

# Bidding-Based Policy Composition for Dynamic Task Streams

**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, consectetur adipiscing elit. Ut purus elit, vestibulum ut, placerat ac, adipiscing vitae, felis. Curabitur dictum gravida mauris. Nam arcu libero, nonummy eget, consectetur 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 (?).

# Bidding-Based Policy Composition for Dynamic Task Streams

**Anonymous authors**

Paper under double-blind review

## Abstract

1 We study multi-objective reinforcement learning settings in which objectives appear  
 2 or disappear at runtime. We propose a modular framework where each objective is  
 3 supported by a selfish local policy, and coordination is achieved through a novel *auc-*  
 4 *tion*-based mechanism: policies bid for the right to execute their actions, with bids  
 5 reflecting the urgency of the current state. The highest bidder selects the action, en-  
 6 abling a dynamic and interpretable trade-off among objectives. To succeed, each policy  
 7 must not only optimize its own objective, but also reason about the presence of other  
 8 goals and learn to produce calibrated bids that reflect relative priority. When objec-  
 9 tives change, the system adapts by simply adding or removing the corresponding poli-  
 10 cies. Moreover, when objectives arise from the same parameterized family—like the  
 11 class of reachability objectives parameterized by target states—identical copies of a pa-  
 12 rameterized policy can be deployed. In our implementation, the policies are trained  
 13 concurrently using proximal policy optimization (PPO). We evaluate on Atari Assault  
 14 and a gridworld-based path-planning task with dynamic targets. Our method achieves  
 15 substantially better performance and reduced sample complexity than a single policy  
 16 trained with PPO.

## 17 1 Introduction

18 A majority of real-world control problems require fulfilling more than one objectives simultaneously,  
 19 where the objectives could be partly contradictory to each other. We consider the setting where  
 20 objectives appear or disappear at deployment time, and no prior information about their arrival is  
 21 available. Our running example is a coffee serving robot in an office building, where new coffee  
 22 requests could appear from any place at any time, old requests could disappear even before being  
 23 served, and the robot must react to updated objectives as quickly as possible. In general, any number  
 24 of requests could be active at any time, and we do not assume any prior knowledge about their  
 25 distribution. Our goal is to design lightweight, adaptive policies that update their behavior in light  
 26 of the evolving objectives. Despite having been studied in the planning literature (?), this class of  
 27 problems lack support in (model-free) reinforcement learning (RL), and the existing solutions from  
 28 the planning setting do not readily extend.

29 We propose a novel policy adaptation algorithm for this class of problems. The heart of our ap-  
 30 proach is a compositional design of policies for multi-objective RL problems, where each objective  
 31 is served using an independent local policy, and coordination is achieved via an online auction-based  
 32 mechanism: policies bid for the right to execute their actions, with bids reflecting the urgency of the  
 33 current state for fulfilling the respective objectives. The highest bidder selects the action, enabling a  
 34 dynamic and interpretable trade-off among objectives. For example, imagine a situation that the cof-  
 35 fee serving robot is approaching an intersection of two corridors, there are two active requests and  
 36 two independent local policies for serving them, the first policy needs to go left, while the second  
 37 policy needs to go right. Clearly a trade-off is unavoidable, and which of the policies is prioritized  
 38 is decided based on auctions.

39 Given this compositional design, every time an objective appears or disappears, only the respon-  
 40 sible policy needs to be added or removed, eliminating the need for modifying the whole system.  
 41 Moreover, when the objectives come from the same parameterized family, like the set of all possible  
 42 coffee requests parameterized by the request locations, we can design a single universal policy for  
 43 this entire parameterized family, so that new additions require adding an identical copy of the policy,  
 44 offer us instant adaptation. In contrast, a monolithic policy that would serve all objectives lacks this  
 45 flexibility, and an attempt to train such a policy with variable number of objectives causes a dramatic  
 46 degradation of training stability as well as the performance (the loss); we demonstrate this in our  
 47 experiments.

48 We show how classical RL policies can be extended with bidding capabilities. While the obvious  
 49 first step is to extend the action space with numeric bid values, there are three key challenges:

50 **Challenge I: enforcing honest bids.** We must ensure that policies bid only in proportion to their  
 51 urgency, because otherwise, there is a risk of obtaining “dishonest” policies that constantly try to  
 52 block others by overbidding. To circumvent this issue, during training, we require the policies to  
 53 pay penalties proportional to their bid values. **KM: I suggest, we just stick to one kind of bidding,**  
 54 **instead of three, if in the end all-pay remains superior in all experiments.** This way, it is against the  
 55 interest of policies to bid too high, because otherwise, the net benefit of achieving their objectives  
 56 would be lost.

57 **Challenge II: achieving environment awareness.** The bidding tactic of policies should not only  
 58 depend on the current state and the own objective, but also needs to account for other objectives for  
 59 maximal effectiveness. For instance, if two coffee requests appear at nearby locations, the respective  
 60 policies need not bid too high to compete against the opponent. Dually, if two requests arrive from  
 61 opposite directions, the bid of the more urgent policy must be high enough to counteract the oppo-  
 62 nent. Therefore, the policies must account for the objectives of opponents to bid effectively. Our  
 63 approach is to model the local policy synthesis problem as an instance of multi-agent nonzerosum  
 64 games, where each agent is responsible for an individual objective, and needs to learn a policy that  
 65 fulfills its objective against the opponents.

66 **Challenge III: unbounded objective count.** Since the number of objectives seen during deploy-  
 67 ment cannot be predicted during training, the game-theoretic approach for environment-aware poli-  
 68 cies falls apart, because now we face a variable number of opponents for which no solution is known.  
 69 We use an attention network that transforms an arbitrary number of opponent objectives to a fixed  
 70 encoding, which is then fed as the input to the RL policy. Each local policy is equipped with its  
 71 own attention network, which is co-trained with the actual policy for an increasing number of op-  
 72 ponent objectives. **KM: If we want to make the unbounded objectives a central feature, we should**  
 73 **restrict ourselves to objectives from the same family. Otherwise, the attention pooling and the game-**  
 74 **theoretic training approach does not work. While this is indeed a restrictive case, it is well-motivated**  
 75 **and I do not see it as a weakness.**

76 We implemented our framework using the proximal policy optimization (PPO) as the base learn-  
 77 ing algorithm. Using the Atari Assault and a gridworld-based path-planning task, we demonstrate  
 78 the superior training stability and the performance of our policies in comparison with the baseline  
 79 monolithic policies. **TODO: Summarize the main take-away: can we give some statistics, like “the**  
 80 **training time was 2x faster” or the “loss was 2x smaller”?**

## 81 1.1 Related work

82 A large body of work in multi-objective reinforcement learning (MORL) relies on *scalarization*,  
 83 aggregating multiple reward functions into a single scalar objective so that standard single-objective  
 84 RL algorithms can be applied. The simplest scalarization method is a weighted sum of individual  
 85 rewards (Gass & Saaty, 1955), though richer nonlinear scalarization functions have also been pro-  
 86 posed (Van Moffaert et al., 2013). A key limitation of scalarization is that the relative importance  
 87 induced by the aggregation function may not align with the designer’s true intent. This mismatch

can initiate a tedious debugging cycle, particularly in large-scale systems (Hayes et al., 2022). In contrast, our approach achieves a trade-off between reward components without collapsing them into a fixed scalar objective.

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

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

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

## 2 Preliminaries: Multi-Objective MDPs

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

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

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

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

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 state-action pairs and the current state to a distribution over actions. A *path* on  $\mathcal{M}$  induced by  $\pi$  is a 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 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 probability of moving to the state  $s_{t+1}$  from  $s_t$  due to action  $a_t$  is positive, i.e.,  $T(s_t, a_t)(s_{t+1}) > 0$ . A path can be either finite or infinite, and we will write  $Paths(\mathcal{M}, \pi)$  to denote the set of all infinite paths for  $\mathcal{M}$  induced by  $\pi$ . Given a finite path  $\rho = (s_0, a_0) \dots (s_t, a_t)$ , the probability that  $\rho$  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 a probability measure over the set of all infinite paths in  $\mathcal{M}$  using standard constructions, which can be found in the literature (Baier & Katoen, 2008). Given a measurable set of paths  $\Omega$  and a function  $f : Paths(\mathcal{M}, \pi) \rightarrow \mathbb{R}$ , we will write  $\mathbb{P}^{\mathcal{M}, \pi}[\Omega]$  and  $\mathbb{E}^{\mathcal{M}, \pi}[f]$  to denote, respectively, the probability measure of  $\Omega$  and the expected value of  $f$  evaluated over random infinite paths.

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

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

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

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

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

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

### 3.1 The Framework

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

**Loser-Rewarded:** the winning policy pays the bidding penalty and the losing policies earn bidding 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  $\Delta \in \{\text{Loser-Rewarded}, \text{Winner-Pays}, \text{All-Pay}\}$ , 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,
- $\mathbf{A} = \{A^1, \dots, A^m\}$  is a set with  $A^i$  being the set of Agent  $i$ ’s actions,



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

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

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

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

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

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

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

- 256 •  $\widetilde{T}(s, b^1, \dots, b^m) := \text{Uniform}(\{(s, 0, i)\}_{i \in I})$ ,
- 257 •  $\widetilde{R}^i(s, b^1, \dots, b^m, (s, 0, i^*)) := \begin{cases} -b^i & i = i^* \vee \Delta = \text{All-Pay}, \\ +b^i & i \neq i^* \wedge \Delta = \text{Loser-Rewarded}, \\ 0 & i \neq i^* \wedge \Delta = \text{Winner-Pays}, \end{cases}$

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

- 259 •  $\widetilde{T}((s, t, i^*), a^1, \dots, a^m) := \begin{cases} T(s, a^{i^*}) \times ((t+1) \bmod \tau) \times \{i^*\} & t < \tau-1, \\ T(s, a^{i^*}) & t = \tau-1, \end{cases}$
- 260 •  $\widetilde{R}^i((s, t, i^*), a^1, \dots, a^m, (s', t+1, i^*)) := R^i(s, a^i, s')$ .

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

### 3.3 Flavors of Cooperation through Bidding

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

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

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

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

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

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

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

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

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

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

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

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

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

- $S$  is the set of states,



- 307 •  $A$  is the set of actions,
  - 308 •  $T : S \times A \rightarrow \mathbb{D}(S)$  is the transition function mapping a state-action pair to a distribution over
  - 309 states,
  - 310 •  $R : S \times A \times S \rightarrow \mathbb{R}^m$  is the reward function with each output component corresponding to the
  - 311 different objectives, and
  - 312 •  $\mu_0 \in \mathbb{D}(S)$  is the initial state distribution.
- 313 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
- 314 state-action pairs and the current state to a distribution over actions.
- 315 **Definition 4 (MAB-MDP).** Let  $\mathcal{M} = (S, A, T, R, \mu_0)$  be an MO-MDP with  $m$  objectives and
- 316 let  $b \in \mathbb{Z}_{>0}$  be the bid upper bound. Also, define  $M = \{1, \dots, m\}$  be indices of the  $m$  agents
- 317 corresponding to the  $m$  objectives along with  $\perp$  representing a null agent. Lastly, let  $B = \{0, \dots, b\}$
- 318 be the range of bids and  $\rho > 0$  be the bid penalty factor. We define the multi-agent bidding Markov
- 319 decision process (MAB-MDP) as a tuple  $\mathcal{B}_{\mathcal{M}} = (\hat{S}, \hat{A}, \hat{T}, P, \hat{R}, \hat{\mu}_0)$  where
- 320 •  $\hat{S} = M \times S$  is the new state space augmented with the index of the agent that won the previous
  - 321 round of bidding,
  - 322 •  $\hat{A} = A^m \times B^m$  represents the action space of the  $m$  agents in which each agent selects an action
  - 323 from  $A$  and a bid from  $B$ ,
  - 324 •  $\hat{T} : \hat{S} \times \hat{A} \rightarrow \mathbb{D}(\hat{S})$  is the new transition function defined as,

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

- 325 where  $B_{\max} := \{i \mid b_i = \max\{b_1, \dots, b_m\}\}$  is the set of agents with maximal bids. The tuple
- 326  $(T(s, a_i), i)$  represents the distribution over  $\hat{S}$  induced by the original transition function  $T$  such
- 327 that the second component is fixed, and the weighted sum represents taking the weighted sums of
- 328 the distributions over  $\hat{S}$ .
- 329 •  $P : \hat{A} \times M \rightarrow \mathbb{R}^m$  is the bidding penalty for the  $m$  agents and the second component is the index
  - 330 of the agent that won the bidding.
  - 331 •  $\hat{R} : \hat{S} \times \hat{A} \times \hat{S} \rightarrow \mathbb{R}^m$  is the reward function for the  $m$  agents with

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

- 332 where  $i \in M$  is the index of the agent that won the bid and chose the action.
- 333 •  $\hat{\mu}_0 := (\mu_0, 1)$  is the initial state distribution over  $\hat{S}$  induced by  $\mu_0$  and the second component is
  - 334 fixed to be 1 without loss of generality.

335 Given an MAB-MDP  $\mathcal{B}_{\mathcal{M}}$ , a *policy* for each agent indexed by  $i \in \{1, \dots, m\}$  takes a similar form:

336  $\pi_i : (\hat{S} \times \hat{A})^* \times \hat{S} \rightarrow \hat{A}$ . Intuitively, a state  $(i, s) \in \hat{S}$  encodes the agent that won the bidding and

337 chose the action to reach  $s$  in the previous step. At each step, each of the agents choose an action and

338 a bid, and an action amongst the set of highest bidders is chosen uniformly at random. The reward

339 function includes a penalty term that captures the desired bidding mechanism.

## 340 5 Implementation and evaluation

### 341 5.1 Implementation

342 Talk about:

- 343 1. different bidding mechanisms
- 344 2. choice of penalty factor

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]

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

4. use with off-the-shelf RL algorithms

## 5.2 Environments

### 5.2.1 MovingTargetsGridworld

Important to mention that we want to maximize min(targets reached).

### 5.2.2 Atari Assault

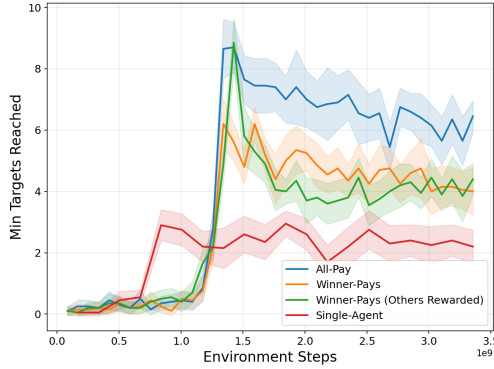
## 5.3 Baselines

1. Weighted sum of rewards with standard RL algorithms

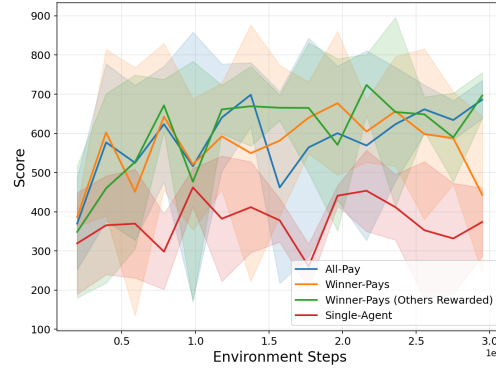
2. Deep W learning implemented on top of DQN

## 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

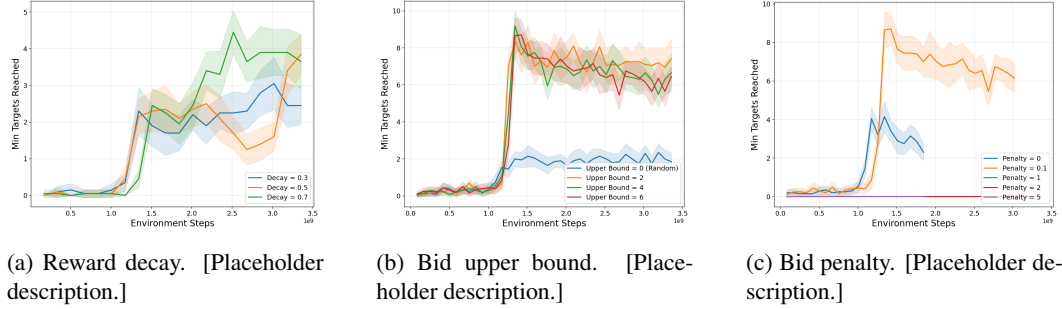


Figure 2: Ablation studies on the MovingTargetsGridworld environment.

## 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.
- Guy Avni, Martin Kurečka, Kaushik Mallik, Petr Novotný, and Suman Sadhukhan. Bidding games on markov decision processes with quantitative reachability objectives. In *Proceedings of the 24th International Conference on Autonomous Agents and Multiagent Systems*, pp. 161–169, 2025.
- Christel Baier and Joost-Pieter Katoen. *Principles of model checking*. MIT press, 2008.
- Lucian Buşoniu, Robert Babuška, and Bart De Schutter. Multi-agent reinforcement learning: An overview. *Innovations in multi-agent systems and applications-1*, pp. 183–221, 2010.
- Woohyeon Byeon, Giseung Park, Jongseong Chae, Amir Leshem, and Youngchul Sung. Multi-objective reinforcement learning with max-min criterion: A game-theoretic approach. *arXiv preprint arXiv:2510.20235*, 2025.
- Saul Gass and Thomas Saaty. The computational algorithm for the parametric objective function. *Naval research logistics quarterly*, 2(1-2):39–45, 1955.
- Conor F Hayes, Roxana Rădulescu, Eugenio Bargiacchi, Johan Källström, Matthew Macfarlane, Mathieu Reymond, Timothy Verstraeten, Luisa M Zintgraf, Richard Dazeley, Fredrik Heintz, et al. A practical guide to multi-objective reinforcement learning and planning: Cf hayes et al. *Autonomous Agents and Multi-Agent Systems*, 36(1):26, 2022.
- Mark Humphrys. W-learning: Competition among selfish q-learners. 1995.
- Andrew J Lazarus, Daniel E Loeb, James G Propp, Walter R Stromquist, and Daniel H Ullman. Combinatorial games under auction play. *Games and Economic Behavior*, 27(2):229–264, 1999.
- Beatriz M Méndez-Hernández, Erick D Rodríguez-Bazan, Yailen Martinez-Jimenez, Pieter Libin, and Ann Nowé. A multi-objective reinforcement learning algorithm for jssp. In *International Conference on Artificial Neural Networks*, pp. 567–584. Springer, 2019.

- 391 Thanh Thi Nguyen, Ngoc Duy Nguyen, Peter Vamplew, Saeid Nahavandi, Richard Dazeley, and  
392 Chee Peng Lim. A multi-objective deep reinforcement learning framework. *Engineering Appli-*  
393 *cations of Artificial Intelligence*, 96:103915, 2020.
- 394 Giseung Park, Woohyeon Byeon, Seongmin Kim, Elad Havakuk, Amir Leshem, and Youngchul  
395 Sung. The max-min formulation of multi-objective reinforcement learning: From theory to a  
396 model-free algorithm. *arXiv preprint arXiv:2406.07826*, 2024.
- 397 Matteo Pirota, Simone Parisi, and Marcello Restelli. Multi-objective reinforcement learning with  
398 continuous pareto frontier approximation. In *Proceedings of the AAAI conference on artificial*  
399 *intelligence*, volume 29, 2015.
- 400 Juan C Rosero, Nicolás Cardozo, and Ivana Dusparic. Multi-objective deep reinforcement learn-  
401 ing optimisation in autonomous systems. In *2024 IEEE International Conference on Autonomic*  
402 *Computing and Self-Organizing Systems Companion (ACSOS-C)*, pp. 97–102. IEEE, 2024.
- 403 Stuart J Russell and Andrew Zimdars. Q-decomposition for reinforcement learning agents. In  
404 *Proceedings of the 20th international conference on machine learning (ICML-03)*, pp. 656–663,  
405 2003.
- 406 Umer Siddique, Paul Weng, and Matthieu Zimmer. Learning fair policies in multi-objective (deep)  
407 reinforcement learning with average and discounted rewards. In *International Conference on*  
408 *Machine Learning*, pp. 8905–8915. PMLR, 2020.
- 409 Kristof Van Moffaert and Ann Nowé. Multi-objective reinforcement learning using sets of pareto  
410 dominating policies. *The Journal of Machine Learning Research*, 15(1):3483–3512, 2014.
- 411 Kristof Van Moffaert, Madalina M Drugan, and Ann Nowé. Scalarized multi-objective reinforce-  
412 ment learning: Novel design techniques. In *2013 IEEE symposium on adaptive dynamic pro-*  
413 *gramming and reinforcement learning (ADPRL)*, pp. 191–199. IEEE, 2013.
- 414 Kristof Van Moffaert, Tim Brys, Arjun Chandra, Lukas Esterle, Peter R Lewis, and Ann Nowé. A  
415 novel adaptive weight selection algorithm for multi-objective multi-agent reinforcement learning.  
416 In *2014 International joint conference on neural networks (IJCNN)*, pp. 2306–2314. IEEE, 2014.
- 417 Harm Van Seijen, Mehdi Fatemi, Joshua Romoff, Romain Laroche, Tavian Barnes, and Jeffrey  
418 Tsang. Hybrid reward architecture for reinforcement learning. *Advances in neural information*  
419 *processing systems*, 30, 2017.