

Divide and Coordinate: A Multi-Policy Framework for Multi-Objective Reinforcement Learning

Anonymous authors
Paper under double-blind review

Keywords: Multi-Objective RL, Multi-Agent RL

Summary

The summary appears on the cover page. Although it can be identical to the abstract, it does not have to be. One might choose to omit the stated contributions in the Summary, given that they will be stated in the box below. The original abstract may also be extended to two paragraphs. The authors should ensure that the contents of the cover page fit entirely on a single page. The cover page does **not** count towards the 8–12 page limit.

 Lorem ipsum dolor sit amet, consectetuer adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetuer id, vulputate a, magna. Donec vehicula augue eu neque. Pellentesque habitant morbi tristique senectus et netus et malesuada fames ac turpis egestas. Mauris ut leo. Cras viverra metus rhoncus sem. Nulla et lectus vestibulum urna fringilla ultrices. Phasellus eu tellus sit amet tortor gravida placerat. Integer sapien est, iaculis in, pretium quis, viverra ac, nunc. Praesent eget sem vel leo ultrices bibendum. Aenean faucibus. Morbi dolor nulla, malesuada eu, pulvinar at, mollis ac, nulla. Curabitur auctor semper nulla. Donec varius orci eget risus. Duis nibh mi, congue eu, accumsan eleifend, sagittis quis, diam. Duis eget orci sit amet orci dignissim rutrum.

Contribution(s)

1. Provide a succinct but precise list of the contribution(s) of the paper. Use contextual notes to avoid implications of contributions more significant than intended and to clarify and situate the contribution relative to prior work (see the examples below). If there is no additional context, enter “None”. Try to keep each contribution to a single sentence, although multiple sentences are allowed when necessary. If using complete sentences, include punctuation. If using a single sentence fragment, you may omit the concluding period. A single contribution can be sufficient, and there is no limit on the number of contributions. Submissions will be judged mostly on the contributions claimed on their cover pages and the evidence provided to support them. Major contributions should not be claimed in the main text if they do not appear on the cover page. Overclaiming can lead to a submission being rejected, so it is important to have well-scoped contribution statements on the cover page.

Context: None

2. The submission template for submissions to RLJ/RLC 2025

Context: Built from previous RLC/RLJ, ICLR, and TMLR submission templates

3. [Example of one contribution and corresponding contextual note for the paper “Policy gradient methods for reinforcement learning with function approximation” (?).]

This paper presents an expression for the policy gradient when using function approximation to represent the action-value function.

Context: Prior work established expressions for the policy gradient without function approximation (?).

Divide and Coordinate: A Multi-Policy Framework for Multi-Objective Reinforcement Learning

Anonymous authors

Paper under double-blind review

Abstract

1

2 1 Introduction

3 **TODO:** I am putting this here, it will go at the end of the introduction.

4 Prior works proposed a composition technique based on Q-learning. Each local policy π^i for each
5 individual reward function R^i would be designed using a Q-learning agent that disregards all reward
6 functions other than its own. Along the Q-function, it also learns a W-function which maps every
7 state to a numeric importance score (Humphrys, 1995). Intuitively, if $W(s)$ is high, then it is highly
8 important for the local policy to be able to execute its action in the state s . The composition of
9 policies happens at runtime, when at each state s , if $W^i(s)$ is the W-value of the i -th local policy,
10 for $i \in [1; m]$, and if $i^* = \arg \max_i W^i(s)$, then we select the action proposed by the policy π^{i^*} at
11 the current state s . It has been demonstrated that, interestingly, W-learning generates selfish local
12 policies that end up cooperating in practice. Subsequently, this framework has been extended to
13 deep learning and applied to realistic applications (Rosero et al., 2024).

14 A limitation of W-learning is that it assumes that all local policies will be honest while broadcasting
15 their W-values: if any of the policies is dishonest, i.e., emits a higher W-value than the actual, then it
16 will get undue advantages in executing its actions, potentially compromising the global performance.
17 To put it in game theoretic terminologies, the local policies are not “strategyproof.” This could be a
18 serious issue if, e.g., the local policies are obtained through different third-party vendors.

19 1.1 Related work

- 20 1. Fairly comprehensive reference survey of multi-objective RL: ?
- 21 2. Reference for definitions of multi agent RL: ?.
- 22 3. W-learning and more recent Deep W learning: Humphrys (1995); Rosero et al. (2024)
- 23 4. General multi-objective deep RL works. ? introduce a multi-policy DQN algorithm to learn
24 multiple policies in parallel such that they have access to any possible linear weighting of the
25 objectives. Also has a good lit review from where I got a few of these other citations.
- 26 5. Multi objective Q learning (tabular): ?, ? for pareto Q learning (maintains a set of policies and
27 updates them), local search after normal learning: ?
- 28 6. Use gradient approximation to iteratively build parametrized policies representing pareto fron-
29 tier ?
- 30 7. Maximin (fairness) MORL: ??, fair deep RL: ?.
- 31 8. These two are only vaguely related and look at decomposing an objective into multiple smaller
32 objectives: Q-function decomposition for a single agent with multiple reward components: ? and
33 on a similar note hybrid reward architecture: ?

34 **2 Preliminaries: Multi-Objective MDPs**

35 Mainstream RL algorithms consider Markov decision processes (MDP) equipped with a *single* re-
 36 ward function, pertaining to a single task or *objective* for the system. In reality, a majority of real-
 37 world applications of RL requires satisfying multiple, partly contradictory objectives. We model
 38 such multi-objective decision-making problems using multi-objective MDPs (MO-MDP), as for-
 39 mally defined below. Intuitively, an MO-MDP has the exact same syntax as a regular MDP, except
 40 that it now has multiple reward functions pertaining to the different objectives. We formalize MO-
 41 MDP below. We will use the notation $\mathbb{D}(\Sigma)$ to represent the set of all probability distributions over
 42 a given alphabet Σ .

43 **Definition 1** (MO-MDP). A multi-objective Markov decision process (MO-MDP) with $m \in \mathbb{Z}_{>0}$
 44 objectives is specified by a tuple $\mathcal{M} = (S, A, T, \mathbf{R}, \mu_0)$, where

- 45 • S is the set of states,
- 46 • A is the set of actions,
- 47 • $T : S \times A \rightarrow \mathbb{D}(S)$ is the transition function mapping a state-action pair to a distribution over the
 48 successor states,
- 49 • $\mathbf{R} = \{R^i : S \times A \times S \rightarrow \mathbb{R}_{\geq 0}\}_{i \in [1; m]}$ is the set of reward functions, and
- 50 • $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

51 The notions of policies and paths induced by them are exactly the same as in classical MDPs, which
 52 we briefly recall below. First, we introduce some notation. Given an alphabet Σ , we will write
 53 Σ^* and Σ^ω to denote the set of every finite and infinite word over Σ , respectively, and will write
 54 $\Sigma^\infty = \Sigma^* \cup \Sigma^\omega$. Given a word $w = \sigma_0 \sigma_1 \dots \in \Sigma^\infty$, and given a $t \geq 0$ that is not larger than the
 55 length of w , we will write w_t and $w_{0:t}$ to denote respectively the t -th element of w , i.e., $w_t = \sigma_t$,
 56 and the prefix of w up to the t -th element, i.e., $w_{0:t} = \sigma_0 \dots \sigma_t$.

57 A *policy* in an MO-MDP \mathcal{M} is a function $\pi : (S \times A)^* \times S \rightarrow \mathbb{D}(A)$ that maps a history of
 58 state-action pairs and the current state to a distribution over actions. A *path* on \mathcal{M} induced by π is a
 59 sequence $\rho = (s_0, a_0)(s_1, a_1), \dots \in (S \times A)^\infty$ such that for every $t \geq 0$, (1) the probability that the
 60 action a_{t+1} is picked by π based on the history is positive, i.e., $\pi(\rho_{0:t}, s_{t+1})(a_{t+1}) > 0$, and (2) the
 61 probability of moving to the state s_{t+1} from s_t due to action a_t is positive, i.e., $T(s_t, a_t)(s_{t+1}) >$
 62 0. A path can be either finite or infinite, and we will write $Paths(\mathcal{M}, \pi)$ to denote the set of all
 63 infinite paths fo \mathcal{M} induced by π . Given a finite path $\rho = (s_0, a_0) \dots (s_t, a_t)$, the probability that ρ
 64 occurs is given by: $\mu_0(s_0) \cdot \prod_{k=0}^{t-1} T(s_k, a_k)(s_{k+1}) \cdot \pi(\rho_{0:k}, s_{k+1})(a_{k+1})$. This can be extended to
 65 a probability measure over the set of all infinite paths in \mathcal{M} using standard constructions, which can
 66 be found in the literature (Baier & Katoen, 2008). Given a measurable set of paths Ω and a function
 67 $f : Paths(\mathcal{M}, \pi) \rightarrow \mathbb{R}$, we will write $\mathbb{P}^{\mathcal{M}, \pi}[\Omega]$ and $\mathbb{E}^{\mathcal{M}, \pi}[f]$ to denote, respectively, the probability
 68 measure of Ω and the expected value of f evaluated over random infinite paths.

69 We will use the standard discounted reward objectives, where we fix $\gamma \in (0, 1)$ as a given dis-
 70 counting factor. Let $\rho = (s_0, a_0)(s_1, a_1), \dots \in Paths(\mathcal{M}, \pi)$ be an infinite path induced by
 71 π . Define the discounted sum function, mapping ρ to the discounted sum of the associated re-
 72 wards: $f_{ds}^i(\rho) := \sum_{t=0}^{\infty} \gamma^t \cdot R^i(s_t, a_t)$. The *i-value* of the policy ρ for \mathcal{M} is the expected
 73 value of the discounted sum of the i -th reward we can secure by executing ρ on \mathcal{M} , written as
 74 $val^{\mathcal{M}, i}(\pi) = \mathbb{E}^{\mathcal{M}, \pi}[f_{ds}^i]$. The *optimal* policy for R^i for a given $i \in [1; m]$ is the policy that maxi-
 75 mizes the i -value. When the reward index i is unimportant, we will refer to every element of the set
 76 $\{val^{\mathcal{M}, i}\}_{i \in [1; m]}$ as a *value component*.

77 When the MO-MDP \mathcal{M} is clear from the context, we will drop it from all notation and will simply
 78 write $Paths(\pi)$, \mathbb{P}^π , \mathbb{E}^π , and val^i .

79 It is known that *memoryless* (aka, stationary) policies suffice for maximizing single discounted re-
 80 ward objectives, where a policy π is called memoryless if the proposed action only depend on the

81 current state. In other words, given every pair of finite paths ρ, ρ' both ending at the same state, the
82 probability distributions $\pi(\rho)$ and $\pi(\rho')$ are identical.

83 Unlike classical single-objective MDPs, the optimal policy synthesis problem for MO-MDP requires
84 fixing one of many possible optimality criteria. Many possibilities exist, including pareto optimality,
85 requiring a solution where none of the value components could be unanimously improved without
86 hurting the others; weighted social welfare, requiring a weighted sum of the value components be
87 maximized; and fairness, requiring the minimum attained value by any value component is maxi-
88 mized. **TODO:** Give some citations for each category.

89 **3 Auction-Based Compositional RL on Multi-Objective MDPs**

90 We consider the compositional approach to policy synthesis for MO-MDPs, where we will design a
91 selfish, *local* policy maximizing each individual value component, towards the fulfillment of some
92 required global coordination requirements. The main crux is in the composition process, where
93 each local policy may propose a different action, but the composition must decide one of the actions
94 that will be actually executed. Importantly, the composition must be implementable in a distributed
95 manner, meaning we will *not* use any global policy that would pick an action by analyzing all local
96 policies and their reward functions. **TODO:** running example

97 **3.1 The Framework**

98 We present a novel *auction*-based RL framework for compositional policy synthesis for MO-MDPs.
99 In our framework, not only do the local policies emit actions, but also they *bid* for the privilege of
100 executing their actions for a given number of time steps $\tau \in \mathbb{N}_{>0}$ in future. The bids are all non-
101 negative real numbers, and the highest bidder's actions get executed for the following τ consecutive
102 steps, with bidding ties being resolved uniformly at random. The policy whose actions are executed
103 is referred to as the *winning* policy, and it must pay a bidding *penalty* that equals to its bid amount;
104 this is to discourage overbidding. The policies whose actions are not executed are called the *losing*
105 policies, and we consider three different settings for the “payment” they must make:

106 **Loser-Rewarded:** the winning policy pays the bidding penalty and the losing policies earn bid-
107 ding rewards equal to their respective bid values;

108 **Winner-Pays:** the winning policy pays the bidding penalty and the losing policies are unaffected
109 (i.e., neither earn bidding rewards nor pay bidding penalties);

110 **All-Pay:** all policies pay bidding penalties equal to their respective bid values.

111 While penalizing the winner discourages overbidding, the situation with the losers is more subtle.
112 In the **Loser-Rewarded** setting, by rewarding the losers, we encourage policies to bid positively if
113 the current state has some importance to them; this way, if they lose the bidding, they will get some
114 positive reward. In the **All-Pay** setting, by penalizing all policies, we discourage policies to bid at all
115 unless it is absolutely important. The **Winner-Pays** setting balances these two: by neither rewarding
116 nor penalizing the losers, we neither encourage nor discourage policies to bid. In Section 3.3, we
117 will see how these three settings induce different kinds of coordination through bidding.

118 For each policy, the bidding penalty or reward gets, respectively, subtracted or added to the *nominal*
119 reward obtained from the reward functions of the given MO-MDP, and the resulting reward is called
120 the *net* reward.

121 In summary, through this novel bidding mechanism, each policy can adjust its bid in proportion to the
122 importance for it to execute its action in the current state, and the associated bidding penalty/reward
123 aims to incentivize policies to be truthful. By making the highest bidder active, it is effectively
124 guaranteed that the most important policy is executed. This way, we obtain a purely decentralized
125 scheme to coordinate local policies in a given MO-MDP.

126 *Remark 1* (On the parameter τ). The parameter τ controls how frequently the agent changes its
 127 policies. In practice, if τ is too small, the switching could be too frequent for any of the objectives
 128 to be fulfilled. For example, **TODO: running example...**

129 **3.2 The Design Problem and Learning Algorithms**

130 We consider the following learning task for our auction-based compositional framework:

131 *Given an MO-MDP, a constant $\tau > 0$, and $\Delta \in \{\text{Loser-Rewarded, Winner-Pays, All-Pay}\}$,*
 132 *compute local policies that are optimal for the net rewards obtained in the mode Δ , given that all*
 133 *other local policies behave selfishly towards maximizing their own net rewards.*

134 We will show how the above learning problem boils down to solving a standard learning problem in
 135 the multi-agent setting, formalized using a decentralized MDP (DEC-MDP) as defined below. The
 136 only difference between a DEC-MDP and an MO-MDP (see Definition 1) is that now each reward
 137 function R^i is owned by the Agent i , who now controls a separate set of actions A^i .

138 **Definition 2** (DEC-MDP). A decentralized Markov decision process (DEC-MDP) with $m \in \mathbb{Z}_{>0}$
 139 agents is specified by a tuple $\mathcal{M} = (S, \mathbf{A}, T, \mathbf{R}, \mu_0)$, where

- 140 • S is the set of states,
- 141 • $\mathbf{A} = \{A^1, \dots, A^m\}$ is a set with A^i being the set of Agent i 's actions,
- 142 • $T : S \times A^1 \times \dots \times A^m \rightarrow \mathbb{D}(S)$ is the transition function mapping a state-action pair to a
 143 distribution over the successor states,
- 144 • $\mathbf{R} = \{R^i : S \times A^1 \times \dots \times A^m \times S \rightarrow \mathbb{R}_{\geq 0}\}_{i \in [1; m]}$ is the set of reward functions, and
- 145 • $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

146 The definitions of policies and paths readily extend from MO-MDP to DEC-MDP.

147 Given a DEC-MDP, the goal is to compute an ensemble of local (memoryless) policies for all individual agents, such that for every $i \in [1; m]$, the i -value cannot be increased by a unanimous change
 148 of the local policy π^i . In other words, the goal is to find a set of selfish local policies that are in a
 149 Nash equilibrium. This is an extensively studied problem in the literature. **TODO: Do a little bit of**
 150 **literature survey...**

152 Our focus is not in improved algorithms for DEC-MDP, but rather to show how the local policy synthesis
 153 problem for the MO-MDP \mathcal{M} in our auction-based framework reduces to the multi-agent policy
 154 synthesis problem in a DEC-MDP $\widetilde{\mathcal{M}}$. Intuitively, for every state s of \mathcal{M} , $\widetilde{\mathcal{M}}$ creates two kinds
 155 of copies, ones where bidding happens and are represented simply as s , and ones of the form (s, t, i^*)
 156 that keeps track of the time t elapsed since the last bidding, and the winner i^* of the last bidding. Fur-
 157 thermore, bidding is facilitated by extending the action space of \mathcal{M} to include all real-valued bids,
 158 and each agent in $\widetilde{\mathcal{M}}$ has an identical copy of this extended action space. After bidding in a state s ,
 159 the winner i^* is selected, and the state moves to $(s, 0, i^*)$. From this point onward, only Agent i^* se-
 160 lects actions a^0, a^1, \dots, a^τ to produce the sequence $(s^1, 1, i^*), (s^2, 2, i^*), \dots, (s^{\tau-1}, \tau-1, i^*), s^\tau$,
 161 after which the next bidding happens, and the process repeats. Finally, the bidding penalties or bid-
 162 ding rewards are only paid during the transition $s \rightarrow (s, 0, i^*)$, otherwise, the rewards are inherited
 163 from the original MO-MDP.

164 We formalize this below. Given an MO-MDP $\mathcal{M} = (S, A, T, \mathbf{R}, \mu_0)$, a constant $\tau > 0$, and
 165 the mode $\Delta \in \{\text{Loser-Rewarded, Winner-Pays, All-Pay}\}$, we define the DEC-MDP $\widetilde{\mathcal{M}} =$
 166 $(\widetilde{S}, \widetilde{\mathbf{A}}, \widetilde{T}, \widetilde{\mathbf{R}}, \widetilde{\mu}_0)$ where

- 167 • $\widetilde{S} := S \cup S \times [0; \tau - 1] \times [1; m]$,
- 168 • $\widetilde{\mathbf{A}} := \{\widetilde{A}^i\}_{i \in [1; m]}$ where $\widetilde{A}^i := A \cup \mathbb{R}$,
- 169 • $\widetilde{\mu}_0 := \mu_0 \times \{0\}$,

170 and for every current state $s \in \tilde{S}$ and every current action $(b^1, \dots, b^m) \in \mathbb{R}^m$, writing the highest
171 bidders as $I = \{i \in [1; m] \mid \forall j \in [1; m]. b^i \geq b^j\}$,

172 • $\tilde{T}(s, b^1, \dots, b^m) := \text{Uniform}(\{(s, 0, i)\}_{i \in I})$,

$$173 \quad \bullet \quad \tilde{R}^i(s, b^1, \dots, b^m, (s, 0, i^*)) := \begin{cases} -b^i & i = i^* \vee \Delta = \text{All-Pay}, \\ +b^i & i \neq i^* \wedge \Delta = \text{Loser-Rewarded}, \\ 0 & i \neq i^* \wedge \Delta = \text{Winner-Pays}, \end{cases}$$

174 whereas if the current state is of the form $(s, t, i^*) \in \tilde{S}$, for every action $(a^1, \dots, a^m) \in A^m$,

$$175 \quad \bullet \quad \tilde{T}((s, t, i^*), a^1, \dots, a^m) := \begin{cases} T(s, a^{i^*}) \times ((t+1) \bmod \tau) \times \{i^*\} & t < \tau-1, \\ T(s, a^{i^*}) & t = \tau-1, \end{cases}$$

176 • $\tilde{R}^i((s, t, i^*), a^1, \dots, a^m, (s', t+1, i^*)) := R^i(s, a^i, s')$.

177 **KM: A soundness theorem would be good, but what can we say concretely?**

178 3.3 Flavors of Cooperation through Bidding

179 We provide theoretical insights into the global behavior that emerges out of the auction-based inter-
180 actions between the local policies. For the sake of theoretical guarantees, and to be able to convey
181 the main essence of our results, we choose the simplest bare bone setting:

182 **Assumption 1.** The given MO-MDP has finite state and action spaces, and for every (memoryless)
183 policy, the bottom strongly connected component (BSCC) of the resulting Markov chain (MC) is a
184 sink state where no reward is earned. Furthermore, the time parameter $\tau = 1$, meaning the bidding
185 takes place at each time step before selecting the action.

186 Firstly, since the MO-MDP is finite, for each individual reward function, *deterministic* memoryless
187 policy suffices. **TODO: give some citation**

188 The following two types of global behaviors are of particular interest:

189 **Social welfare** is the sum (equivalently, the average) of the i -values for all i . We may ask: is the
190 emergent global behavior guaranteed to achieve the maximal social welfare?

191 **Fairness** is measured by the disparity between different i -values, i.e., $\max_{i,j \in [1;m]} |val^i - val^j|$.
192 Fairness is maximized when the disparity is minimized. We may ask: is the emergent global behav-
193 ior guaranteed to achieve the maximal fairness?

194 **Theorem 1.** Suppose the MO-MDP is such that at each state s and for every action a , there exists
195 at most a single $i \in [1; m]$ such that the optimal policy for R^i selects a at s . Then, the **Loser-**
196 **Rewarded** setting maximizes the social welfare.

197 *Proof sketch.* First, consider the simple one-shot game, where the agents bid just one time to select
198 an action, and the reward is based on the resulting single probabilistic transition. Suppose for the
199 index $i \in [1; m]$, the expected reward from using the action $a \in A$ is E_a^i , and define $E_+^i :=$
200 $\max_{a \in A} E_a^i$ and $E_-^i := \min_{a \in A} E_a^i$.

201 We claim that the optimal bid b_*^i for policy i equals $(E_+^i - E_-^i)/2$, and upon winning the bidding
202 the optimal action is $a_+ = \arg \max_{a \in A} E_a^i$. Notice that no matter whether policy i becomes the
203 winner or the loser, its net reward is at least $(E_+^i + E_-^i)/2$: if it wins and chooses a_+ , after paying
204 the bidding penalty, the net reward is $E_+^i - (E_+^i - E_-^i)/2 = (E_+^i + E_-^i)/2$; if it loses, no matter
205 what action the opponent chooses, its nominal reward is at least E_-^i , and after the bidding reward,
206 the net reward is $E_-^i + (E_+^i - E_-^i)/2 = (E_+^i + E_-^i)/2$. If policy i bids $b^i < b_*^i$, then upon
207 losing, its net reward will be $E_-^i + b^i < E_-^i + b_*^i = (E_+^i + E_-^i)/2$. If it bids $b^i > b_*^i$, then upon
208 winning, its net reward will be $E_+^i - b^i < E_+^i - b_*^i = (E_+^i + E_-^i)/2$. Therefore, the optimal bid is
209 $b_*^i = (E_+^i - E_-^i)/2$, which is what each selfish policy is expected to select.

210 Suppose, policy i is the winner. Then, for every $j \neq i$, $b_*^i \geq b_*^j$, i.e., $(E_+^i - E_-^i)/2 \geq (E_+^j - E_-^j)/2$.
 211 Simplifying, we get $E_+^i + E_-^j \geq E_-^i + E_+^j$. It follows that $E_+^i + \sum_{j \neq i} E_-^j \geq E_-^i + \sum_{j \neq i} E_+^j \geq$
 212 $E_-^i + E_+^k + \sum_{j \neq i, k} E_-^j$ for every for every $k \neq i$. Since the MO-MDP is purely competitive, there
 213 will be at least a single k such that a given action is optimal for k , and therefore the claim follows
 214 for the single-shot case.

215 Now, for the general multi-shot case, we inductively apply the above principle in the Bellman equa-
 216 tion, which extends the claim to paths of arbitrary length. The convergence of the Bellman iteration
 217 is guaranteed because it is a contraction mapping (since $\gamma < 1$). \square

218 4 A Multi-Agent Bidding Approach for Multi-Objective RL

219 **Definition 3** (MO-MDP). A multi-objective Markov decision process (MO-MDP) with $m \in \mathbb{Z}_{>0}$
 220 objectives is specified by a tuple $\mathcal{M} = (S, A, T, R, \mu_0)$, where

- 221 • S is the set of states,
- 222 • A is the set of actions,
- 223 • $T : S \times A \rightarrow \mathbb{D}(S)$ is the transition function mapping a state-action pair to a distribution over
 224 states,
- 225 • $R : S \times A \times S \rightarrow \mathbb{R}^m$ is the reward function with each output component corresponding to the
 226 different objectives, and
- 227 • $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

228 A *policy* in an MO-MDP \mathcal{M} is a function $\pi : (S \times A)^* \times S \rightarrow \mathbb{D}(A)$ that maps a history of
 229 state-action pairs and the current state to a distribution over actions.

230 **Definition 4** (MAB-MDP). Let $\mathcal{M} = (S, A, T, R, \mu_0)$ be an MO-MDP with m objectives and
 231 let $b \in \mathbb{Z}_{>0}$ be the bid upper bound. Also, define $M = \{1, \dots, m\}$ be indices of the m agents
 232 corresponding to the m objectives along with \perp representing a null agent. Lastly, let $B = \{0, \dots, b\}$
 233 be the range of bids and $\rho > 0$ be the bid penalty factor. We define the multi-agent bidding Markov
 234 decision process (MAB-MDP) as a tuple $\mathcal{B}_M = (\hat{S}, \hat{A}, \hat{T}, P, \hat{R}, \hat{\mu}_0)$ where

- 235 • $\hat{S} = M \times S$ is the new state space augmented with the index of the agent that won the previous
 236 round of bidding,
- 237 • $\hat{A} = A^m \times B^m$ represents the action space of the m agents in which each agent selects an action
 238 from A and a bid from B ,
- 239 • $\hat{T} : \hat{S} \times \hat{A} \rightarrow \mathbb{D}(\hat{S})$ is the new transition function defined as,

$$\hat{T}((_, s), (\mathbf{a}, \mathbf{b})) := \frac{1}{|B_{\max}|} \sum_{i \in B_{\max}} (T(s, a_i), i)$$

240 where $B_{\max} := \{i \mid b_i = \max\{b_1, \dots, b_m\}\}$ is the set of agents with maximal bids. The tuple
 241 $(T(s, a_i), i)$ represents the distribution over \hat{S} induced by the original transition function T such
 242 that the second component is fixed, and the weighted sum represents taking the weighted sums of
 243 the distributions over \hat{S} .

- 244 • $P : \hat{A} \times M \rightarrow \mathbb{R}^m$ is the bidding penalty for the m agents and the second component is the index
 245 of the agent that won the bidding.
- 246 • $\hat{R} : \hat{S} \times \hat{A} \times \hat{S} \rightarrow \mathbb{R}^m$ is the reward function for the m agents with

$$\hat{R}_k((_, s_0), (\mathbf{a}, \mathbf{b}), (i, s)) := R_k(s_0, a_i, s) - P_k((\mathbf{a}, \mathbf{b}), i)$$

247 where $i \in M$ is the index of the agent that won the bid and chose the action.

- 248 • $\hat{\mu}_0 := (\mu_0, 1)$ is the initial state distribution over \hat{S} induced by μ_0 and the second component is
 249 fixed to be 1 without loss of generality.

250 Given an MAB-MDP $\mathcal{B}_{\mathcal{M}}$, a *policy* for each agent indexed by $i \in \{1, \dots, m\}$ takes a similar form:
251 $\pi_i : (\hat{\mathcal{S}} \times \hat{\mathcal{A}})^* \times \hat{\mathcal{S}} \rightarrow \hat{\mathcal{A}}$. Intuitively, a state $(i, s) \in \hat{\mathcal{S}}$ encodes the agent that won the bidding and
252 chose the action to reach s in the previous step. At each step, each of the agents choose an action and
253 a bid, and an action amongst the set of highest bidders is chosen uniformly at random. The reward
254 function includes a penalty term that captures the desired bidding mechanism.

255 **5 Implementation and evaluation**

256 **5.1 Implementation**

257 Talk about:

- 258 1. different bidding mechanisms
- 259 2. choice of penalty factor
- 260 3. action window (remarking that we could additionally allow agents to choose length of action
261 window)
- 262 4. use with off-the-shelf RL algorithms

263 **5.2 Environments**

264 **5.2.1 MovingTargetsGridworld**

265 Important to mention that we want to maximize $\min(\text{targets reached})$.

266 **5.2.2 Atari Assault**

267 **5.3 Baselines**

- 268 1. Weighted sum of rewards with standard RL algorithms
- 269 2. Deep W learning implemented on top of DQN

270 **5.4 Performance comparison with baselines**

271 Include plots of training steps vs performance of our algorithms vs baselines on both environments

272 **5.5 Interpretability**

273 Include plots of distribution of control steps amongst agents, table of average, median, max, min of
274 bids of agents

275 **5.6 Modularity**

276 Plots of performance in gridworld with increasing number of objectives

277 **5.7 Ablations**

278 Impact of max bid, penalty factor

279 **References**

280 Christel Baier and Joost-Pieter Katoen. *Principles of model checking*. MIT press, 2008.

281 Mark Humphrys. W-learning: Competition among selfish q-learners. 1995.

- 282 Juan C Rosero, Nicolás Cardozo, and Ivana Dusparic. Multi-objective deep reinforcement learn-
283 ing optimisation in autonomous systems. In *2024 IEEE International Conference on Autonomic
284 Computing and Self-Organizing Systems Companion (ACSOS-C)*, pp. 97–102. IEEE, 2024.