

498

499

500

501

502

505 **Options** Options were shown to improve the sampling efficiency of learning agents through faster 506 credit assignment (Mann & Mannor, 2014; Solway et al., 2014), better exploration (Baranes & 507 Oudeyer, 2013; Bellemare et al., 2020), and transfer of knowledge across tasks (Konidaris & Barto, 508 2007; Alikhasi & Lelis, 2024). However, previous methods for learning options require the user to 509 design them before learning starts (Sutton et al., 1999) or to provide considerable information as input to the process, such as the option duration (Frans et al., 2017; Tessler et al., 2017) or the number 510 of options learned (Bacon et al., 2017; Igl et al., 2020). Other methods rely on the agent interaction 511 with the current environment (Achiam et al., 2018; Machado et al., 2018; Jinnai et al., 2020) or with 512 other earlier environments, as in transfer learning approaches (Konidaris & Barto, 2007; Alikhasi & 513 Lelis, 2024). We present a novel way of learning options as they are not learned from the agent's 514 experience nor designed by the user, but harnessed from foundation models. While we use options to 515 define a search space, future work will explore their use as functions neural policies can call. 516

517 We provide additional related works on "foundation models as policies", "foundation models for 518 planning", and "foundation models as search guidance" in Appendix A.

- 519
- 520 521

## 6 CONCLUSIONS

522

523 If given a single problem to solve, reinforcement learning agents start their learning process from 524 scratch. They have to learn by interacting with the environment even the most basic skills to solve 525 the problem. In this paper, we introduced INNATECODER, a system that equips learning agents with skills, in the form of programmatic options, before the agent starts to interact with the environment. 526 This is achieved by extracting programmatic options from foundation models. We hypothesized that 527 even if the model is unable to write programs encoding strong policies for a problem, sub-programs 528 of the generated program could encode helpful options. We tested our hypothesis in MicroRTS and 529 Karel, two domains in which programmatic policies represent the current state of the art. The policies 530 INNATECODER generated outperformed, often by a large margin, a baseline that did not attempt to 531 learn options; a baseline that learned the options while learning how to solve the problem; a baseline 532 that learned the options from programs sampled directly from the domain-specific language; and the 533 foundation model that attempted to generate programmatic policies directly. We also showed that 534 some of the policies INNATECODER generated were competitive or outperformed the winners of previous MicroRTS competitions, including programmatic policies written by human programmers 536 and a deep reinforcement learning agent that used a larger computational budget than we allowed 537 INNATECODER to use. These results place INNATECODER as the current state-of-the-art in both Karel and MicroRTS. Our experiments also showed that INNATECODER's scheme of using programmatic 538 options to induce semantic spaces allows it to benefit from thousands of options, while most previous work can benefit only from dozens or at most hundreds of options (Eysenbach et al., 2019).

# 540 REFERENCES

546

- Joshua Achiam, Harrison Edwards, Dario Amodei, and Pieter Abbeel. Variational option discovery algorithms. *arXiv preprint arXiv:1807.10299*, 2018.
- Mahdi Alikhasi and Levi H. S. Lelis. Unveiling options with neural network decomposition. In *The Twelfth International Conference on Learning Representations*, 2024.
- 547 Pierre-Luc Bacon, Jean Harb, and Doina Precup. The option-critic architecture. In *Proceedings of* 548 *the AAAI conference on artificial intelligence*, volume 31, 2017.
- Adrien Baranes and Pierre-Yves Oudeyer. Active learning of inverse models with intrinsically motivated goal exploration in robots. *Robotics and Autonomous Systems*, 61:49–73, 2013. ISSN 0921-8890.
- Osbert Bastani, Yewen Pu, and Armando Solar-Lezama. Verifiable reinforcement learning via policy
   extraction. In *Proceedings of the International Conference on Neural Information Processing Systems*, pp. 2499–2509, 2018.
- Marc G. Bellemare, Salvatore Candido, Pablo S. Castro, Jun Gong, Marlos C. Machado, Subhodeep
   Moitra, Sameera Ponda, and Ziyu Wang. Autonomous navigation of stratospheric balloons using
   reinforcement learning. *Nature*, 588:77–82, 2020.
- GW Brown. Iterative solution of games by fictitious play. Activity Analysis of Production and Allocation, pp. 374–376, 1951.
- Rudy Bunel, Matthew Hausknecht, Jacob Devlin, Rishabh Singh, and Pushmeet Kohli. Leveraging
   grammar and reinforcement learning for neural program synthesis. In *International Conference on Learning Representations*, 2018.
- Tales Henrique Carvalho, Kenneth Tjhia, and Levi H. S. Lelis. Reclaiming the source of programmatic policies: Programmatic versus latent spaces. In *The Twelfth International Conference on Learning Representations*, 2024.
- Sergio Jiménez Celorrio, Javier Segovia Aguas, and Anders Jonsson. A review of generalized planning. *Knowl. Eng. Rev.*, 34:e5, 2019. doi: 10.1017/S0269888918000231. URL https://doi.org/10.1017/S0269888918000231.
- Angelica Chen, David M. Dohan, and David R. So. Evoprompting: Language models for code-level neural architecture search. *CoRR*, abs/2302.14838, 2023. doi: 10.48550/ARXIV.2302.14838. URL https://doi.org/10.48550/arXiv.2302.14838.
- 575 Xinyun Chen, Chang Liu, and Dawn Song. Execution-guided neural program synthesis. In *International Conference on Learning Representations*, 2018.
   577
- Benjamin Eysenbach, Abhishek Gupta, Julian Ibarz, and Sergey Levine. Diversity is all you need: Learning skills without a reward function. In *International Conference on Learning Representations* (*ICLR*), 2019. URL https://arxiv.org/abs/1802.06070.
- Kevin Frans, Jonathan Ho, Xi Chen, Pieter Abbeel, and John Schulman. Meta learning shared
   hierarchies. *arXiv preprint arXiv:1710.09767*, 2017.
- 583 584 Scott Goodfriend. A competition winning deep reinforcement learning agent in microrts, 2024.
- Pei-Fu Guo, Ying-Hsuan Chen, Yun-Da Tsai, and Shou-De Lin. Towards optimizing with large language models. *CoRR*, abs/2310.05204, 2023. doi: 10.48550/ARXIV.2310.05204. URL https://doi.org/10.48550/arXiv.2310.05204.
- Shengyi Huang, Santiago Ontañón, Chris Bamford, and Lukasz Grela. Gym-µrts: Toward affordable
   full game real-time strategy games research with deep reinforcement learning. In 2021 IEEE
   *Conference on Games (CoG)*, pp. 1–8. IEEE, 2021.
- Idress Husien and Sven Schewe. Program generation using simulated annealing and model checking.
   In *International Conference on Software Engineering and Formal Methods*, pp. 155–171. Springer, 2016.

594 595 596	Maximilian Igl, Andrew Gambardella, Jinke He, Nantas Nardelli, N Siddharth, Wendelin Böhmer, and Shimon Whiteson. Multitask soft option learning. In <i>Conference on Uncertainty in Artificial Intelligence</i> , pp. 969–978, 2020.
597 598 599 600	Jeevana Priya Inala, Osbert Bastani, Zenna Tavares, and Armando Solar-Lezama. Synthesizing programmatic policies that inductively generalize. In <i>International Conference on Learning Representations</i> , 2020.
601 602 603 604	Yuu Jinnai, Jee W. Park, Marlos C. Machado, and George Konidaris. Exploration in reinforcement learning with deep covering options. In <i>International Conference on Learning Representations</i> , 2020.
605 606	S. Kirkpatrick, C. D. Gelatt, and M. P. Vecchi. Optimization by Simulated Annealing. <i>Science</i> , 220 (4598):671–680, May 1983. doi: 10.1126/science.220.4598.671.
607 608 609	Martin Klissarov and Marlos C. Machado. Deep laplacian-based options for temporally-extended exploration. In <i>International Conference on Machine Learning</i> , 2023.
610 611 612 613	Martin Klissarov, Pierluca D'Oro, Shagun Sodhani, Roberta Raileanu, Pierre-Luc Bacon, Pascal Vincent, Amy Zhang, and Mikael Henaff. Motif: Intrinsic motivation from artificial intelligence feedback. In <i>The Twelfth International Conference on Learning Representations</i> , 2024.
614 615	George Konidaris and Andrew Barto. Building portable options: Skill transfer in reinforcement learning. In <i>International Joint Conference on Artificial Intelligence</i> , pp. 895–900, 2007.
616 617 618	J. R. Koza. Genetic Programming: On the Programming of Computers by Means of Natural Selection. MIT Press, Cambridge, MA, 1992.
619 620 621	Marc Lanctot, Vinicius Zambaldi, Audrunas Gruslys, Angeliki Lazaridou, Karl Tuyls, Julien Pérolat, David Silver, and Thore Graepel. A unified game-theoretic approach to multiagent reinforcement learning. <i>Advances in neural information processing systems</i> , 30, 2017.
622 623 624 625	Joel Lehman, Jonathan Gordon, Shawn Jain, Kamal Ndousse, Cathy Yeh, and Kenneth O. Stanley. Evolution through large models. <i>CoRR</i> , abs/2206.08896, 2022. doi: 10.48550/ARXIV.2206.08896. URL https://doi.org/10.48550/arXiv.2206.08896.
626 627 628 629 630	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 <i>IEEE</i> <i>International Conference on Robotics and Automation, ICRA 2023, London, UK, May 29 - June</i> 2, 2023, pp. 9493–9500. IEEE, 2023. doi: 10.1109/ICRA48891.2023.10160591. URL https: //doi.org/10.1109/ICRA48891.2023.10160591.
631 632 633 634 635	Fei Liu, Xi Lin, Zhenkun Wang, Shunyu Yao, Xialiang Tong, Mingxuan Yuan, and Qingfu Zhang. Large language model for multi-objective evolutionary optimization. <i>CoRR</i> , abs/2310.12541, 2023a. doi: 10.48550/ARXIV.2310.12541. URL https://doi.org/10.48550/arXiv. 2310.12541.
636 637 638	Guan-Ting Liu, En-Pei Hu, Pu-Jen Cheng, Hung-Yi Lee, and Shao-Hua Sun. Hierarchical pro- grammatic reinforcement learning via learning to compose programs. In <i>Proceedings of the</i> <i>International Conference on Machine Learning</i> , 2023b.
639 640 641 642	Shengcai Liu, Caishun Chen, Xinghua Qu, Ke Tang, and Yew-Soon Ong. Large language models as evolutionary optimizers. <i>CoRR</i> , abs/2310.19046, 2023c. doi: 10.48550/ARXIV.2310.19046. URL https://doi.org/10.48550/arXiv.2310.19046.
643 644 645	Xiao Liu, Wubing Chen, and Mao Tan. Fidelity-induced interpretable policy extraction for reinforce- ment learning. <i>CoRR</i> , abs/2309.06097, 2023d.
646 647	Marlos C. Machado, Marc G. Bellemare, and Michael Bowling. A laplacian framework for option discovery in reinforcement learning. In <i>International Conference on Machine Learning (ICML)</i> , 2017.

648	Marlos C. Machado, Clemens Rosenbaum, Xiaoxiao Guo, Miao Liu, Gerald Tesauro, and Murray
649	Campbell. Eigenoption discovery through the deep successor representation. In International
650	Conference on Learning Representations, 2018.
651	

- Timothy Mann and Shie Mannor. Scaling up approximate value iteration with options: Better policies
  with fewer iterations. In *Proceedings of the International Conference on Machine Learning*, volume 32 of *Proceedings of Machine Learning Research*, pp. 127–135, 2014.
- Julian R. H. Mariño, Rubens O. Moraes, Tassiana C. Oliveira, Claudio Toledo, and Levi H. S. Lelis.
   Programmatic strategies for real-time strategy games. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pp. 381–389, 2021.
- Leandro C. Medeiros, David S. Aleixo, and Levi H. S. Lelis. What can we learn even from the weakest? Learning sketches for programmatic strategies. In *Proceedings of the AAAI Conference on Artificial Intelligence*, pp. 7761–7769. AAAI Press, 2022.
- Elliot Meyerson, Mark J. Nelson, Herbie Bradley, Arash Moradi, Amy K. Hoover, and Joel Lehman. Language model crossover: Variation through few-shot prompting. *CoRR*, abs/2302.12170, 2023. doi: 10.48550/ARXIV.2302.12170. URL https://doi.org/10.48550/arXiv.2302.
   12170.
- Stephanie Milani, Zhicheng Zhang, Nicholay Topin, Zheyuan Ryan Shi, Charles A. Kamhoua,
  Evangelos E. Papalexakis, and Fei Fang. MAVIPER: learning decision tree policies for interpretable
  multi-agent reinforcement learning. In *Machine Learning and Knowledge Discovery in Databases - European Conference*, volume 13716 of *Lecture Notes in Computer Science*, pp. 251–266. Springer,
  2022.
- Rubens O. Moraes and Levi H. S. Lelis. Searching for programmatic policies in semantic spaces.
   2024.
- Rubens O. Moraes, David S. Aleixo, Lucas N. Ferreira, and Levi H. S. Lelis. Choosing well
  your opponents: How to guide the synthesis of programmatic strategies. In *Proceedings of the International Joint Conference on Artificial Intelligence*, pp. 4847–4854, 2023.
- Muhammad Umair Nasir, Sam Earle, Julian Togelius, Steven James, and Christopher W. Cleghorn.
   LImatic: Neural architecture search via large language models and quality-diversity optimization.
   *CoRR*, abs/2306.01102, 2023. doi: 10.48550/ARXIV.2306.01102. URL https://doi.org/
   10.48550/arXiv.2306.01102.
- Santiago Ontañón. Combinatorial multi-armed bandits for real-time strategy games. *Journal of Artificial Intelligence Research*, 58:665–702, 2017.
- Joon Sung Park, Joseph C. O'Brien, Carrie Jun Cai, Meredith Ringel Morris, Percy Liang, and
   Michael S. Bernstein. Generative agents: Interactive simulacra of human behavior. In *Proceedings of the Annual ACM Symposium on User Interface Software and Technology*, pp. 2:1–2:22. ACM,
   2023.
  - Richard E Pattis. *Karel the robot: a gentle introduction to the art of programming.* John Wiley & Sons, 1994.

689

- Doina Precup, Richard S. Sutton, and Satinder Singh. Theoretical results on reinforcement learning with temporally abstract options. pp. 382–393, Berlin, Heidelberg, 1998. Springer Berlin Heidelberg.
- Wenjie Qiu and He Zhu. Programmatic reinforcement learning without oracles. In *The Tenth International Conference on Learning Representations*, 2022.
- Eui Chul Shin, Illia Polosukhin, and Dawn Song. Improving neural program synthesis with inferred execution traces. *Advances in Neural Information Processing Systems*, 31, 2018.
- 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. *CoRR*, abs/2305.11014, 2023. doi: 10.48550/ARXIV.2305.11014. URL https://doi.org/10.48550/arXiv.2305.11014.

702	Ishika Singh, Valts Blukis, Arsalan Mousavian, Ankit Goyal, Danfei Xu, Jonathan Tremblay, Dieter
703	Fox, Jesse Thomason, and Animesh Garg. Progprompt: Generating situated robot task plans using
705	large language models. In IEEE International Conference on Robotics and Automation, ICRA 2023, London, UK, May 20, June 2, 2023, pp. 11523, 11530, IEEE, 2023, doi: 10.1100/ICPA48801
706	2023.10161317. URL https://doi.org/10.1109/ICRA48891.2023.10161317.
707	2020.10101017. CRE heepo., , doi:org, 10.1100, 101110091.2020.10101017.
708	Alec Solway, Carlos Diuk, Natalia Córdova, Debbie Yee, Andrew G. Barto, Yael Niv, and Matthew M.
709	Botvinick. Optimal behavioral hierarchy. PLOS Computational Biology, 10:1–10, 2014.
710	Richard S Sutton, Doina Precup, and Satinder Singh. Between MDPs and semi-MDPs: A framework
711	for temporal abstraction in reinforcement learning. Artificial intelligence, 112(1-2):181-211, 1999.
712	Chen Tessler, Shahar Givony, Tom Zahavy, Daniel Mankowitz, and Shie Mannor. A deep hierarchical
713	approach to lifelong learning in minecraft. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , volume 31, 2017.
715	
716	Nikolaas Tinbergen. The Study of Instinct. Oxford University Press, 1951.
718	Dweep Trivedi, Jesse Zhang, Shao-Hua Sun, and Joseph J. Lim. Learning to synthesize programs as
719	interpretable and generalizable policies. In Advances in Neural Information Processing Systems,
720	pp. 25146–25163, 2021.
721	A. Verma, V. Murali, R. Singh, P. Kohli, and S. Chaudhuri. Programmatically interpretable rein-
722	forcement learning. In Proceedings of the International Conference on Machine Learning, pp.
723	5052–5061, 2018.
724	Abbinay Verma Hoang Le Visong Yue and Swarat Chaudhuri Imitation-projected programmatic
725 726	reinforcement learning. Advances in Neural Information Processing Systems, 32, 2019.
727	Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and
728	Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. CoRR,
729	abs/2305.16291, 2023a. doi: 10.48550/ARXIV.2305.16291. URL https://doi.org/10.
730	48550/arXiv.2305.16291.
731	Zihao Wang, Shaofei Cai, Anji Liu, Xiaojian Ma, and Yitao Liang. Describe, explain, plan and select:
732	Interactive planning with large language models enables open-world multi-task agents. CoRR,
733	abs/2302.01560, 2023b. doi: 10.48550/ARXIV.2302.01560. URL https://doi.org/10.
734	48550/arXiv.2302.01560.
736	Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen.
737	Large language models as optimizers. CoRR, abs/2309.03409, 2023. doi: 10.48550/ARXIV.2309.
738	<b>03409. URL</b> https://doi.org/10.48550/arXiv.2309.03409.
739	Zirui Zhao, Wee Sun Lee, and David Hsu. Large language models as commonsense knowledge for
740	large-scale task planning. CoRR, abs/2305.14078, 2023. doi: 10.48550/ARXIV.2305.14078. URL
741	https://doi.org/10.48550/arXiv.2305.14078.
742	Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li
743	Lewei Lu, Xiaogang Wang, Yu Qiao, Zhaoxiang Zhang, and Jifeng Dai. Ghost in the minecraft:
744	Generally capable agents for open-world environments via large language models with text-based
740	knowledge and memory. CoRR, abs/2305.17144, 2023. doi: 10.48550/ARXIV.2305.17144. URL
747	https://doi.org/10.48550/arXiv.2305.17144.
748	
749	A ADDITIONAL RELATED WORK
750	

Foundation Models as Policies Foundation models have been used to perceive, plan, and act (Park et al., 2023), often decomposing long-horizon goals into subtasks (Wang et al., 2023b), and/or integrating additional agent features such as memory (Zhu et al., 2023) and/or automatic learning curricula (Wang et al., 2023a). By contrast, INNATECODER uses the model as a pre-processing step to generate programmatic options, which make it more accessible due to the limited number of model calls. Foundation models have also been used to learn reward functions that are later used to

train agents (Klissarov et al., 2024). Similarly to our work, the model is used in a pre-processing
step. However, in contrast with our work, it learns reward functions, while we learn programmatic
options. Moreover, it needs a "diverse" set of states, which are generated by existing and proficient
agents; INNATECODER learns in a zero-shot setting and only uses the generated options themselves
to generate a set of states used to filter out the options encoding non-novel behaviors.

Foundation Models for Planning Previous work has used foundation models for generating programmatic policies in the context of planning. For example, in generalized planning (GP), the goal is to synthesize programs that solve classical planning problems (Celorrio et al., 2019); foundation models have shown promise for GP (Silver et al., 2023). Foundation models has also been used successfully in the context of code generation for decision making in robotics (Liang et al., 2023; Singh et al., 2023). In contrast to these works, we do not attempt to use the model's generated program as a policy, but we extract options from them and use these options to induce a search space.

769 770 771

772

773

774

775

776

777 778 779

780 781

782

761

**Foundation Models as Search Guidance** Foundation models have also been used to guide search algorithms. This includes methods for solving optimization problems (Yang et al., 2023; Guo et al., 2023) and to guide Monte Carlo tree search (Zhao et al., 2023). Foundation models were also used in genetic operators (Lehman et al., 2022; Liu et al., 2023c; Meyerson et al., 2023; Chen et al., 2023), including multi-objective (Liu et al., 2023a) and quality-diversity algorithms (Nasir et al., 2023). These works are resource-intensive due to calling the model during the search (Liu et al., 2023a). This contrasts with our work, which uses the model a small number of times in a pre-processing step.

- **B** DEEP REINFORCEMENT LEARNING COMPARISON
- B.1 MICRORTS

To compare INNATECODER with a Deep Reinforcement Learning algorithm, we used the Gym- $\mu$ RTS Huang et al. (2021). We evaluated INNATECODER with PPO Gridnet self-play using an encodedecode model. We chose this model because is the closest one to ours, as they both learn through self play. The DRL agents proposed by Huang et al. were specifically designed and tested for a map of size 16 × 16. We used the Gym- $\mu$ RTS with the same settings presented in the repository, changing only the budget and the UnitTable used in the experiments. We trained both algorithms with a budget of 300 million steps in the MicroRTS.

790 We trained 15 DRL models for the BasesWorkers map  $(16 \times 16)$ , using a Xeon 2.90GHz, 64GB of 791 memory, and a dedicated Nvidia A10. Also, we performed 15 individual runs for INNATECODERFor 792 each pair (DRL-INNATECODER) we ran 10 matches of the policies using the same evaluation used by 793 Huang et al.. Each individual result is shown in Figure 6. Each run shows the number of losses, ties, 794 and wins that INNATECODERachieved against DRL. The best score obtained by DRL is 5 losses and 795 5 wins presented in the individual run number 5 of graphs. In contrast to the policies INNATECODER generates, which can be used to play at any of the two locations of the map, the DRL-PPO agent is 796 trained for a fixed position. This provides an advantage to DRL-PPO, as INNATECODER does not 797 specialize in a given location of the map. Figure 7 shows the average results, where the whiskers show 798 the 95% confidence interval. INNATECODER wins more than 70% of the matches with the PPO agent. 799 Moreover, while the DRL agent needs around 3 days to train for 300 million steps, INNATECODER 800 can perform the same number of training in less than 36 hours of computation using a single CPU. 801

802 803

804

B.2 KAREL

For Karel the Robot, we compared INNATECODER with Hierarchical Programmatic Reinforcement Learning (HPRL) (Liu et al., 2023b), which uses neural and programmatic representations of policies.
We evaluate the version of the system that uses PPO (HPRL-PPO), one of the best variants of HPRL. Table 2 shows the comparison between INNATECODER, HPRL-PPO, SHC (Carvalho et al., 2024), and CEM (Trivedi et al., 2021). These tests were performed with a budget of 10<sup>5</sup> episodes. INNATECODER (IC in the table) obtained the highest average return in all tasks.



Average Winrate InnateCoder versus DRL



Figure 7: Average results of INNATECODER against DRL in basesWorkers16x16A map, with 95% confidence interval.

# C CAN INNATECODER IMPROVE WITH MORE INFORMATION?

We evaluated whether INNATECODER's sample efficiency can scale with the amount of information we provide in the prompt used to generate the programmatic options. Figure 8 shows the results of two versions of INNATECODER: one that uses prompts with more information (IC+) and one that uses prompts with less information (IC-). The prompts are given in Section J.2 (more information) and Section J.3 (less information). The key difference between IC+ and IC- is that, in the former, we explain in the prompt how the agent can maximize its return. For example, in FourCorners we wrote "gent has to place one marker in each of the four corner cells of the grid". In contrast, in IC-we wrote "the robot will receive different reward values depending on its interactions inside the grid". Providing more information was never worse, and it was significantly better in two cases: Seeder and Snake. The ability to improve with more information is important because it allows the user of INNATECODER to achieve stronger results by crafting prompts that encode domain knowledge.

864	Task	HPRL-PPO	SHC	CEM	IC
865	Tubk	III KE I I O	bile	CLIM	10
866	StairClimberSparse	<b>1.000</b> ±0.00	<b>1.000</b> ±0.00	$0.601 \pm 0.44$	<b>1.000</b> ±0.00
867	MazeSparse	<b>1.000</b> ±0.00	<b>1.000</b> ±0.00	$0.097 \pm 0.03$	<b>1.000</b> ±0.00
868	TopOff	<b>1.000</b> ±0.00	<b>1.000</b> ±0.00	$0.812 \pm 0.29$	<b>1.000</b> ±0.00
000	FourCorners	1.000 ±0.00	1.000 ±0.00	$0.332 \pm 0.29$	1.000 ±0.00
869	Harvester	$0.924 \pm 0.13$	$0.906 \pm 0.07$	$0.487 \pm 0.17$	1.000 ±0.00
870	CleanHouse	$0.826 \pm 0.21$	$0.598 \pm 0.41$	$0.127 \pm 0.02$	1.000 ±0.00
871	DoorKey	$0.389 \pm 0.09$	$0.449 \pm 0.04$	$0.203 \pm 0.11$	<b>0.493</b> ±0.03
872	OneStroke	$0.784 \pm 0.11$	$0.908 \pm 0.01$	$0.683 \pm 0.23$	<b>0.932</b> ±0.01
873	Seeder	$0.539 \pm 0.17$	$0.779 \pm 0.10$	$0.339 \pm 0.14$	<b>0.931</b> ±0.09
874	Snake	$0.283 \pm 0.18$	$0.217 \pm 0.08$	$0.053 \pm 0.02$	<b>0.391</b> ±0.15
875					

Table 2: Mean and standard error of final episodic return of INNATECODER HPRL-PPO, SHC, and CEM in Karel problems.



Figure 8: Evaluating INNATECODER with more (IC+) and less (IC-) information provided in the prompt used to harness programmatic options from the foundation model. We used GPT 40 in this experiment. Average episodic return (maximum is 1.0 for all tasks) per number of episodes. The plots show average episodic return of 30 independent runs (seeds) and the 95% confidence interval.

D CAN INNATECODER IMPROVE WITH MODEL SIZE?



917 Figure 9: The plots show the average winning rate of 30 independent runs (seeds) and the 95% confidence interval.

```
918
            # InnateCoder's policy for the 64x64 map 1 # Llama's policy for the 64x64 map
         1
919
         2
             for (Unit u)
                                                       2
920
         3
                                                       3
                u.attackIfInRange()
         4
                for (Unit u)
                                                       4
921
         5
                                                       5
                    u.attackIfInRange()
922
         6
                                                       6
                    u.train(Heavy, EnemyDir, 50)
923
         7
                                                       7
                 u.harvest(10)
924
         8
                if (u.OpponentHasUnitInPlayerRange()) 8
         0
                                                       0
                    pass
925
         10
                                                      10
                 else
926
         11
                    u.train(Worker, EnemyDir, 3)
                                                      11
927
        12
                                                      12
                    u.train(Ranged, EnemyDir, 15)
928
        13
                                                      13
                u.attack(Strongest)
         14
                for (Unit u)
                                                      14
929
        15
                    u.harvest(5)
930
        16
                for (Unit u)
931
        17
                    u.build(Barracks, EnemyDir, 8)
932
        18
                u.moveToUnit(Enemy, Strongest)
        19
                if (u.HasUnitInOpponentRange())
933
        20
                    for (Unit u)
934
        21
                        u.moveToUnit(Enemv, Farthest)
935
        22
                         u.train(Light, EnemyDir, 6)
936
        23
                u.train(Worker, Down, 6)
```

```
for(Unit u)
 for(Unit u)
  u.train(Worker, Down, 5)
 for(Unit u)
  u.moveToUnit(Ally, Closest)
   u.harvest(5)
for(Unit u)
  u.build(Barracks, Up, 1)
  u.train(Ranged, EnemyDir, 10)
   u.attackIfInRange()
 for(Unit u)
  u.moveToUnit(Enemy, Weakest)
   u.attack(Weakest)
```

Figure 10: Left: A programmatic policy INNATECODER generated for the largest 64×64 map. This policy defeats the last three winners of the MicroRTS Competition: COAC, Mayari, and RAISocketAI. Right: one of the policies Llama 3.1 generated for the same map.

940 941 942

943

944

945

946

937 938

939

We evaluated GPT 3.5-turbo and GPT 40 on the MicroRTS tasks. Figure 9 shows the results, where we report the average winning rate and the 95% confidence intervals of 30 independent runs (seeds). Interestingly, we do not notice a significant change in winning rate as we move from the larger GPT 40 to the smaller GPT 3.5.

#### E **SAMPLES OF PROGRAMS**

951

952

953

954

955

956

957

958

Figure 10 shows an example of a programmatic policy INNATECODER generated for the BloodBath 64×64 map (left), and a program Llama 3.1 generated for the same map (right). The policy INNATE-CODER generated defeats the last three winners of the MicroRTS competition: COAC, Mayari, and RAISocketAI. This policy presents non-trivial features. For example, lines 7-11 will train Worker and Ranged units only if the player is not engaged in combat. This means that this policy focuses on economy and on Ranged units in the early stages of the game. Later in the game, the agent will save its resources to train Light units (line 21). Light units can be trained and move more quickly than Ranged units. While a Light unit is trained in 80 time steps of the game, a Ranged unit requires 100 time steps to be trained; also, a Light unit can move one cell every 8 time steps, while Ranged units move one cell every 10 time steps. Light units allow for a faster return of the resources invested.

959 The program the foundation model generates represents a weak policy (program shown on the right 960 of Figure 10). However, even weak policies can contain pieces of code—options—that can be helpful while searching for strong policies to play the game. For example, lines 2 and 3 offer a prioritization 961 scheme for investing resources to train Worker units because it will iterate over all units until it finds 962 the Base, which will be used to train up to 5 workers. Lines 2 and 3 are often found in strong policies. 963 For example, the program shown on the left of Figure 10 shows a similar structure in lines 1 and 10. 964

965

#### 966 967

#### **OPTION USAGE IN PROGRAMMATIC POLICIES** F

968 Figure 11 shows an example of a program INNATECODER generated as the final policy (the one that the system produces as output) for the LetMeOut map  $(16 \times 8)$ , where four options are used. The 969 colored lines represent the options. Lines 3 and 4 form one option (blue), while lines 8 and 9 (red), 11 970 and 12 (purple), and 13 and 14 (green) form the other three options. In this representative example, 8 971 of 14 lines are options, which represents 57% of the lines. We also analyzed 10 independent runs of

972	1	for(Unit u)
973	2	u.harvest(1)
974	3	if(HasNumberOfUnits(Barracks,1))
975	4	u.train(Ranged,Down,10)
976	5	for(Unit u)
977	6	for(Unit u)
078	7	u.moveToUnit(Ally,MostHealthy)
070	8	for(Unit u)
979	9	u.attack(Closest)
980	10	for(Unit u)
981	11	for(Unit u)
982	12	u.train(Worker,EnemyDir,3)
983	13	for(Unit u)
984	14	u.moveToUnit(Ally,Weakest)

Figure 11: Programmatic policy INNATECODER generated for the LetMeOut map ( $16 \times 8$ ). The colored lines represent options from one execution of INNATECODER that used 5000 options.

990 INNATECODER in the  $9 \times 8$  map with an initial option set extracted from 120 policies, which were 991 generated with the Llama 3.1 model. We found that on average, 63% of the options in the initial 992 set are used in a best response during the self-play training process. Recall that a best response is 993 the solution INNATECODER finds to a given MDP within the self-play algorithm. The maximum 994 percentage of options used in the 10 runs was 73% and the minimum was 46%. Even if an option is 995 not used in the policy INNATECODER outputs, the option can still have played an important role in 996 allowing INNATECODER arrive at the output it produced. This is because an option could have been part of one of the best responses encountered during self play, and these best responses provide the 997 signal needed to guide the search toward stronger policies for playing the game (Moraes et al., 2023). 998

1000 The policy shown in Figure 11 also provides an explanation for the difference in performance 1001 between the versions of INNATECODER that are initialized with a pool of 1400 or fewer options 1002 and the versions that are initialized with a pool of 5000 or more options (see Figure 5). The options 1003 highlighted in Figure 11 are in the set of 5000 options of an INNATECODER run, but they are not in the set of 1400 options of another run of the algorithm. These options are clearly helpful because they 1004 appear in the output policy. Sampling 5000 options instead of 1400 increases the chances of adding 1005 some of these helpful options to the library, thus explaining the gap we see between INNATECODER 1006 with 1400 options or fewer and INNATECODER with 5000 options or more in Figure 5. 1007

1008

999

985 986

987

988 989

#### 1009

# 1010 G VARYING THE VALUE OF $\epsilon$

1011 1012

Figure 12 shows the performance of INNATECODER as we vary the value of  $\epsilon$ . Recall that  $\epsilon$  dictates 1013 how much of the search is done in the syntax space and how much of the search is done in the 1014 semantic space. Also, recall that smaller values of  $\epsilon$  mean that the semantic space is used more often. 1015 There are two groups of lines at 80000 games played:  $\epsilon$ -values of 0.3, 0.4, 0.5, and 0.6 (top lines) 1016 and 0.1, 0.2, and 0.7 (bottom lines). These results suggest that INNATECODER is robust to the choice 1017 of  $\epsilon$  value, as a wide range of values produce competitive results among themselves. The results also 1018 suggest that if the value of  $\epsilon$  is too small, then the search does not sufficiently explore the syntax space to eventually find programs that were not originally in the library of options. If the value of  $\epsilon$  is 1019 too large, then the search is not making use of the helpful options in the library as often as it could. 1020

1021 1022

#### H DOMAIN-SPECIFIC LANGUAGES (DSLS)

1023 1024 1025

In this section, we present the DSLs used in our experiments for both Karel the Robot and MicroRTS.



of the language does not have. Considering these key differences, it is unlikely that the data from
 Mariño et al. (2021) have influenced our results more than the vast collection of programs written in
 various programming languages that are present in the corpus used to train these foundation models.

1080	
1081	$S \rightarrow SS$ for (Unit u) S if (B) then S
1082	$\int \int \partial f(D) + her C = 0 = 0 = C = C = 0$
1083	II(B)  then  S  erse  S   C   X
1084	$B  ightarrow has { m Number Of Units}(T,N) \mid { m opponent Has Number Of Units}(T,N)$
1085	$\mid$ hasLessNumberOfUnits $(T,N)\mid$ haveQtdUnitsAttacking $(N)$
1086	hasUnitWithinDistanceFromOpponent $(N)$
1087	hasNumberOfWorkersHarvesting $(N)$
1088	is Type $(T)$   is Builder()
1089	1011 growth = 0 + 10 + 10 + 10 + 10 + 10 + 10 + 10
1090	
1091	opponentHasUnitThatKillsUnitInOneAttack()
1092	hasUnitInOpponentRange()
1093	opponentHasUnitInPlayerRange()
1094	canHarvest()
1095	$C \rightarrow \text{build}(T, D, N) \mid \text{train}(T, D, N) \mid \text{moveToUnit}(T_n, O_n) \mid \text{attack}(O_n)$
1096	harvest(N)   attackIfInRange()   moveAway()
1097	$T \rightarrow \text{Desc}   $
1098	$I \rightarrow Base   Barracks   Ranged   Heavy$
1099	Light Worker
1100	$N \to 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
1101	10   15   20   25   50   100
1102	$D \rightarrow \text{EnemvDir}   Up   Down   Right   Left$
1103	$O \rightarrow \text{Strongost} \mid \text{Weakest} \mid C \mid \text{esect} \mid \text{Farthest}$
1104	$O_p \rightarrow \text{Strongest}$   weakest   crosest   ratchest
1105	LessHealthy   MostHealthy   Random
1106	$T_p \to \operatorname{Ally} \mid \operatorname{Enemy}$
1107	

1115

1117

#### 

The KAREL problem sets (Trivedi et al., 2021; Liu et al., 2023b) are divided into two parts— KAREL
and KAREL-HARD. KAREL consists of six different tasks, while KAREL-HARD includes four
additional tasks that are comparatively more difficult to solve. In this section, we describe the initial
state and the return function of each task in both KAREL and KAREL-HARD problem sets.

1116 I.1 KAREL

**StairClimber.** In this task, the agent operates within a  $12 \times 12$  grid containing stairs formed by walls. The goal for the agent is to reach a marker located above its position on the stairs. The initial positions of the agent and the marker are randomly initialized on the stairs. The agent receives an episodic return of 1 if it successfully reaches the marker, and 0 otherwise. Moving to a position outside the contour of the stairs results in a return of -1.

FourCorners. In this task, the goal for the agent is to place a marker in the four corners of a  $12 \times 12$ grid. The initial position of the agent is randomly initialized near the wall. The return is calculated as the number of corners with one marker divided by four.

1126

1127**TopOff.** In this task, the agent is always initialized on the bottom left of a  $12 \times 12$  grid, and the1128markers are placed randomly on the bottom row of the grid. The goal for the agent to place markers1129on top of the markers present in that row. The return is calculated as the number of markers topped1130off divided by the number of markers present in the grid.

1131

**Maze.** A maze, formed by walls, is randomly configured on a  $12 \times 12$  grid, and a marker and an agent are randomly placed within it. The goal for the agent is to reach the marker, that yields an episodic return of 1. Otherwise, the episodic return is 0.

1134 1135 1136 1137 **CleanHouse.** Markers and walls are randomly placed in a  $22 \times 14$  grid, referred to as the apartment. The position of the agent is also initialized randomly. Its goal is to pick all the markers inside the apartment. The return is calculated as the number of picked markers divided by the total number of markers present.

**Harvester.** In this task, the agent operates within an  $8 \times 8$  grid filled with markers. The agent is placed randomly on the bottom row, with the goal of collecting all markers in the grid. The return is calculated as the number of collected markers divided by the total number of cells in the grid.

1143

1145

1138

1144 I.2 KAREL-HARD

**1146 DoorKey.** The environment, consisting of an  $8 \times 8$  grid is divided into two chambers by a vertical 1147 bar with a door. The position of the agent is initialized randomly, and a marker is placed randomly in 1148 each chamber. The goal for the agent is to pick the marker in the left chamber, which will open the 1149 door in the vertical bar, and then pick the marker in the right chamber. Picking each marker results in 1150 a return of 0.5.

1151

1155

1152**OneStroke.** In this task, the goal for the agent is to visit every cell of an  $8 \times 8$  grid without repetition.1153Once a cell is visited, it transforms into a wall, and the episode terminates if the agent hits a wall.1154The return is calculated as the number of visited cells divided by the total number of cells in the grid.

**Seeder.** The environment is given by an  $8 \times 8$  empty grid with an agent initialized at a random position. The goal for the agent is to place a marker in each cell of the grid. The return is calculated as the number of placed markers divided by the total number of cells.

1159 1160

1161Snake. In this task, an  $8 \times 8$  grid is initialized with an agent and a marker at random positions.1162The agent behaves like the head of a snake, with its body growing after collecting a marker. Once a<br/>marker is collected, another marker is placed at a random position. This process continues until the<br/>agent collects 20 markers. The goal for the agent is to collect the markers without hitting its own<br/>body. The return is calculated as the number of collected markers divided by 20.

- 1166
- <sup>1167</sup> J KAREL PROMPTS
- 1168 1169

In this section, we present the prompts used to generate a program that encodes a policy. These are the prompts we used multiple times to obtain a set of programmatic policies. This section is divided into three subsections. Section J.1 includes the complete prompt used to obtain a program from the model to solve the 'Seeder' task. For each task, we provide two types of prompts: one with more information and another with less information. Section J.2 contains the details of the environment with more information for each task, while Section J.3 contains the prompts with less information for each task.

1176 1177

#### 1178 J.1 COMPLETE PROMPT FOR THE SEEDER TASK

The prompt for obtaining a program that encodes a policy for the 'Seeder' task is shown below. Note that the first paragraph of the prompt explains the environment. Only this paragraph is changed from one task to the next.

1183 1184

Consider an 8x8 gridded 2D environment in Karel the Robot, where an agent has to place
one marker in every cell of the grid. The grid is initially empty and the initial position of the
agent is randomly assigned at the beginning of each episode.

The following is the Context Free Grammar (CFG) for the Karel domain:
$P \rightarrow \min\{\mathbf{S}\}$
$C \rightarrow W(H) = (D) (C)   C C   A   D = D = AT D = N (C)   E (D) (C)   E (D) (C) = C (C)$
$D \rightarrow W_{\text{HEE}}(D) \{0\} \{0\} \{0\} \{1\} \{1\} \{1\} \{1\} \{1\} \{1\} \{1\} \{1\} \{1\} \{1$
$B \to H \mid \text{norm}$
$H \rightarrow$ frontlsClear   leftlsClear   rightlsClear   markersPresent   noMarkersPresent
$A \rightarrow \text{move} \mid \text{turnLeft} \mid \text{turnRight} \mid \text{putMarker} \mid \text{pickMarker}$
$N \rightarrow 1 \mid 2 \mid 3 \mid \mid \text{infinity}$
The CFG is explained below in the "CFG Explanation" section:
CFG Explanation:
P: The main program named as "run" which contains Statement S
S: Consists of statements such as WHILE, IF, IFELSE, REPEAI. Can have multiple statements "S S" or action "A"
B: Percention H or not perception H that returns true or false
H: Some boolean variables that provide the idea of the environment with true or false
A: Action to be taken by the agent
N: A positive integer that indicates the number of repetitions
frontIsClear: Checks if the next cell towards the direction the agent is facing is inside the grid
leftIsClear: Checks if the next cell to the left of the direction the agent is facing is inside the
grid
rightIsClear: Checks if the next cell to the right of the direction the agent is facing is inside
the grid move: The agent moves towards its front
turn left: The agent turns left
turnRight: The agent turns right
putMarker: The agent puts a marker in the current cell
pickMarker: The agent picks a marker from the current cell
: It is not part of the CFG. It has been used to indicate all positive numbers in between.
To write a program from the given CFG, the following "Program Writing Guidelines" must
be followed:
Program Writing Guidelines:
1. In the CFG, 'infinity' means that the value of N can be up to infinity. So do not write any
part of the program like "R=infinity"
2. DO NOT write "" in the program, since it is not a part of the CFG
3. The program must start with "run"
now your tasks are the following 3:
1. Read carefully about the details of the environment, the CFG and its explanation.
2. Follow the CFG and the program writing guidelines and write a program of maximum 8
lines that will gain the maximum reward inside this given environment.
3. Write the program inside iprogram <sup>2</sup> <sub>i</sub> /program <sup>2</sup> <sub>i</sub> tag.

#### 1242 J.2 ENVIRONMENT DETAILS WITH MORE INFORMATION 1243

#### 1244 **STAIRCLIMBER**

1245 1246

1247

1248

1249

1250 1251

1252 1253 1254

1255

1256

1257

1259

1260 1261

Consider a 12x12 gridded 2D environment in Karel the Robot, where an agent has to reach a marker by climbing up along a stair. The grid contains a stair-like structure and the initial position of the agent is randomly assigned near the stair with the marker placed at the higher end, at the beginning of each episode.

#### FOURCORNERS

Consider a 12x12 gridded 2D environment in Karel the Robot, where an agent has to place one marker in each of the four corner cells of the grid. The grid is initially empty and the initial position of the agent is randomly assigned near the wall of the grid at the beginning of each episode.

#### TOPOFF

Consider a 12x12 gridded 2D environment in Karel the Robot, where an agent has to place markers on top of other cells that already have markers, and then reach at the rightmost cell of the bottom row. The markers are initialized randomly at the bottom row. The position of the agent is fixed at the leftmost cell of the bottom row at the beginning of each episode.

#### Maze

Consider a 12x12 gridded 2D environment in Karel the Robot, where an agent has to pick a marker following a path surrounded by walls. Some cells of the grid are filled with walls, collectively referred to as the maze. The agent, the marker and the maze are randomly placed at the beginning of each episode.

#### **CLEANHOUSE**

Consider a 22x14 gridded 2D environment in Karel the Robot, where an agent has to pick some markers that are placed randomly inide the grid. There are also some cells filled with obstacles. The position of the agent is fixed whereas, the markers are initialized at random cells at the beginning of each episode.

#### HARVESTER

Consider an 8x8 gridded 2D environment in Karel the Robot, where an agent has to pick one marker from every cell of the grid. The grid is initially filled with markers and the initial position of the agent is randomly assigned at the beginning of each episode.

#### DOORKEY

Consider an 8x8 gridded 2D environment divided into two chambers by walls in Karel the Robot, where an agent has to pick a marker from the left chamber and put it over another marker placed at the right chamber. The agent cannot access the right chamber without picking the marker from the left chamber. The position of the agent, the marker of the left

- 1286 1287
- 1290
- 1291
- 1293 1294 1295

1	2	9	6
1	2	9	7
1	2	9	8

1302 1303

1304

1305

1306 1307

1309

1310 1311

1312

1313

1314 1315 1316

1317 1318 1319

1320

1321

chamber and the marker of the right chamber are initialized randomly at the beginning of each episode.

#### **ONESTROKE**

Consider an 8x8 gridded 2D environment in Karel the Robot, where an agent has to visit as many cells as possible in one attempt. Once the agent visits a cell, it will be filled with a wall. The position of the agent is initialized randomly at the beginning of each episode.

#### SEEDER

Consider an 8x8 gridded 2D environment in Karel the Robot, where an agent has to place one marker in every cell of the grid. The grid is initially empty and the initial position of the agent is randomly assigned at the beginning of each episode.

#### SNAKE

Consider an 8x8 gridded 2D environment in Karel the Robot, where an agent has to pick a marker multiple times inside the grid. The position of the marker will be changed once the agent picks the marker. Each time the agent picks a marker, it will be attached to its body. For instance, if the agent picks 3 markers, it will have 3 markers added to its back. The agent has to pick as many markers as possible without hitting the markers attached to its body. The position of the agent and the marker are initialized randomly at the beginning of each episode.

#### J.3 ENVIRONMENT DETAILS WITH LESS INFORMATION

#### **STAIRCLIMBER**

Consider a 2D 12x12 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, there is a stair-like structure and the agent is placed near the stair. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### FOURCORNERS

Consider a 2D 12x12 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, no markers are present in any cell of the grid. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

1341 1342

1344

1345

1347

1348

#### 1350 TOPOFF

Consider a 2D 12x12 grid with an agent placed at a fixed cell at the leftmost cell of the bottom row, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, there are markers present in some cells at the bottom row of the grid. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### MAZE

Consider a 2D 12x12 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, there are few walls and a marker inside the grid. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### CLEANHOUSE

Consider a 2D 22x14 grid with an agent placed at a fixed cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, there are some markers and obstacles randomly placed inside the grid. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### HARVESTER

Consider a 2D 8x8 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, a marker is present in each cell of the grid.

The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

## DOORKEY

1396	
1397	Consider a 2D 8x8 grid divided into two chambers with an agent placed randomly in the left
1398	chamber, which will interact within the grid. Each cell may contain a marker, be empty, or be
1399	blocked by a wall. Initially, for this particular problem, there is one marker in each chamber
1/00	of the grid.
1401	The robot will receive different reward values depending on its interactions inside the grid.
1401	The interaction of the robot will be decided by a program written in the domain-specific
1402	language to be provided below as a context-free grammar (CFG). The goal is to generate a
1403	program that will maximize the sum of rewards the robot obtains by following that program.

# 1404 ONESTROKE

Consider a 2D 8x8 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, no markers are present in any cell of the grid. Cells that will be visited by the agent will be filled with obstacles. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### SEEDER

Consider a 2D 8x8 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, no markers are present in any cell of the grid. The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific.

The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a program that will maximize the sum of rewards the robot obtains by following that program.

#### SNAKE

Consider a 2D 8x8 grid with an agent placed randomly in any cell, which will interact within the grid. Each cell may contain a marker, be empty, or be blocked by a wall. Initially, for this particular problem, there is one marker inside the grid. The position of the marker changes depending on a certain behaviour of the agent.
The robot will receive different reward values depending on its interactions inside the grid. The interaction of the robot will be decided by a program written in the domain-specific language to be provided below as a context-free grammar (CFG). The goal is to generate a

program that will maximize the sum of rewards the robot obtains by following that program.

## K MICRORTS MAPS

Figure 13 shows the three MicroRTS maps used in our experiments. In these maps, the geometric shapes with blue borders represent the units of one player, while those bounded by red borders represent the units of the other player. Neutral objects do not have colored borders (e.g., resources in light green and walls in dark green).





# <sup>1512</sup> L MICRORTS PROMPTS

In this section, we present the prompt used to obtain one program that encodes a policy for MI-CRORTS. This is the prompt that we use multiple times to obtain a set of programmatic policies.
Section L.1 shows the complete prompt for the 'NoWhereToRun (9x8)' map, as a complete example.
Section L.2 shows the description of the environment used in the prompts of each map.

#### 1519 L.1 COMPLETE PROMPT FOR THE NOWHERETORUN (9x8) MAP 1520

The prompt for obtaining a program that encodes policies for the 'NoWhereToRun (9x8)' map is shown below. Note that the first paragraph of the prompt mentions the details of the environment for a given map. Only this paragraph describing the environment is updated to obtain the program for each separate map. We use a subset of the MicroLanguage, where conditionals are removed. Our goal was to have a language that would be easier for the foundation model to use.

1526	
1527 1528	2-dimensional array with the following structure:
1529	– There are a total of 8 neutral resource cells situated along the central column of
1530	the map, dividing the map into two parts. Each resource cell contains 10 units of
1531	resources.
1532	- The base B1 of player 1 is located at index $(1,1)$ , which is located on the left side of
1533	the map.
1534 1535	<ul> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> </ul>
1536	- Each player controls one base, which initially has 5 units of resources.
1537	- The only unit a player controls at the beginning of the game is the base.
1538	Consider this Context Free Grammar (CEG) describing a programming language for writing
1539	programs encoding strategies of microRTS. The CFG is shown in the $\langle CFG \rangle \langle /CFG \rangle$
1540	tag bellow:
1541	
1543	$\langle CFG \rangle$
1544	$S \rightarrow SS \mid \text{for(Unit u)} S \mid C \mid \lambda$
1545	$C \rightarrow u \text{ build}(T, D, N) \mid u \text{ train}(T, D, N) \mid u \text{ moveToUnit}(T, O)$
1546	$C \rightarrow u.$ bund $(I, D, N) \mid u.$ train $(I, D, N) \mid u.$ into the root int $(I_p, O_p)$
1547	$ u.attack(O_p)  u.narvest(N)  u.attackfiftRange()  u.noveAway()$
1548	$I \rightarrow Base   Barracks   Ranged   Heavy$
1549	Light   Worker
1550	$N  o 0 \mid 1 \mid 2 \mid 3 \mid 4 \mid 5 \mid 6 \mid 7 \mid 8 \mid 9$
1551	10   15   20   25   50   100
1552	$D \rightarrow \text{EnemyDir} \mid \text{Up} \mid \text{Down} \mid \text{Right} \mid \text{Left}$
1554	$O_p \rightarrow \text{Strongest} \mid \text{Weakest} \mid \text{Closest} \mid \text{Farthest}$
1555	LessHealthy   MostHealthy   Random
1556	$T_p \rightarrow \text{Ally} \mid \text{Enemy}$
1557	
1558	This language allows posted loops. It contains several command emiented functions $(C)$
1559	The Command functions ('C' in the CEG) are described below:
1560	1. u build(T. D. N): Ruilde N units of type T on a call logated in the D direction of the
1561	unit. The u build function is used to build Barracks and Base
1562	2 u troin(T D N): Trains N units of type T on a call located in the D direction of the
156/	2. u.u.am(1, D, 1v). Trains IV units of type 1 of a cen located in the D direction of the structure responsible for training them. For example, the instruction u train(Heavy
1565	Down, 1) will allow the agent to train at most 1 heavy unit in the down direction of

	the Barrack, while the instruction u.train(Heavy, EnemyDir, 20) will allow to train
	at most 20 towards the direction of the opponent. The number used in the function
	calls could play a big role in the strategy the program encodes. The ultrain function is used to train Worker, Panged Light, and Heavy units
	is used to train worker, Ranged, Light, and reavy units.
	ing a criterion O_p.
	4. u.attack(O_p): Sends N Worker units to harvest resources.
	5. u.harvest(N): Sends N Worker units to harvest resources. For example, u.harvest(5) will send 5 workers to harvest resources.
	6. u.attackIfInRange(): Commands a unit to stay idle and attack if an opponent unit comes within its attack range.
	7. u.moveAway(): Commands a unit to move in the opposite direction of the player's base.
'T' re	presents the types a unit can assume.
'D' re	presents the directions available used in action functions.
'О_р' 'Т р'	is a set of criteria to select an opponent unit based on their current state.
'N' is	the number of units that can be any integers from 0 to 10, or 15, or 20, or 25, or 50, or
100.	
Thef	allowing 4 are some guidelines for writing the playing strategy.
I ne I	Showing 4 are some guidelines for writing the playing strategy:
	1. There is NO NEED TO write classes, or initiate objects such as Unit, Worker, etc. There is also NO NEED TO write comments
	2 Use curly braces like C/C++/Java while writing any 'for' block. Start the curly
	braces in the same line of the block.
	3. A strategy must be written inside one or multiple 'for' blocks.
	4. This language does not have if-statements.
Now	your tasks are the following 8:
	1. Understand the command (C) functions from above and try to relate them in the context of playing strategies for a real-time strategy game.
	2. Write a program in the microRTS language encoding a very strong game-playing strategy for the map described above. You must follow the guidelines of writing the playing strategy while writing your program.
	3. You must not use any symbols (for example &&,   , etc.) that the CFG does not accept. You have to strictly follow the CFG while writing the program.
	4. Look carefully, the methods of non-terminal symbols C have prefixes 'u.' in the examples since they are methods of the object 'Unit u'. You should follow the patterns of the examples.
	5. Write only the pseudocode inside ' $< strateau > < /strateau >$ ' tag.
	6. Do not write unnecessary symbols of the CFG such as, $(S \rightarrow)$ , $(\rightarrow)$ , etc.
	7. Check the program and ensure it does not violate the rules of the CFG or the
	guidelines for writing the strategy.
	8. The for loops in this language iterate over all units and the instructions inside the for

# 1620 L.2 ENVIRONMENT DETAILS

# 1622 NOWHERETORUN (9x8)1623

<ul> <li>Consider a 9x8 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There are a total of 8 neutral resource cells situated along the central column of the map, dividing the map into two parts. Each resource cell contains 10 units of resources.</li> <li>The base B1 of player 1 is located at index (1,1), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> IBLEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (2,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>The tare a 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each player 1 is located at index (5,1</li></ul>	Consid 2-dim	
<ul> <li>There are a total of 8 neutral resource cells situated along the central column of the map, dividing the map into two parts. Each resource cell contains 10 units of resources.</li> <li>The base B1 of player 1 is located at index (1,1), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> <b>JBLEGAME</b> (24x24) <b>Consider a</b> 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 12-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (20,21) and (3,21), also located on both sides of the wall. <ul> <li>Each player controls two bases, which initially have 5 units of resource each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> <b>DISTANTRESOURCES</b> (32x32) <b>Consider a</b> 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map.</li> </ul> <b>Consider a</b> 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map.</li> <li>There are total of 12 neutral resource cells R loc</li></ul></li></ul>		der a 9x8 gridded map of microRTS, a real-time strategy game. Consider this map as a
<ul> <li>There are a total of a neutral resource cells stuated along the central column of the map, dividing the map into two parts. Each resource cell contains 10 units of resources.</li> <li>The base B1 of player 1 is located at index (1,1), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> <b>IBLEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game.</b> Consider this map as 12-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources. <ul> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> <b>DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game.</b> Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map.</li> <li>Th</li></ul></li></ul>	- 4111	The second straig with the following structure:
<ul> <li>The base B1 of player 1 is located at index (1,1), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> /// IEEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as 12-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 2 are located at indices (20,21) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>There are a total of 12 neutral resource</li></ul></li></ul>		- There are a total of 8 neutral resource cells situated along the central column of the map dividing the map into two parts. Each resource cell contains 10 units of
<ul> <li>The base B1 of player 1 is located at index (1,1), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> IBLEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (20,21) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. <ul> <li>The base B1 of player 1 is located at index (6,14), which is loc</li></ul></li></ul>		resources.
<ul> <li>the map.</li> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> <li>//BLEGAME (24x24)</li> <li>Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as 1 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 28 vorkers beside each base. So a total of 4 workers for each of the players.</li> </ul> </li> <li>DISTANTRESOURCES (32x32)</li> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul> </li> </ul>		- The base B1 of player 1 is located at index $(1,1)$ , which is located on the left side of
<ul> <li>The base B2 of player 2 is located at index (7,6), which is located on the right side of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> IBLEGAME (24X24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as 12-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources. <ul> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the resources. <ul> <li>The base B1 of playe</li></ul></li></ul></li></ul>		the map.
<ul> <li>of the map.</li> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> JBLEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources. <ul> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		- The base B2 of player 2 is located at index (7,6), which is located on the right side
<ul> <li>Each player controls one base, which initially has 5 units of resources.</li> <li>The only unit a player controls at the beginning of the game is the base.</li> </ul> JBLEGAME (24x24) Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left corners of the map. Each resource cells R located at the top-right and bottom-left side of the map.</li> <li>The base B1 of player 2 is located at index (25,17), which is</li></ul>		of the map.
<ul> <li>The only unit a player controls at the beginning of the game is the base.</li> <li>JBLEGAME (24x24)</li> <li>Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> </li> <li>DISTANTRESOURCES (32x32)</li> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul> </li> </ul>		- Each player controls one base, which initially has 5 units of resources.
<ul> <li>JBLEGAME (24x24)</li> <li>Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> </li> <li>DISTANTRESOURCES (32x32)</li> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul> </li> </ul>		- The only unit a player controls at the beginning of the game is the base.
<ul> <li>Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the right side of the map.</li> <li>Each player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player 2 is located at index (25,17), which is located on the right corners of the map.</li> <li>Each player 2 is located at index (b, 14), which is located on the right corners.</li> <li>The base B2 of player 2 is located at index (b, 14), which is located on the right side of the map.</li> </ul>	DUBLE	Game (24x24)
<ul> <li>Consider a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. <ul> <li>The base B1 of player 1 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		
<ul> <li>There is a wall in the middle of the map consisting of two columns that has a small passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (2,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>	Consid a 2-di	der a 24x24 gridded map of microRTS, a real-time strategy game. Consider this map as mensional array with the following structure:
<ul> <li>passage of 4 cells. The small passage consists of 4 resource cells each having only 1 resource.</li> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32X32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		- There is a wall in the middle of the map consisting of two columns that has a small
<ul> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		passage of 4 cells. The small passage consists of 4 resource cells each having only 1
<ul> <li>There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right corners of the map respectively where each of them contains 10 units of resources.</li> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>Each player 2 is located at index (25,17), which is located on the right side of the map.</li> </ul></li></ul>		resource.
<ul> <li>The bases of player 1 are located at indices (3,2) and (20,2), located on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the right side of the map.</li> <li>Each player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player 2 is located at index side side side side of the map.</li> </ul></li></ul>		- There are 28 resource cells at the top-left, top-right, bottom-left and bottom-right
<ul> <li>The bases of player 1 are focated at indices (3,2) and (20,2), focated on both sides of the wall.</li> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul></li></ul>		The bases of player 1 are located at indices $(3, 2)$ and $(20, 2)$ located on both sides of
<ul> <li>The bases of player 2 are located at indices (20,21) and (3,21), also located on both sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul></li></ul>		the wall.
<ul> <li>sides of the wall.</li> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		- The bases of player 2 are located at indices (20.21) and (3.21), also located on both
<ul> <li>Each player controls two bases, which initially have 5 units of resources each.</li> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> </ul> DISTANTRESOURCES (32x32) Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources. <ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul></li></ul>		sides of the wall.
<ul> <li>There are 2 workers beside each base. So a total of 4 workers for each of the players.</li> <li>DISTANTRESOURCES (32X32)</li> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>		- Each player controls two bases, which initially have 5 units of resources each.
<ul> <li>DISTANTRESOURCES (32x32)</li> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul> </li> </ul>		- There are 2 workers beside each base. So a total of 4 workers for each of the players.
<ul> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>		1 2
<ul> <li>Consider a 32x32 map of microRTS, a real-time strategy game. Consider this map as a 2-dimensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>		NTRESOURCES (20x22)
<ul> <li>2-dimensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	<b>VDIST</b>	ANTRESOURCES (32x32)
<ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	VDIST/ Consi	ANTRESOURCES (32x32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a
<ul> <li>at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	DISTA Consi 2-dim	ANTRESOURCES (32x32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure:
<ul> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	/DIST/ Consi 2-dim	ANTRESOURCES (32x32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure: – There are two L-shaped obstacles on the map, each with a passage of 4 cells located
<ul> <li>corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	/DIST/ Consi 2-dim	ANTRESOURCES (32x32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure: - There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.
<ul> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	/DIST/ Consi 2-dim	ANTRESOURCES (32X32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure: – There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides. – There are a total of 12 neutral resource cells R located at the top-right and bottom-left
<ul> <li>of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	VDIST/ Consi 2-dim	ANTRESOURCES (32x32) der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure: – There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides. – There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.
<ul> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	/DIST/ Consi 2-dim	<ul> <li>ANTRESOURCES (32x32)</li> <li>der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side</li> </ul>
<ul> <li>side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	VDISTA Consi 2-dim	<ul> <li>ANTRESOURCES (32X32)</li> <li>der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> </ul>
<ul> <li>Each player controls one Base, which initially has 20 units of resources.</li> <li>There is one worker for each player besides their bases.</li> </ul>	/DIST/ Consi 2-dim	<ul> <li>ANTRESOURCES (32x32)</li> <li>der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right interval.</li> </ul>
<ul> <li>There is one worker for each player besides their bases.</li> </ul>	DIST/ Consi 2-dim	<ul> <li>ANTRESOURCES (32x32)</li> <li>der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure: <ul> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> </ul> </li> </ul>
	/DIST/ Consi/ 2-dim	<ul> <li>ANTRESOURCES (32X32)</li> <li>der a 32x32 map of microRTS, a real-time strategy game. Consider this map as a ensional array with the following structure:</li> <li>There are two L-shaped obstacles on the map, each with a passage of 4 cells located at the middle of left and right sides.</li> <li>There are a total of 12 neutral resource cells R located at the top-right and bottom-left corners of the map. Each resource center contains 20 units of resources.</li> <li>The base B1 of player 1 is located at index (6,14), which is located on the left side of the map.</li> <li>The base B2 of player 2 is located at index (25,17), which is located on the right side of the map.</li> <li>Each player controls one Base, which initially has 20 units of resources.</li> </ul>

#### 1674 BLOODBATH (64x64) 1675

6	
	Consider a 64x64 gridded map of microRTS, a real-time strategy game. Consider this map as
	<ul> <li>There are total 4 neutral resource cells situated close to the top-left, top-right, bottom-left and bottom-right sides of the map respectively. Each resource cell contains 40</li> </ul>
	units of resources.
	- There are obstacles in between each of the 4 resource centers.
	- The base B1 of player 1 is located at index (53, 55), which is located on the bottom-
	right side of the map.
	- The base B2 of player 2 is located at index (2, 6), which is located on the top-left side of the map.
	- Each player controls one Base each, which initially has 5 units of resources.
	– There is no worker for both player 1 and 2 in the initial map setup.
	- The only unit a player controls at the beginning of the game is the Base.
М	Self-Play Learning Algorithms
Self- Resp IBR $\pi_1$ to proc The	play algorithms attempt to approximate an optimal policy for two-player games. Iterated Bestonse (IBR) (Lanctot et al., 2017) is perhaps the simplest self-play algorithm that we could use starts with an arbitrary policy $\pi_0$ in $[D]$ for one of the players and approximates a best response $\sigma_0$ for the other player. Then it approximates a best response $\pi_2$ to $\pi_1$ for the first player. This ess is repeated a number of times $m$ , which is normally determined by a computational budge last resulting policy $\pi_m$ is returned as IBR's approximate optimal policy for the game.
The ons FP) Il pi amo	self-play process IBR follows generates a sequence of policies for the players, but IBR only iders the latest policy while computing a best response. Other algorithms, such as Fictitious Play (Brown, 1951), compute best responses to a policy that mixes the best responses computed in revious iterations. The use of more policies allows the method to find optimal policies even in es with cyclic dynamics such as Rock, Paper, and Scissors.
The nult nste een redu or p	algorithm we use in our experiments, Local Learner (2L) (Moraes et al., 2023), also consider iple policies in its self-play loop. This allows 2L to use more information than IBR. However ad of including all policies seen in the process, like FP does, it selects a subset of the policies in the process. Using all policies can be computationally wasteful, as many of the policies are indant". 2L selects which policies to use based on the information collected during the search programmatic policies. We refer the reader to the work of Moraes et al. (2023) for a detailed anation of 2L.
1	