# 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

## 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 $\pi^i$ 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 $\pi^{i^*}$ 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.

### 1.1  Related work

20  A large body of work in multi-objective reinforcement learning (MORL) relies on scalarization, ag-
21  gregating multiple reward functions into a single scalar objective so that standard single-objective
22  RL algorithms can be applied. The simplest scalarization method is a weighted sum of individual
23  rewards (Gass & Saaty, 1955), though richer nonlinear scalarization functions have also been pro-
24  posed (Van Moffaert et al., 2013). A key limitation of scalarization is that the relative importance
25  induced by the aggregation function may not align with the designer's true intent. This mismatch
26  can initiate a tedious debugging cycle, particularly in large-scale systems (Hayes et al., 2022). In
27  contrast, our approach achieves a trade-off between reward components without collapsing them
28  into a fixed scalar objective.

29  Other works pursue trade-offs by fixing a specific optimality criterion. Common choices include
30  Pareto optimality (Van Moffaert & Nowé, 2014) and its approximations (Pirotta et al., 2015), as
31  well as fairness-based criteria across reward functions (Park et al., 2024; Byeon et al., 2025; Siddique
32  et al., 2020). These approaches typically learn a single monolithic policy that satisfies the chosen
33  criterion. By contrast, our objective is to learn independent, selfish local policies for each reward
34  component and compose them at runtime in a principled manner, thereby preserving modularity
35  while still achieving a coherent global trade-off.

36  Relatively few works study distributed local policies for multiple rewards. A notable example is W-
37  learning (Humphrys, 1995) and its deep RL extension (Rosero et al., 2024), where separate selfish
38  policies are trained alongside meta-policies (W-functions) that assign each state a score reflecting
39  its urgency. At runtime, the policy with the highest score is selected. Other approaches employ
40  alternative aggregation mechanisms, such as ranked voting over actions (Méndez-Hernández et al.,
41  2019), or fixed aggregation rules like summing action values across agents (Russell & Zimdars,
42  2003). While conceptually related, our approach is technically simpler: it relies on an engineered
43  reward structure that enables the use of standard learning algorithms (e.g., PPO) without additional
44  meta-policies or complex aggregation schemes. Furthermore, to the best of our knowledge, we are
45  the first to introduce the incremental MORL setting, in which reward components can be added or
46  removed at runtime.

47  The idea of bidding-based selfish policies originates from analogous techniques for multi-objective
48  path planning problems on finite graphs (Avni et al., 2024), as well as from the broader literature
49  on bidding games (Lazarus et al., 1999; Avni et al., 2019; 2025). These works study strategic
50  interaction in finite arenas, where adversarial players bid for the right to determine the next move
51  from a shared action space in pursuit of their objectives. Although these works provide strong
52  theoretical guarantees, they do not naturally extend to infinite arenas. Moreover, players in such
53  games are typically budget-constrained, and the central question concerns the minimum budget
54  required to win. In contrast, we consider infinite arenas and eliminate explicit budget constraints
55  by incorporating bidding rewards and penalties directly into the learning framework.

## 56  2  Preliminaries: Multi-Objective MDPs

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

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

67  • $S$ is the set of states,

68  • $A$ is the set of actions,

69  • $T : S \times A \to \mathbb{D}(S)$ is the transition function mapping a state-action pair to a distribution over the
70    successor states,

71  • $\mathbf{R} = \{R^i : S \times A \times S \to \mathbb{R}_{\geq 0}\}_{i \in [1;m]}$ is the set of reward functions, and

72  • $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

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

79  A *policy* in an MO-MDP $\mathcal{M}$ is a function $\pi : (S \times A)^* \times S \to \mathbb{D}(A)$ that maps a history of
80  state-action pairs and the current state to a distribution over actions. A *path* on $\mathcal{M}$ induced by $\pi$ is a
81  sequence $\rho = (s_0, a_0)(s_1, a_1), \ldots \in (S \times A)^\infty$ such that for every $t \geq 0$, (1) the probability that the
82  action $a_{t+1}$ is picked by $\pi$ based on the history is positive, i.e., $\pi(\rho_{0:t}, s_{t+1})(a_{t+1}) > 0$, and (2) the
83  probability of moving to the state $s_{t+1}$ from $s_t$ due to action $a_t$ is positive, i.e., $T(s_t, a_t)(s_{t+1}) >$

84　0. A path can be either finite or infinite, and we will write $Paths(\mathcal{M}, \pi)$ to denote the set of all
85　infinite paths fo $\mathcal{M}$ induced by $\pi$. Given a finite path $\rho = (s_0, a_0) \ldots (s_t, a_t)$, the probability that $\rho$
86　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
87　a probability measure over the set of all infinite paths in $\mathcal{M}$ using standard constructions, which can
88　be found in the literature (Baier & Katoen, 2008). Given a measurable set of paths $\Omega$ and a function
89　$f : Paths(\mathcal{M}, \pi) \to \mathbb{R}$, we will write $\mathbb{P}^{\mathcal{M}, \pi}[\Omega]$ and $\mathbb{E}^{\mathcal{M}, \pi}[f]$ to denote, respectively, the probability
90　measure of $\Omega$ and the expected value of $f$ evaluated over random infinite paths.

91　We will use the standard discounted reward objectives, where we fix $\gamma \in (0, 1)$ as a given dis-
92　counting factor. Let $\rho = (s_0, a_0)(s_1, a_1), \ldots \in Paths(\mathcal{M}, \pi)$ be an infinite path induced by
93　$\pi$. Define the discounted sum function, mapping $\rho$ to the discounted sum of the associated re-
94　wards: $f_{\mathrm{ds}}^i(\rho) \coloneqq \sum_{t=0}^{\infty} \gamma^t \cdot R^i(s_t, a_t)$. The *i-value* of the policy $\rho$ for $\mathcal{M}$ is the expected
95　value of the discounted sum of the $i$-th reward we can secure by executing $\rho$ on $\mathcal{M}$, written as
96　$val^{\mathcal{M}, i}(\pi) = \mathbb{E}^{\mathcal{M}, \pi}[f_{\mathrm{ds}}^i]$. The *optimal* policy for $R^i$ for a given $i \in [1; m]$ is the policy that maxi-
97　mizes the $i$-value. When the reward index $i$ is unimportant, we will refer to every element of the set
98　$\{val^{\mathcal{M}, i}\}_{i \in [1; m]}$ as a *value component*.

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

101　It is known that *memoryless* (aka, stationary) policies suffice for maximizing single discounted re-
102　ward objectives, where a policy $\pi$ is called memoryless if the proposed action only depend on the
103　current state. In other words, given every pair of finite paths $\rho, \rho'$ both ending at the same state, the
104　probability distributions $\pi(\rho)$ and $\pi(\rho')$ are identical.

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

## 3　Auction-Based Compositional RL on Multi-Objective MDPs

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

### 3.1　The Framework

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

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

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

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

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

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

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

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

### 3.2 The Design Problem and Learning Algorithms

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

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

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

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

- $S$ is the set of states,

- $\mathbf{A} = \{A^1, \ldots, A^m\}$ is a set with $A^i$ being the set of Agent $i$'s actions,

- $T : S \times A^1 \times \ldots \times A^m \to \mathbb{D}(S)$ is the transition function mapping a state-action pair to a distribution over the successor states,

- $\mathbf{R} = \{R^i : S \times A^1 \times \ldots \times A^m \times S \to \mathbb{R}_{\geq 0}\}_{i \in [1;m]}$ is the set of reward functions, and

- $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

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

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 of the local policy $\pi^i$. In other words, the goal is to find a set of selfish local policies that are in a Nash equilibrium. This is an extensively studied problem in the literature. **TODO:** Do a little bit of literature survey...

174 Our focus is not in improved algorithms for DEC-MDP, but rather to show how the local policy syn-
175 thesis problem for the MO-MDP $\mathcal{M}$ in our auction-based framework reduces to the multi-agent pol-
176 icy synthesis problem in a DEC-MDP $\widetilde{\mathcal{M}}$. Intuitively, for every state $s$ of $\mathcal{M}$, $\widetilde{\mathcal{M}}$ creates two kinds
177 of copies, ones where bidding happens and are represented simply as $s$, and ones of the form $(s, t, i^*)$
178 that keeps track of the time $t$ elapsed since the last bidding, and the winner $i^*$ of the last bidding. Fur-
179 thermore, bidding is facilitated by extending the action space of $\mathcal{M}$ to include all real-valued bids,
180 and each agent in $\widetilde{\mathcal{M}}$ has an identical copy of this extended action space. After bidding in a state $s$,
181 the winner $i^*$ is selected, and the state moves to $(s, 0, i^*)$. From this point onward, only Agent $i^*$ se-
182 lects actions $a^0, a^1, \ldots, a^\tau$ to produce the sequence $(s^1, 1, i^*), (s^2, 2, i^*), \ldots, (s^{\tau-1}, \tau - 1, i^*), s^\tau$,
183 after which the next bidding happens, and the process repeats. Finally, the bidding penalties or bid-
184 ding rewards are only paid during the transition $s \to (s, 0, i^*)$, otherwise, the rewards are inherited
185 from the original MO-MDP.

186 We formalize this below. Given an MO-MDP $\mathcal{M} = (S, A, T, \mathbf{R}, \mu_0)$, a constant $\tau > 0$, and
187 the mode $\Delta \in \{\mathsf{Loser\text{-}Rewarded}, \mathsf{Winner\text{-}Pays}, \mathsf{All\text{-}Pay}\}$, we define the DEC-MDP $\widetilde{\mathcal{M}} =$
188 $(\widetilde{S}, \widetilde{\mathbf{A}}, \widetilde{T}, \widetilde{\mathbf{R}}, \widetilde{\mu}_0)$ where

189 • $\widetilde{S} := S \cup S \times [0; \tau - 1] \times [1; m]$,

190 • $\widetilde{\mathbf{A}} := \{\widetilde{A}^i\}_{i \in [1;m]}$ where $\widetilde{A}^i := A \cup \mathbb{R}$,

191 • $\widetilde{\mu}_0 := \mu_0 \times \{0\}$,

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

194 • $\widetilde{T}(s, b^1, \ldots, b^m) := \mathrm{Uniform}(\{(s, 0, i)\}_{i \in I})$,

195 • $\widetilde{R}^i(s, b^1, \ldots, b^m, (s, 0, i^*)) := \begin{cases} -b^i & i = i^* \vee \Delta = \mathsf{All\text{-}Pay}, \\ +b^i & i \neq i^* \wedge \Delta = \mathsf{Loser\text{-}Rewarded}, \\ 0 & i \neq i^* \wedge \Delta = \mathsf{Winner\text{-}Pays}, \end{cases}$

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

197 • $\widetilde{T}((s, t, i^*), a^1, \ldots, a^m) := \begin{cases} T(s, a^{i^*}) \times ((t+1) \mod \tau) \times \{i^*\} & t < \tau - 1, \\ T(s, a^{i^*}) & t = \tau - 1, \end{cases}$

198 • $\widetilde{R}^i((s, t, i^*), a^1, \ldots, a^m, (s', t+1, i^*)) := R^i(s, a^i, s')$.

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


## 3.3 Flavors of Cooperation through Bidding

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

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

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

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

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

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

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

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

223  We claim that the optimal bid $b_*^i$ for policy $i$ equals $(E_+^i - E_-^i)/2$, and upon winning the bidding
224  the optimal action is $a_+ = \arg\max_{a\in A} E_a^i$. Notice that no matter whether policy $i$ becomes the
225  winner or the loser, its net reward is at least $(E_+^i + E_-^i)/2$: if it wins and chooses $a_+$, after paying
226  the bidding penalty, the net reward is $E_+^i - (E_+^i - E_-^i)/2 = (E_+^i + E_-^i)/2$; if it loses, no matter
227  what action the opponent chooses, its nominal reward is at least $E_-^i$, and after the bidding reward,
228  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
229  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
230  winning, its net reward will be $E_+^i - b^i < E_+^i - b_*^i = (E_+^i + E_-^i)/2$. Therefore, the optimal bid is
231  $b_*^i = (E_+^i - E_-^i)/2$, which is what each selfish policy is expected to select.

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

237  Now, for the general multi-shot case, we inductively apply the above principle in the Bellman equa-
238  tion, which extends the claim to paths of arbitrary length. The convergence of the Bellman iteration
239  is guaranteed because it is a contraction mapping (since $\gamma < 1$). **KM:** I am not sure about this
240  extension. $\square$

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

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

244  • $S$ is the set of states,

245  • $A$ is the set of actions,

246  • $T : S \times A \to \mathbb{D}(S)$ is the transition function mapping a state-action pair to a distribution over
247    states,

248  • $R : S \times A \times S \to \mathbb{R}^m$ is the reward function with each output component corresponding to the
249    different objectives, and

250  • $\mu_0 \in \mathbb{D}(S)$ is the initial state distribution.

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

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

258 • $\hat{S} = M \times S$ is the new state space augmented with the index of the agent that won the previous
259     round of bidding,

260 • $\hat{A} = A^m \times B^m$ represents the action space of the $m$ agents in which each agent selects an action
261     from $A$ and a bid from $B$,

262 • $\hat{T} : \hat{S} \times \hat{A} \to \mathbb{D}(\hat{S})$ is the new transition function defined as,

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

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

267 • $P : \hat{A} \times M \to \mathbb{R}^m$ is the bidding penalty for the $m$ agents and the second component is the index
268     of the agent that won the bidding.

269 • $\hat{R} : \hat{S} \times \hat{A} \times \hat{S} \to \mathbb{R}^m$ is the reward function for the $m$ agents with

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

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

271 • $\hat{\mu}_0 \coloneqq (\mu_0, 1)$ is the initial state distribution over $\hat{S}$ induced by $\mu_0$ and the second component is
272     fixed to be 1 without loss of generality.

273 Given an MAB-MDP $\mathcal{B}_\mathcal{M}$, a *policy* for each agent indexed by $i \in \{1, \ldots, m\}$ takes a similar form:
274 $\pi_i : (\hat{S} \times \hat{A})^* \times \hat{S} \to \hat{A}$. Intuitively, a state $(i, s) \in \hat{S}$ encodes the agent that won the bidding and
275 chose the action to reach $s$ in the previous step. At each step, each of the agents choose an action and
276 a bid, and an action amongst the set of highest bidders is chosen uniformly at random. The reward
277 function includes a penalty term that captures the desired bidding mechanism.

## 278   5   Implementation and evaluation

### 279   5.1   Implementation

280 Talk about:

281 1. different bidding mechanisms

282 2. choice of penalty factor

283 3. action window (remarking that we could additionally allow agents to choose length of action
284     window)

285 4. use with off-the-shelf RL algorithms

### 286   5.2   Environments

### 287   5.2.1   MovingTargetsGridworld

288 Important to mention that we want to maximize $\min(\text{targets reached})$.
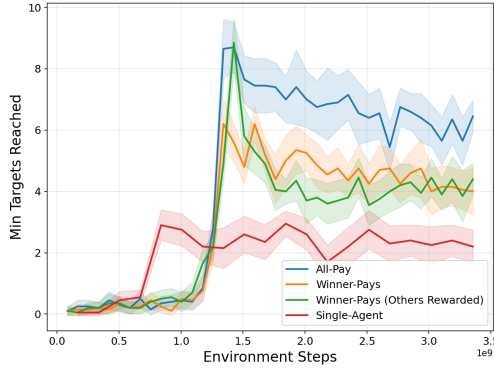
### 289   5.2.2   Atari Assault

### 290   5.3   Baselines

291 1. Weighted sum of rewards with standard RL algorithms

292 2. Deep W learning implemented on top of DQN

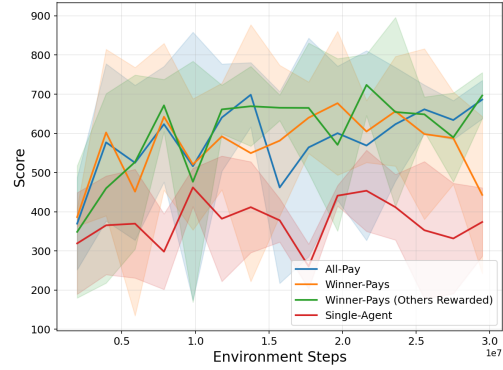Table 1: Performance (mean with 95% CI) averaged over the last 5 evaluation checkpoints.

| Algorithm | Gridworld (Min Targets Reached) | Assault (Score) |
|---|---|---|
| All-Pay | 6.05 [5.74, 6.36] | 634.80 [591.14, 678.46] |
| Winner-Pays | 4.07 [3.78, 4.36] | 578.04 [521.02, 635.06] |
| Winner-Pays (Others Rewarded) | 4.20 [3.92, 4.48] | 662.60 [619.09, 706.11] |
| Single-Agent | 2.31 [2.08, 2.54] | 384.72 [343.09, 426.35] |

## 5.4 Performance comparison with baselines

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



(a) MovingTargetsGridworld. [Placeholder: describe convergence behavior, relative performance of mechanisms, and any notable differences in sample efficiency.]

(b) Atari Assault. [Placeholder: describe convergence behavior, relative performance of mechanisms, and any notable differences in sample efficiency.]

Figure 1: Learning curves for different bidding mechanisms across both environments.

## 5.5 Interpretability

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

## 5.6 Modularity

Plots of performance in gridworld with increasing number of objectives

## 5.7 Ablations

Impact of max bid, penalty factor

# References

Guy Avni, Thomas A Henzinger, and Ventsislav Chonev. Infinite-duration bidding games. *Journal of the ACM (JACM)*, 66(4):1–29, 2019.

Guy Avni, Kaushik Mallik, and Suman Sadhukhan. Auction-based scheduling. In *International Conference on Tools and Algorithms for the Construction and Analysis of Systems*, pp. 153–172. Springer, 2024.
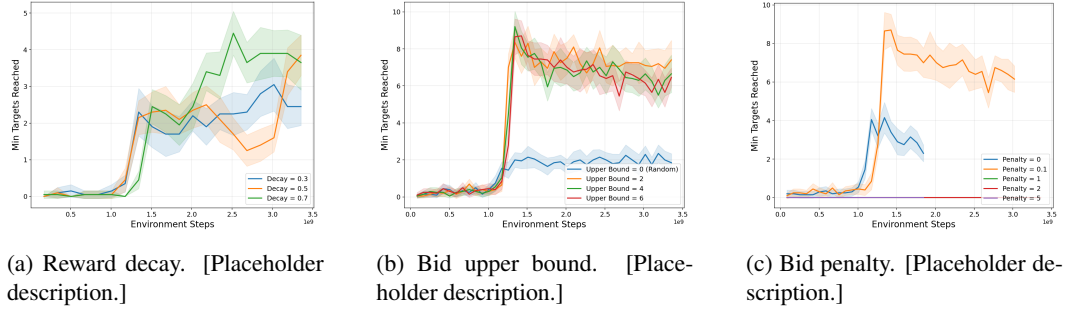
(a) Reward decay. [Placeholder description.]

(b) Bid upper bound. [Placeholder description.]

(c) Bid penalty. [Placeholder description.]

Figure 2: Ablation studies on the MovingTargetsGridworld environment.

308 Guy Avni, Martin Kurečka, Kaushik Mallik, Petr Novotný, and Suman Sadhukhan. Bidding games
309     on markov decision processes with quantitative reachability objectives. In *Proceedings of the 24th*
310     *International Conference on Autonomous Agents and Multiagent Systems*, pp. 161–169, 2025.

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

312 Lucian Buşoniu, Robert Babuška, and Bart De Schutter. Multi-agent reinforcement learning: An
313     overview. *Innovations in multi-agent systems and applications-1*, pp. 183–221, 2010.

314 Woohyeon Byeon, Giseung Park, Jongseong Chae, Amir Leshem, and Youngchul Sung. Multi-
315     objective reinforcement learning with max-min criterion: A game-theoretic approach. *arXiv*
316     *preprint arXiv:2510.20235*, 2025.

317 Saul Gass and Thomas Saaty. The computational algorithm for the parametric objective function.
318     *Naval research logistics quarterly*, 2(1-2):39–45, 1955.

319 Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane,
320     Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz,
321     et al. A practical guide to multi-objective reinforcement learning and planning: Cf hayes et al.
322     *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.

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

324 Andrew J Lazarus, Daniel E Loeb, James G Propp, Walter R Stromquist, and Daniel H Ullman.
325     Combinatorial games under auction play. *Games and Economic Behavior*, 27(2):229–264, 1999.

326 Beatriz M Méndez-Hernández, Erick D Rodríguez-Bazan, Yailen Martinez-Jimenez, Pieter Libin,
327     and Ann Nowé. A multi-objective reinforcement learning algorithm for jssp. In *International*
328     *Conference on Artificial Neural Networks*, pp. 567–584. Springer, 2019.

329 Thanh Thi Nguyen, Ngoc Duy Nguyen, Peter Vamplew, Saeid Nahavandi, Richard Dazeley, and
330     Chee Peng Lim. A multi-objective deep reinforcement learning framework. *Engineering Appli-*
331     *cations of Artificial Intelligence*, 96:103915, 2020.

332 Giseung Park, Woohyeon Byeon, Seongmin Kim, Elad Havakuk, Amir Leshem, and Youngchul
333     Sung. The max-min formulation of multi-objective reinforcement learning: From theory to a
334     model-free algorithm. *arXiv preprint arXiv:2406.07826*, 2024.

335 Matteo Pirotta, Simone Parisi, and Marcello Restelli. Multi-objective reinforcement learning with
336     continuous pareto frontier approximation. In *Proceedings of the AAAI conference on artificial*
337     *intelligence*, volume 29, 2015.

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

Stuart J Russell and Andrew Zimdars. Q-decomposition for reinforcement learning agents. In *Proceedings of the 20th international conference on machine learning (ICML-03)*, pp. 656–663, 2003.

Umer Siddique, Paul Weng, and Matthieu Zimmer. Learning fair policies in multi-objective (deep) reinforcement learning with average and discounted rewards. In *International Conference on Machine Learning*, pp. 8905–8915. PMLR, 2020.

Kristof Van Moffaert and Ann Nowé. Multi-objective reinforcement learning using sets of pareto dominating policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.

Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforcement learning: Novel design techniques. In *2013 IEEE symposium on adaptive dynamic programming and reinforcement learning (ADPRL)*, pp. 191–199. IEEE, 2013.

Kristof Van Moffaert, Tim Brys, Arjun Chandra, Lukas Esterle, Peter R Lewis, and Ann Nowé. A novel adaptive weight selection algorithm for multi-objective multi-agent reinforcement learning. In *2014 International joint conference on neural networks (IJCNN)*, pp. 2306–2314. IEEE, 2014.

Harm Van Seijen, Mehdi Fatemi, Joshua Romoff, Romain Laroche, Tavian Barnes, and Jeffrey Tsang. Hybrid reward architecture for reinforcement learning. *Advances in neural information processing systems*, 30, 2017.