

# 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  $\pi^i$  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.

## 19 2 Preliminaries: Multi-Objective MDPs

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

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

- 30 •  $S$  is the set of states,
- 31 •  $A$  is the set of actions,
- 32 •  $T : S \times A \rightarrow \mathbb{D}(S)$  is the transition function mapping a state-action pair to a distribution over  
33 the successor states,

- 34 •  $\{R^i : S \times A \times S \rightarrow \mathbb{R}_{\geq 0}\}_{i \in [1;m]}$  is the set of reward functions, and  
 35 •  $\mu_0 \in \mathbb{D}(S)$  is the initial state distribution.

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

42 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  
 43 state-action pairs and the current state to a distribution over actions. A *path* on  $\mathcal{M}$  induced by  $\pi$  is a  
 44 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  
 45 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  
 46 probability of moving to the state  $s_{t+1}$  from  $s_t$  due to action  $a_t$  is positive, i.e.,  $T(s_t, a_t)(s_{t+1}) >$   
 47 0. A path can be either finite or infinite, and we will write  $Paths(\mathcal{M}, \pi)$  to denote the set of all  
 48 infinite paths fo  $\mathcal{M}$  induced by  $\pi$ . Given a finite path  $\rho = (s_0, a_0)\dots(s_t, a_t)$ , the probability that  $\rho$   
 49 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  
 50 a probability measure over the set of all infinite paths in  $\mathcal{M}$  using standard constructions, which can  
 51 be found in the literature (Baier & Katoen, 2008). Given a measurable set of paths  $\Omega$  and a function  
 52  $f : Paths(\mathcal{M}, \pi) \rightarrow \mathbb{R}$ , we will write  $\mathbb{P}^{\mathcal{M}, \pi}[\Omega]$  and  $\mathbb{E}^{\mathcal{M}, \pi}[f]$  to denote, respectively, the probability  
 53 measure of  $\Omega$  and the expected value of  $f$  evaluated over random infinite paths.

54 We will use the standard discounted reward objectives, where we fix  $\gamma \in [0, 1]$  as a given discounting  
 55 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  
 56 discounted sum function, mapping  $\rho$  to the discounted sum of the associated rewards:  $f_{ds}^i(\rho) :=$   
 57  $\sum_{t=0}^{\infty} \gamma^t \cdot R^i(s_t, a_t)$ . The *i-value* of the policy  $\rho$  for  $\mathcal{M}$  is the expected value of the discounted sum  
 58 of the *i*-th reward we can secure by executing  $\rho$  on  $\mathcal{M}$ , written as  $val^{\mathcal{M}, i}(\pi) = \mathbb{E}^{\mathcal{M}, \pi}[f_{ds}^i]$ . When  
 59 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  
 60 *value component*.

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

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

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

### 73 3 Auction-Based Compositional RL on Multi-Objective MDPs

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

81 **3.1 The Framework**

82 We present a novel *auction*-based RL framework for compositional policy synthesis for MO-MDPs.  
83 In our framework, not only do the local policies emit actions, but also they *bid* for the privilege  
84 of executing their actions for a given number of time steps  $\tau \in \mathbb{N}_{>0}$  in future. The bids are all  
85 nonnegative real numbers, and the highest bidder's actions get executed for the subsequent  $\tau$  steps,  
86 with ties being resolved uniformly at random. The policy whose actions are executed is referred to  
87 as the *active* policy while the rest are called *idle* policies. To discourage overbidding, we require the  
88 active policy to pay a one-time price—modeled as a negative reward or a *penalty*—equal to its bid  
89 value. This way, it is against the interest of a policy to bid more than the total reward it would earn  
90 if it is active in the next  $\tau$  steps; otherwise, the net earning would be negative, while bidding the  
91 zero amount would secure non-negative earnings. Through bidding, each policy can communicate  
92 the importance for it to execute its actions, and the composition mechanism guarantees that the most  
93 important policy is executed.

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

97 We introduce three different variations of the compositional framework, based on three kinds of  
98 reward systems for the idle policies:

- 99 1. Each idle policy earns a reward of the same amount as its bid, in addition to the default reward  
100 from the transitions in the MO-MDP.
- 101 2. Idle policies only get the default reward specified by the original MO-MDP, i.e., they neither earn  
102 a reward nor pay a penalty associated to the bidding.
- 103 3. Just like the active policy, each idle policy pays a penalty of the same amount as its bid.

104 Each of the three variations have their own benefits and pitfalls, which we will describe subsequently  
105 in Section 3.3.

106 **3.2 Learning Local Policies**

107 We show that computing each individual local policy in our framework reduces to finding the optimal  
108 policy in a stochastic game, for which we could use any standard off-the-shelf learning framework.  
109 In particular, given an MO-MDP and given the objective  $i \in [1; m]$ , we construct a stochastic game  
110 that captures all possible interactions of policy  $i$  against the other policies.

111 Suppose we are given the MO-MDP  $\mathcal{M} = (S, A, T, R, \mu_0)$ .

112 **3.3 A Comparative Study of the Three Variations**113 **4 A Multi-Agent Bidding Approach for Multi-Objective RL**

114 **Definition 2** (MO-MDP). A multi-objective Markov decision process (MO-MDP) with  $m \in \mathbb{Z}_{>0}$   
115 objectives is specified by a tuple  $\mathcal{M} = (S, A, T, R, \mu_0)$ , where

- 116 •  $S$  is the set of states,
- 117 •  $A$  is the set of actions,
- 118 •  $T : S \times A \rightarrow \mathbb{D}(S)$  is the transition function mapping a state-action pair to a distribution over  
119 states,
- 120 •  $R : S \times A \times S \rightarrow \mathbb{R}^m$  is the reward function with each output component corresponding to the  
121 different objectives, and
- 122 •  $\mu_0 \in \mathbb{D}(S)$  is the initial state distribution.

123 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  
 124 state-action pairs and the current state to a distribution over actions.

125 **Definition 3 (MAB-MDP).** Let  $\mathcal{M} = (S, A, T, R, \mu_0)$  be an MO-MDP with  $m$  objectives and  
 126 let  $b \in \mathbb{Z}_{>0}$  be the bid upper bound. Also, define  $M = \{1, \dots, m\}$  be indices of the  $m$  agents  
 127 corresponding to the  $m$  objectives along with  $\perp$  representing a null agent. Lastly, let  $B = \{0, \dots, b\}$   
 128 be the range of bids and  $\rho > 0$  be the bid penalty factor. We define the multi-agent bidding Markov  
 129 decision process (MAB-MDP) as a tuple  $\mathcal{B}_{\mathcal{M}} = (\hat{S}, \hat{A}, \hat{T}, P, \hat{R}, \hat{\mu}_0)$  where

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

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

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

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

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

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

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

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

## 150 References

- 151 Christel Baier and Joost-Pieter Katoen. *Principles of model checking*. MIT press, 2008.
- 152 Mark Humphrys. W-learning: Competition among selfish q-learners. 1995.
- 153 Juan C Rosero, Nicolás Cardozo, and Ivana Dusparic. Multi-objective deep reinforcement learning  
 154 optimisation in autonomous systems. In *2024 IEEE International Conference on Autonomic  
 155 Computing and Self-Organizing Systems Companion (ACSOS-C)*, pp. 97–102. IEEE, 2024.