

---

# Hierarchical Decision Making by Generating and Following Natural Language Instructions

---

Anonymous Author(s)

Affiliation

Address

email

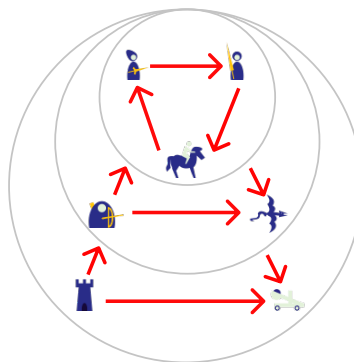
## A Detailed game design

We develop an RTS game based on the MiniRTS framework, aspiring to make it intuitive for humans, while still providing a significant challenge to machines due to extremely high-dimensional observation and actions spaces, partial observability, and non-stationary environment dynamics imposed by the opponent. Below we describe the key game concepts.

### A.1 Game units specifications

**Building units** Our game supports 6 different building types, each implementing a particular function in the game. Any building unit can be constructed by the PEASANT unit type at any available map location by spending a specified amount of resources. Later, the constructed building can be used to construct units. Most of the building types can produce up to one different unit type, except of WORKSHOP, which can produce 3 different unit types. This property of the WORKSHOP building allows various strategies involving bluffing. A full list of available building units can be found in Table 1.

**Army units** The game provides a player with 7 army unit types, each having different strengths and weaknesses. PEASANT is the only unit type that can construct building units and mine resources, so it is essential for advancing to the later stages of the game. We design the attack relationships between each unit type with a *rock-paper-scissors* dynamic—meaning that each unit type has another unit type that it is effective against or vulnerable to. This property means that effective agents must be reactive to their opponent’s strategy. See Fig. 1 for a visualization. Descriptions of army units can be found in Table 2.



**Resource unit** RESOURCE is a stationary and neutral unit type, it cannot be constructed by anyone, and is only created during the map generation phase. PEASANTS of both teams are allowed to mine the same RESOURCE unit, until it is exhausted. Initial capacity is set to 500, and one mine action subtracts 10 points from the RESOURCE. Several RESOURCE units are placed randomly on the map, which gives rise to many strategies around RESOURCE domination.

Figure 1: Our game implements the *rock-paper-scissors* attack graph, where each unit has some units it is effective against and vulnerable to.

## 35 A.2 Game map

36 We represent the game map as a discrete grid of 32x32. Each cell of the grid can either be grass  
37 or water, where the grass cell is passable for any army units, while the water cell prevents all units  
38 except of DRAGON to go through. Having water cells around one’s main base can be leveraged as a  
39 natural protection. We generate maps randomly for each new game, we first place one TOWN HALL  
40 for each player randomly. We then add some water cells onto the map, making sure that there is  
41 at least one path between two opposing TOWN HALLs, but otherwise aiming to create bottlenecks.  
42 Finally, we randomly locate several RESOURCE units onto the map such that they are approximately  
43 equidistant from the players TOWN HALLs.

## 44 B RTS game as an Reinforcement Learning environment

45 Our platform can be also used as an RL environment. In our code base we implement a framework  
46 that allows a straightforward interaction with the game environment in a canonical RL training loop.  
47 Below we detail the environment properties.

### 48 B.1 Observation space

49 We leverage both spatial representation of the map, as well as internal state of the game engine  
50 (e.g. units health points and attacking cool downs, the amount of resources, etc.) to construct an  
51 observation. We carefully address the fog of war, by masking out the regions of the map that have not  
52 been visited. In addition, we remove any unseen enemy units attributes from the observation. The  
53 partial observability of the environment makes it especially challenging to apply RL due to highly  
54 non-stationary state distribution.

### 55 B.2 Action space

56 At each timestep of the environment we predict an action for each of our units, both buildings and  
57 army. The action space is consequently large—for example, any unit can go to any location at each  
58 timestep. Prediction of an unit action proceeds in steps, we first predict an action type (e.g. MOVE or  
59 ATTACK), then, based on the action type, we predict the action outputs. For example, for the BUILD  
60 BUILDING action type the outputs will be the type of the future building and its location on the game  
61 map. We summarize all available action types and their structure in Table 3.

### 62 B.3 Reward structure

63 We support a sparse reward structure, e.g. the reward of 1 is issued to an agent at the end if the game  
64 is won, all the other timesteps receive the reward of 0. Such reward structure makes exploration an  
65 especially challenging given the large dimensionality of the action space and the planning horizon.

## 66 C Data collection

67 We design a data collection task based on ParlAI, a transparent framework to interact with human  
68 workers. We develop separate game control interfaces for both the *instructor* and the *executor* players,

Building name	Description
TOWN HALL	The main building of the game, it allows a player to train PEASANTS and serves as a storage for mined resources.
BARRACK	Produces SPEARMEN.
BLACKSMITH	Produces SWORDMEN.
STABLE	Produces CAVALRY.
WORKSHOP	Produces CATAPULT, DRAGON and ARCHER. The only building that can produce multiple unit types.
GUARD TOWER	A building that can attack enemies, but cannot move.

Table 1: The list of the building units available in the game.

Unit name	Description
PEASANT	Gathers minerals and constructs buildings, not good at fighting.
SPEARMAN	Effective against cavalry.
SWORDMAN	Effective against spearmen.
CAVALRY	Effective against swordmen.
DRAGON	Can fly over obstacles, can only be attacked by archers and towers.
ARCHER	Great counter unit against dragons.
CATAPULT	Easily demolishes buildings.

Table 2: The list of the army units available in the game.

Action Type	Action Output	Input Features
IDLE	NULL	NULL
CONTINUE	NULL	NULL
GATHER	resource_id	resources_features
ATTACK	enemy_unit_id	enemy_units_features
TRAIN UNIT	unit_type	unit_type_features
BUILD BUILDING	unit_type, (x,y)	unit_type_features, map_cells_features
MOVE	(x,y)	map_cells_features

Table 3: We implement a separate action classifier per action type, because each action type needs to model a probability distribution over different objects (Action Output). For example, for the ATTACK action we need estimate a probability distribution over all visible enemy units and predict an enemy unit id, or BUILD BUILDING action needs to model two probability distributions, one over building type to be constructed, and another over all possible  $(x, y)$  discrete location on the map where the future building will be placed.

69 and ask two humans to play the game collaboratively against a rule-based AI opponent. Both player  
70 have the same access to the game observation, but different control abilities.

71 The *instructor* control interface allows the human player to perform the following actions:

- 72 • **Issue** a natural language instruction to the *executor* at any time of the game. We allow any  
73 free-form language instruction.
- 74 • **Pause** the game flow at any time. Pausing allows the *instructor* to analyze the game state  
75 more thoroughly and plan strategically.
- 76 • **Warn** the *executor* player in case they do not follow issued instructions precisely. This option  
77 allows us to improve data quality, by filtering *executors* who do not follow instructions.

78 On the other hand, the *executor* player gets to:

- 79 • **Control** the game units by direct manipulation using computer’s input devices (e.g. mouse).  
80 The *executor* is tasked to complete the current instruction, rather than to win the game.
- 81 • **Ask** the *instructor* for either a new instruction, or a clarification.

82 Each human workers is assigned with either the *instructor* or the *executor* role randomly, thus the  
83 same person can experience the game on both ends over multiple attempts.

## 84 C.1 Quality control

85 To make sure that we collect data of high quality we take the following steps:

86 **Game manual** We provide a detailed list of instructions to a human worker at the beginning of  
87 each game and during the game’s duration. This manual aims to narrate a comprehensive overview  
88 various game elements, such as player roles, army and building units, control mechanics, etc. We  
89 also record several game replays that serve as an introductory guideline to the players.

90 **Onboarding** We implement an onboarding process to make sure that novice players are comfortable  
91 with the game mechanics, so that they can play with other players effectively. For this, we ask a

Strategy Name	Description
SIMPLE	This strategy first sends all 3 initially available PEASANTS to mine to the closest resource, then it chooses one army unit type from SPEARMAN, SWORDSMAN, CAVALRY, ARCHER, or DRAGON, then it constructs a corresponding building, and finally trains 3 units of the selected type and sends them to attack. The strategy then continuously maintains the army size of 3, in case an army unit dies.
MEDIUM	Same as SIMPLE strategy, only the size of the army is randomly selected between 3 and 7.
STRONG	This strategy is adaptive, and it reacts to the opponent’s army. In particular, this strategy constantly scouts the map using one PEASANT and to learn the opponent’s behaviour. Once it sees the opponent’s army it immediately trains a counter army based on the attack graph (see Fig. 1). Then it clones the MEDIUM strategy.
SECOND BASE	This strategy aims to build a second TOWN HALL near the second closest resource and then it uses the double income to build a large army of a particular unit type. The other behaviours is the same as in the MEDIUM strategy.
TOWER RUSH	A non-standard strategy, that first scouts the map in order to find the opponent using a spare PEASANT. Once it finds it, it starts building GUARD TOWERS close to the opponent’s TOWN HALL so they can attack the opponent’s units.
PEASANT RUSH	This strategy sends first 3 PEASANTS to mine, then it keeps producing more PEASANTS and sending them to attack the opponent. The hope of this strategy is to beat the opponent by surprise.

Table 4: The rule-based strategies we use as an opponent to the human players during data collection.

novice player to perform the *executor*’s duties and pair them with a bot that issues a pre-defined set of natural language instructions that implements a simple walkthrough strategy. We allocate enough time for the human player to work on the current instruction, and to also get comfortable with the game flow. We let the novice player play several games until we verify that they pass the required quality bar. We assess the performance of the player by running a set of pattern-matching scripts that verify if the performed control actions correspond to the issued instructions (for example, if an instruction says "build a barrack", we make sure that the player executes the corresponding low-level action). If the human player doesn’t pass our qualification requirements within 5 games, we prevent them from participating in our data collection going forward and filter their games from the dataset.

**Player profile** We track performance of each player, breaking it down by a particular role (e.g. *instructor* or *executor*). We gather various statistics about each player and build a comprehensive player profile. For example, for the *instructor* role we gather data such as overall win rate, the number of instructions issued per game, diversity of issued instructions; for the *executor* role we monitor how well they perform on the issued instruction (using a pattern matching algorithm), the number of warnings they receive from the *instructor*, and many more. We then use this profile to decide whether to upgrade a particular player to playing against stronger opponents (see Appendix C.2) in case they are performing well, or prevent them from participating in our data collection at all otherwise.

**Feedback** We use several initial round of data collection as a source of feedback from the human players. The received feedback helps us to improve the game quality. Importantly, after we finalize the game configuration, we disregard all the previously collected data in our final dataset.

**Final filtering** Lastly, we take another filtering pass against all the collected game replays and eliminate those replays that don’t meet the following requirements:

- A game should have at least 3 natural language instructions issued by the *instructor*.
- A game should have at least 25 low-level control actions issued by the *executor*.

By implementing all the aforementioned safe guards we are able to gather a high quality dataset.

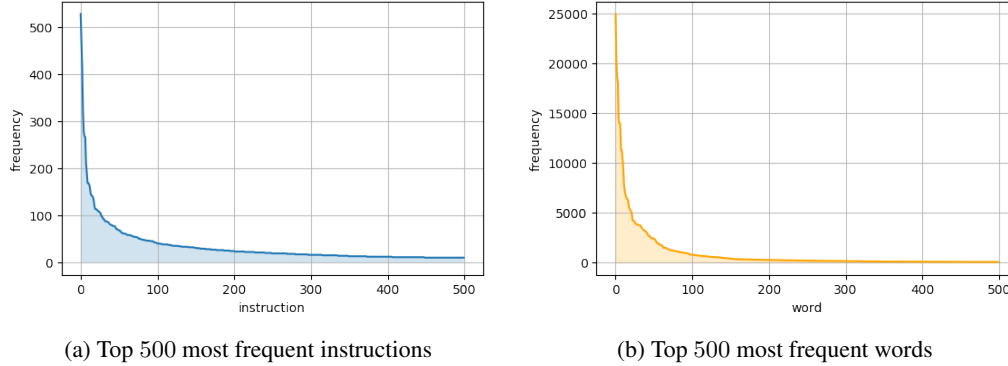


Figure 2: Frequency histograms for the dataset instructions and words.

## C.2 Rule-based bots

We design a set of diverse game strategies that are implemented by our rule-based bots ( Table 4). Our handcrafted strategies explore much of the possibilities that the game can offer, which in turn allows us to gather a multitude of emergent human behaviours in our dataset. Additionally, we employ a resource scaling hyperparameter, which controls the amount of resources a bot gets during mining. This hyperparameter offers a finer control over the bot’s strength, which we find beneficial for onboarding novice human players. We pair a team of two human players (the *instructor* and *executor*) with a randomly sampled instance of a rule-based strategy and the resource scaling hyperparameter during our data collection, so the human player doesn’t know in advance who is their opponent. This property rewards reactive players. We later observe that our models are able to learn the scouting mechanics from the data, which is a crucial skill to be successful in our game.

## D Model architecture

### D.1 Convolutional channels of Spatial Encoder

We use the following set of convolutional channels to extract different bits of information from spatial representation of the current observation.

1. **Visibility:** 3 binary channels for each state of visibility of a cell (VISIBLE, SEEN, and INVISIBLE).
2. **Terrain:** 2 binary channels for each terrain type of a cell (grass or water).
3. **Our Units:** 13 channels for each unit type of our units. Here, a cell contains the number of our units of the same type located in it.
4. **Enemy Units:** similarly 13 channels for visible enemy units.
5. **Resources:** 1 channel for resource units.

Linguistic Phenomena	Example
Counting	<i>Build 3 dragons.</i>
Spatial Reference	<i>Send him to the choke point behind the tower.</i>
Locations	<i>Build one to the left of that tower.</i>
Composed Actions	<i>Attack archers, then peasants.</i>
Cross-instruction anaphora	<i>Use it as a lure to kill them.</i>

Table 5: Complex linguistic phenomena emerge as humans instruct others how to play the game.

## 139 D.2 Action Classifiers

140 At each step of the game we predict actions for each of the player’s units, we do this by performing a  
 141 separate forward pass for ofv the following network for each unit. Firstly, we run an MLP (Fig. 3)  
 142 based action classifier to sample the unit’s ACTION TYPE. We feed the unit’s global summary features  
 143 (see Fig. 3 of the main paper) into the classifier and sample an action type (see Table 3 for the full list  
 144 of possible actions). Then, given the sampled action type we predict the ACTION OUTPUT based on  
 145 the unit’s features, unit dependent instructions features, and the action input features. We provide an  
 146 overview of ACTION OUTPUTS and INPUT FEATURES for each actions in Table 3. In addition, you  
 147 can refer to the diagram Fig. 4.

## 148 E Dataset details

149 Through our data collection we gather a dataset of  
 150 over 76 thousand of instructions and correspond-  
 151 ing executions. We observe a wide variety of dif-  
 152 ferent strategies and their realizations in natural  
 153 language. For example, we observe emergence of  
 154 complicated linguistic constructions (Table 5).

155 We also study the distribution of collected instruc-  
 156 tions. While we notice that some instructions are  
 157 more frequent than others, we still observe a good  
 158 coverage of strategies realizations, which serve as  
 159 a ground for generalization. In Table 7 we provide  
 160 a list of most frequently used instructions, and  
 161 in Fig. 2 shows the overall frequency distribution  
 162 for instructions and words in our dataset.

163 Finally, we provide a random sample of 50 instructions from our dataset in Table 6, where showing  
 164 the diversity and complexity of the collected instructions.

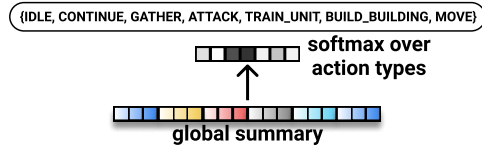


Figure 3: The ACTION TYPE classifier is parameterized as an MLP network to model a softmax distribution over action types based on the unit’s global summary features vector.

Instruction
<p> <i>Build 1 more cavalry.</i>  <i>Attack peasons.</i>  <i>Build barrack in between south pass at new town.</i>  <i>Have all peasants gather minerals next to town hall.</i>  <i>Have all peasants mine ore.</i>  <i>Fight u peaas.</i>  <i>Stop the peasants from mining.</i>  <i>Build a new town hall between the two west minerals patches.</i>  <i>Build 2 more swords.</i>  <i>Use cavalry to attack enemy.</i>  <i>Explore and find miners.</i>  <i>If you see any idle peasants please have them build.</i>  <i>Okay that doesn't work then build them on your side of the wall then.</i>  <i>Create 4 more archers.</i>  <i>Make a new town hall in the middle of all 3.</i>  <i>Attack tower with catas.</i>  <i>Kill cavalry and peasants then their townhall.</i>  <i>Attack enemy peasants with cavalry as well.</i>  <i>Send all peasants to collect minerals.</i>  <i>Attack enemy peasant.</i>  <i>Keep creating peasants and sending them to mine.</i>  <i>Send one catapult to attack the northern guard tower send a dragon for protection.</i>  <i>Send all but 1 peasant to mine.</i>  <i>Mine with the three peasants.</i>  <i>Use that one to scout and don't stop.</i>  <i>Bring scout back to base to mine.</i>  <i>You'll need to attack them with more peasants to kill them.</i>  <i>Build a barracks.</i>  <i>Send all peasants to find a mine and mine it.</i>  <i>Start mining there with your 3.</i>  <i>Make four peasants.</i>  <i>Move archers west then north.</i>  <i>Attack with cavalry.</i>  <i>Make two more workers.</i>  <i>Make 2 more calvary and send them over with the other ones.</i>  <i>Return to base with scout.</i>  <i>Build 2 peasants at the new mine.</i>  <i>If attacked retreat south.</i>  <i>Make the rest gather minerals too.</i>  <i>All peasants flee the enemy.</i>  <i>Attack the peasants in the area.</i>  <i>Attack the last archer with all peasants on the map.</i> </p>

Table 6: Examples of randomly sampled instructions.

Instruction	Frequency	Instruction	Frequency
<i>Attack.</i>	527	<i>Send idle peasants to mine.</i>	68
<i>Send all peasants to mine.</i>	471	<i>Attack that peasant.</i>	68
<i>Build a workshop.</i>	414	<i>Send all peasants to mine min-</i>	65
		<i>erals.</i>	
<i>Retreat.</i>	323	<i>Build a barracks.</i>	64
<i>Build a stable.</i>	278	<i>Build barrack.</i>	62
<i>Send peasants to mine.</i>	267	<i>Return to mine.</i>	62
<i>All peasants mine.</i>	266	<i>Build peasant.</i>	61
<i>Send idle peasant to mine.</i>	211	<i>Build catapult.</i>	61
<i>Build workshop.</i>	191	<i>Create a dragon.</i>	61
<i>Build a dragon.</i>	168	<i>Mine with peasants.</i>	60
<i>Kill peasants.</i>	168	<i>Build 3 peasants.</i>	59
<i>Attack enemy.</i>	166	<i>Defend.</i>	58
<i>Attack peasants.</i>	159	<i>Build cavalry.</i>	58
<i>Build a guard tower.</i>	146	<i>Make an archer.</i>	58
<i>Attack the enemy.</i>	142	<i>Attack dragon.</i>	58
<i>Stop.</i>	141	<i>Send all peasants to collect min-</i>	57
		<i>erals.</i>	
<i>Attack peasant.</i>	139	<i>Defend base.</i>	57
<i>Kill that peasant.</i>	132	<i>Build 2 more peasants.</i>	56
<i>Mine.</i>	119	<i>Build 2 peasants.</i>	55
<i>Build another dragon.</i>	113	<i>Make 2 archers.</i>	55
<i>Make another peasant.</i>	113	<i>Make dragon.</i>	54
<i>Build stable.</i>	112	<i>Build 2 dragons.</i>	54
<i>Make a dragon.</i>	110	<i>Attack dragons.</i>	54
<i>Build a blacksmith.</i>	108	<i>Make a stable.</i>	53
<i>Build a catapult.</i>	108	<i>Make a catapult.</i>	53
<i>Back to mining.</i>	106	<i>Build 6 peasants.</i>	52
<i>Build another peasant.</i>	104	<i>Attack archers.</i>	50
<i>Make a peasant.</i>	98	<i>Kill all peasants.</i>	50
<i>Build a barrack.</i>	97	<i>Build 2 catapults.</i>	50
<i>Build 4 peasants.</i>	93	<i>Idle peasant mine.</i>	49
<i>Have all peasants mine.</i>	92	<i>Make peasant.</i>	48
<i>Build 2 archers.</i>	90	<i>Attack enemy peasant.</i>	48
<i>Build dragon.</i>	87	<i>Attack archer.</i>	48
<i>Attack with peasants.</i>	87	<i>Build another archer.</i>	47
<i>Return to mining.</i>	87	<i>Make 4 peasants.</i>	47
<i>Build a peasant.</i>	86	<i>Make 3 peasants.</i>	47
<i>Idle peasant to mine.</i>	85	<i>Build 2 more archers.</i>	46
<i>Make a workshop.</i>	83	<i>Send idle peasant back to mine.</i>	46
<i>Create a workshop.</i>	81	<i>Make more peasants.</i>	46
<i>Mine with all peasants.</i>	80	<i>Make 2 more peasants.</i>	46
<i>Build 3 more peasants.</i>	79	<i>Build blacksmith.</i>	46
<i>Create another peasant.</i>	79	<i>Collect minerals.</i>	45
<i>Send all idle peasants to mine.</i>	77	<i>Kill.</i>	45
<i>Build 3 archers.</i>	77	<i>Build an archer.</i>	45
<i>Kill peasant.</i>	77	<i>Keep mining.</i>	45
<i>Make another dragon.</i>	76	<i>Keep attacking.</i>	43
<i>Kill him.</i>	72	<i>Attack dragons with archers.</i>	43
<i>Build guard tower.</i>	70	<i>Create a stable.</i>	42
<i>Attack town hall.</i>	70	<i>Make 3 more peasants.</i>	42
<i>Start mining.</i>	69	<i>Attack the peasant.</i>	41

Table 7: The top 100 instructions sorted by their usage frequency.



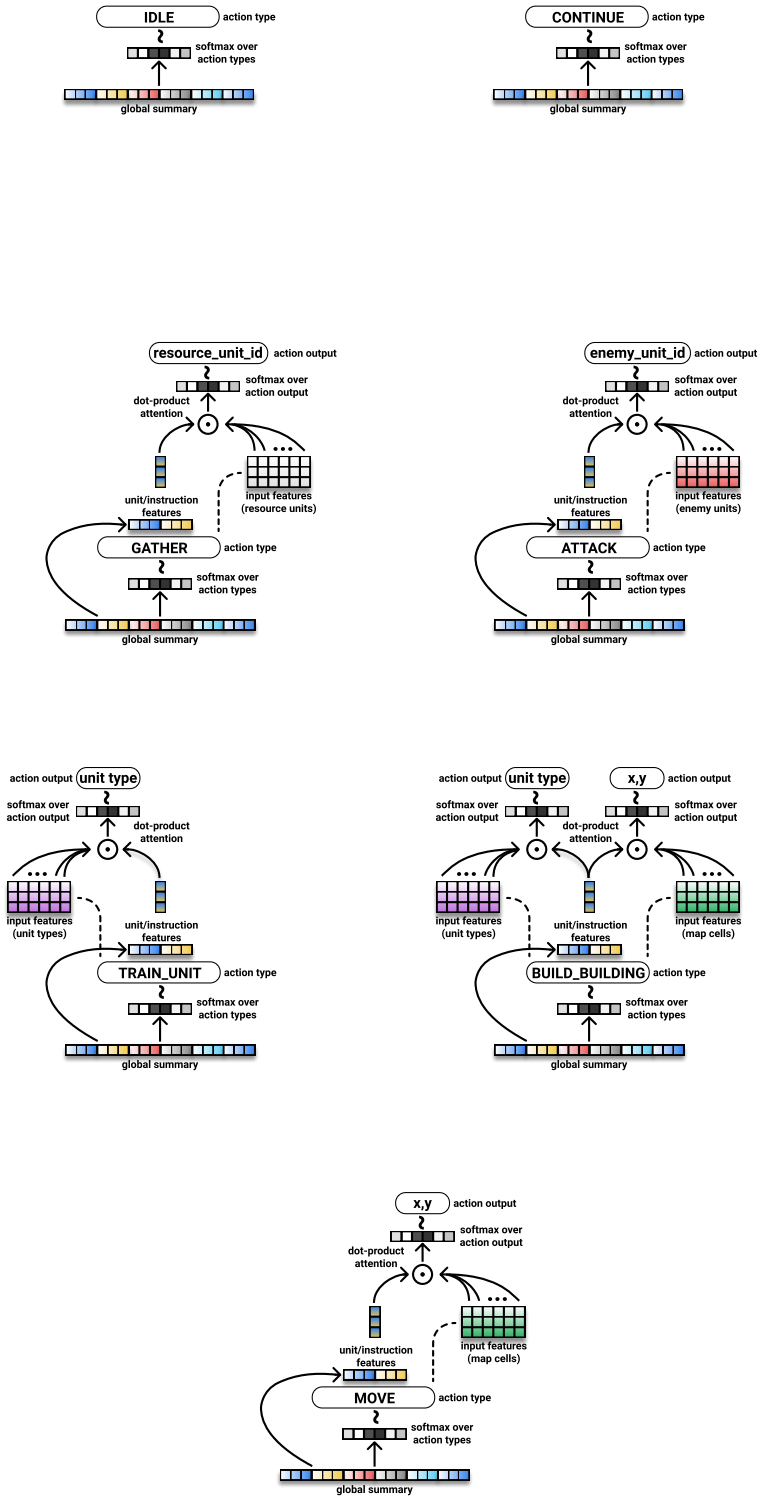


Figure 4: Separate classifiers for each of the available action types.