

Optimizing Empty Container Repositioning and Fleet Deployment via Configurable Semi-POMDPs

Riccardo Poiani, Ciprian Stirbu, Alberto Maria Metelli, Marcello Restelli

Abstract—With the continuous growth of the global economy and markets, resource imbalance has risen to be one of the central issues in real logistic scenarios. In marine transportation, this trade imbalance leads to Empty Container Repositioning (ECR) problems. Once the freight has been delivered from an exporting country to an importing one, the laden will turn into empty containers that need to be repositioned to satisfy new goods requests in exporting countries. In such problems, the performance that any cooperative repositioning policy can achieve strictly depends on the routes that vessels will follow (i.e., fleet deployment). Historically, Operation Research (OR) approaches were proposed to jointly optimize the repositioning policy along with the fleet of vessels. However, the stochasticity of future supply and demand of containers, together with black-box and non-linear constraints that are present within the environment, make these approaches unsuitable for these scenarios. In this paper, we introduce a novel framework, Configurable Semi-POMDPs, to model this type of problems. Furthermore, we provide a two-stage learning algorithm, “Configure & Conquer” (CC), that first configures the environment by finding an approximation of the optimal fleet deployment strategy, and then “conquers” it by learning an ECR policy in this tuned environmental setting. We validate our approach in large and real-world instances of the problem. Our experiments highlight that CC avoids the pitfalls of OR methods and that it is successful at optimizing both the ECR policy and the fleet of vessels, leading to superior performance in world trade environments.

Index Terms—Configurable Environments, Empty Container Repositioning, Fleet Deployment, Reinforcement Learning

I. INTRODUCTION

NOWADAYS marine transportation is crucial for the world’s economy: 80% of the global trade is carried by sea, and most of the world’s marine cargo is transported in containers [1], [2]. In 2004, over 60% of the total amount of goods shipped by sea were containerized, while some routes among economically strong countries are containerized up to 100% [1]. When using such methods to convey goods, problems arise due to the joint combination of the inner nature of container flow with the imbalance of global trade between different regions. For instance, once the freight has been delivered from an exporting country to an importing one, the laden will turn into empty containers that need to be repositioned to satisfy new goods requests in exporting countries. This problem takes the name of Empty Container

Repositioning (ECR) [3]. In ECR, when a vessel arrives at a port to discharge laden, the port has two options: discharge a certain amount of empty containers that are present on the vessel, or load the vessel with empty containers from its stock. The goal is to find *cooperative* repositioning policies that minimize the shortage of demand of empty containers in a given horizon. ECR can be a very costly activity in complex logistic networks and, even if it does not directly generate income, it can account for about 20% of the total costs for shipping companies [3]. Thus, building efficient ECR strategies is a crucial point for real-world logistic scenarios.

As highlighted in previous studies [4], the demand of empty containers that can be satisfied in a given horizon by any ECR policy strictly depends on the routes that the vessels of the given shipping company will follow. From an intuitive point of view, this is clear if we consider a scenario with two routes that have no ports in common. Suppose that most of demand of empty containers concerns ports that are present on the first route. If most of the vessels that the shipping company owns follow the second route, any ECR policy will have poor performance since the cooperation ability of the network to exchange containers is limited by the poor assignments of vessels to routes. In the marine transportation literature, the problem of assigning vessel to routes to maximize some given target function takes the name of Fleet Deployment (FD) [4]. More specifically, in FD, the set of routes is predetermined by the shipping company, and the task is to select the starting port on a route for each vessel that the shipping company owns.

In this paper, we study the adoption of FD techniques to improve ECR policies. Historically, Operation Research (OR) approaches were proposed to jointly optimize the repositioning policy along with the fleet of vessels [5]. However, the stochasticity within the environment together with black-boxed and non-linear constraints that are present in the environment makes them unsuitable for such complex scenarios [2]. Multi-agent Reinforcement Learning (MARL) techniques, on the other hand, have recently achieved success at limiting these problems in ECR settings in which the assignments of vessels to routes is predetermined and given to the learners [2], [6]. They model the problem as a Semi-Partially Observable Markov Decision Process (Semi-POMDP) and propose ad-hoc neural architectures to learn cooperative policies.¹

In this paper, our focus is on a more complex problem, that is *jointly optimizing* the repositioning policy (ECR) together with the fleet of vessels (FD). In this sense, we are jointly

R. Poiani, A.M. Metelli and M. Restelli are with the Dipartimento di Elettronica, Informazione e Bioingegneria, Politecnico di Milano, Milan, Italy. R. Poiani started working on this project while he was an intern at InstaDeep, Paris, France. C. Stirbu works for InstaDeep, Paris, France.

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

¹Notice that the “Semi” component arises from the fact that the ECR problem is intrinsically event-driven: a repositioning action needs to be taken only when a vessel arrives at a given port.

solving the ECR+FD problem. More specifically, from an agents’ perspective, the fleet of vessels can be seen as features of the environment that can be optimized to reach higher performances. In this sense, for single-agent problems, Configurable Markov Decision Processes (Conf-MDPs) [7]–[9] have recently been introduced to extend the Markov Decision Process (MDP) [10] framework to account for environmental configurations. In Conf-MDPs, an *agent* and a *configurator* are responsible for finding the optimal policy-configuration pair. This is clearly related to our application scenario: our agents are in charge of deciding which container repositioning policy to play, whereas the configurator is entitled to select the fleet of vessels. While the early works [7]–[9] focused on the case in which agent and configurator share the same objective, in [11], the setting has been extended to the case in which the configurator and the agent have different (and, possibly, adversarial) goals. Although these approaches have strong theoretical guarantees, how to successfully scale them to more complex domains remains an open question. Indeed, in a *multi-agent cooperative* setting, the dimension of the problem explodes with the number of agents. Moreover, the intrinsic non-stationarity present in multi-agent systems significantly complicates the learning process [12]. We also note that all the previous methods assume the state to be fully observable; in ECR, however, the agents operates under *partial observability*, that introduces an additional challenge.²

The contributions of our works are summarized as follows:

- We introduce the *Configurable Semi-POMDPs* (Conf-Semi-POMDPs), a novel framework whose goal is extending Conf-MDPs to the more complex multi-agent, partially-observable dynamics of the ECR+FD problem (Section IV). In particular, we focus on the cases in which the configuration of the environment (i.e., assignment of vessels to routes) is decided by a central entity (i.e., *configurator*); that, in our case, is the shipping company.
- To solve Conf-Semi-POMDP, we propose a general two-step solution algorithm called “Configure & Conquer” (CC) (Section V). The goal of CC is to build solutions that successfully *scale* the joint optimization process (i.e., policy and configurations) to our large multi-agent systems. CC first optimizes the configurator to output an approximation of the optimal configuration (i.e., fleet deployment) and then “conquers” it by learning a policy in this tuned environmental setting (i.e., the ECR cooperative policy). The main intuition that CC exploits in its configure step is that to compare two distinct configurations one can even leverage suboptimal policies.
- We validate our approach in large and real-world instances of the ECR+FD problem (Section VI). Our experiments show that CC avoids the pitfalls of OR methods and that it is successful at optimizing both the ECR policy and the fleet of vessels, leading to superior performance in world trade environments.

II. PRELIMINARIES

A. Configurable MDPs

A Configurable Markov Decision Process (Conf-MDP) [7] is defined as a tuple $(\mathcal{S}, \mathcal{A}, r, \gamma, \mu, \mathcal{P}, \Pi)$, where \mathcal{S} is the set of states, \mathcal{A} is the set of actions, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function specifying the reward $r(s, a)$ for taking action a in state s , $\gamma \in (0, 1)$ is the discount factor, $\mu \in \Delta(\mathcal{S})$ ³ is the distribution of the initial state, \mathcal{P} and Π are the model and policy spaces respectively. In particular, every $p \in \mathcal{P}$ is a transition function $p : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S})$ that specifies a probability distribution $p(\cdot|s, a)$ over next state upon taking action a in state s , and every $\pi \in \Pi$ is a policy $\pi : \mathcal{S} \rightarrow \Delta(\mathcal{A})$ specifying a probability distribution $\pi(\cdot|s)$ over actions for every state s . In Conf-MDPs, the goal is to find the optimal model-policy pair $(p^*, \pi^*) \in \mathcal{P} \times \Pi$ that maximizes the *expected return*: $\mathcal{J}^{p, \pi} := \mathbb{E} \left[\sum_{t=0}^{+\infty} \gamma^t r(s_t, a_t) | \pi, p, s_0 \sim \mu \right]$, where the expectation is taken w.r.t. the randomness of π , p , and μ . This joint optimization is solved by two cooperating entities: the *agent*, responsible for improving the policy π , and the *configurator*, whose goal is to learn a configuration p .

B. Semi-POMDPs

A Semi-Partially Observable Markov Decision Process (Semi-POMDP) [6] is defined as a tuple $(\mathcal{D}, \mathcal{S}, \mathcal{A}, p, r, \mathcal{O}, o, \gamma, \mu)$, where \mathcal{S} , γ and μ have the same meaning as before. \mathcal{D} is the set of agents, $\mathcal{A} = \mathcal{A}_1 \times \dots \times \mathcal{A}_{|\mathcal{D}|}$ is the set of joint actions the agents can perform, $r : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ is the reward function extended to joint actions, \mathcal{O} is the set of joint observations that are perceived by the agents, $o : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{O})$ is the observation function, that, for every state s and joint action \mathbf{a} provides a probability distribution $o(\cdot|s, \mathbf{a})$ over joint observations. The transition function $p : \mathcal{S} \times \mathcal{A} \rightarrow \Delta(\mathcal{S} \times \mathbb{N})$ provides a probability distribution $p(\cdot, \cdot|s, \mathbf{a})$ over the next state and the time interval k associated to the transition from the current state s to the next one s' . A joint policy π maps a history of observations $\tau = (o_0, \mathbf{a}_0, r_0, o_1, \mathbf{a}_1, r_1, \dots)$ to a distribution over joint actions $\pi(\cdot|\tau)$. The goal consists in finding an optimal policy π^* that maximizes the expected return: $\mathcal{J}^{p, \pi} := \mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^{t_i} r(s_{t_i}, \mathbf{a}_{t_i}) | \pi, p, s_0 \sim \mu \right]$, where $t_0 = 0$ and $t_i = t_{i-1} + k_i$ for $i \geq 1$.

C. Empty Container Repositioning

As highlighted in [2], the ECR problem can be modeled as a graph $\mathcal{G} := (\mathcal{H}, \mathcal{V}, \mathcal{E})$, where $(\mathcal{H}, \mathcal{V}, \mathcal{E})$ are the set of harbor, vessels and routes respectively. More specifically, each harbor $h \in \mathcal{H}$ has a stock of empty containers of maximum capacity C_h . We denote with C_h^t the number of containers available at day t . Each route $e \in \mathcal{E}$ is a directed cycle of consecutive harbors in \mathcal{H} , namely $e_k := (h_1, \dots, h_{|e_k|})$, where $h_1 = h_{|e_k|}$. Routes can intersect with each other. Each vessel $v \in \mathcal{V}$ is associated with a maximum container capacity C_v and a route $e \in \mathcal{E}$. We denote with C_v^t the amount of empty space at day t on vessel v . Finally, we denote with u_v the speed function

²Notice that the optimal repositioning policy is history-dependent.

³We denote with $\Delta(\mathcal{X})$ the set of probability distributions over a set \mathcal{X} .

of vessel v . Given source and destination harbors h_i and h_j , u_v provides a probability distribution over the number of days required by v to reach h_j starting from h_i . Order of goods between ports are described by a stochastic function q . Given two ports $h_i, h_j \in \mathbb{H}$ and a day t , q provides a probability distribution on the number of goods that are requested to be shipped from h_i to h_j at day t . More specifically, h_i can satisfy this demand using stock of empty containers at the previous day (i.e., $C_{h_i}^{t-1}$). Whenever this amount is not enough, a shortage of containers will happen. We denote with L_h^t the total shortage on harbor h at day t . When vessels arrive at harbors, the following happens:

- laden containers for that destination will be discharged to the port; after some days these containers will turn empty containers and will accumulate in the port's stock;
- empty containers can be loaded/discharged on/from the vessel. These are the actions that the ECR policy is responsible for optimizing.

As mentioned in [2], we remark that the behavior with which containers (both full and empty) are loaded/discharged to/from vessels is complex to be modeled. Indeed, this mechanism is subject to black-boxed and non-linear country regulations.

The goal is to find a policy that minimizes the total shortage, namely $\sum_{t,h \in \mathbb{H}} L_h^t$. A full mathematical model is available in Appendix A of [2]. As one can easily verify, ECR can be formalized as a Semi-POMDP. More specifically, we note that agents (i.e., ports) are required to take repositioning actions when vessels arrive. The time that passes in between two subsequent vessel arrivals is non-constant and, in our case, stochastic (i.e., determined by function u_v).

III. RELATED WORKS

Conf-MDPs have been introduced in [7] for finite spaces, and extended in [9] for more complex continuous environments. In these seminal works, the agent is fully responsible for the configuration activity of the environment, which, in turn, results in an auxiliary task to optimize performance. As highlighted in [7], this leads to a clear distinction between Conf-MDPs and multi-task learning [13]. Indeed, in Conf-MDPs, the agent is not interested in learning and gathering experience samples in sub-optimal configuration; its interest is solely toward the optimal policy in the optimal environmental configuration. The configuration activity within the environment, as shown in more recent works [8], [11], can also be carried out by an external entity (i.e., configurator) whose goals can even be adversary w.r.t. the ones of the agent [11]. None of the previous methods, however, have been designed to handle the more complex multi-agent cooperative setting, in which numerous additional challenges are present (e.g., partial observability, highly dimensional states, intrinsic non-stationarity). *Environment design* literature [14], [15] is also related to configurable environments. However, substantial differences are present since these approaches assume that the configurator (interested party) has (partial) access to the agent's best response to a given environment. In addition, the agent's policies, given an environment, are fixed, and, consequently, the optimization process is not joint.

ECR problems [3] were historically solved using OR methods. However, the environment stochasticity, together with the non-linear and black-boxed constraints that are present in the problem, have led researchers to explore MARL solutions [2], [6], [16]. Since the performance that a method can achieve in terms of satisfied demand in a given horizon strictly depends on the routes that the vessels will follow, *network design* and *fleet deployment* techniques have been proposed to jointly optimize the policy along with the fleet of sailing boats [4]. The difference between network design and fleet deployment is that, in network design methods, routes are generated together with assignments of vessel to routes; in FD, instead, the set of routes is fixed and pre-determined. The joint optimization problem has been investigated from an OR point of view in [5]. However, in this more complex case, the pitfalls of OR methods are even amplified by the complexity of the problem. Furthermore, for large networks of vessels and ports, solving the joint problem with a mathematical programming formulation, leads to high computational requirements. For these reasons, hybrid approaches such as [17] have been proposed. In [17] the authors study the network design setting, and propose a method to generate paths for each of the vessels. Each network configuration is evaluated using the value of the objective function of the mixed-integer formulation of the ECR problem in that specific configuration. Then, they use Genetic Algorithms (GA) [18] to optimize for the configuration with the best objective function. This, however, inherits all the pitfalls of the OR approach. As our experiments will show, the plan can diverge from reality, leading to suboptimal solutions.

IV. THE CONFIGURABLE SEMI-POMDP FRAMEWORK

As we have seen, the Conf-MDP framework [7] models scenarios in which a configurator and a single agent cooperate to improve overall performance. In this section, we generalize the formulation to account for the peculiarities of the ECR+FD problem, i.e., the presence of multiple agents, the partial observability, and the semi-Markov property.

Definition 1: A Configurable Semi-POMDP (Conf-Semi-POMDP) is a tuple $(\mathcal{D}, \mathcal{S}, \mathcal{A}, r, \mathcal{O}, o, \gamma, \mu, \mathcal{P}, \Pi)$, where $(\mathcal{D}, \mathcal{S}, \mathcal{A}, r, \mathcal{O}, o, \gamma, \mu)$ is a Semi-POMDP without transition function, and \mathcal{P} and Π are the model and policy spaces.

More precisely, $\Pi = \Pi_1 \times \dots \times \Pi_{|\mathcal{D}|}$ is the set of history-dependent policies that the agents have access to (i.e., the set of ECR repositioning policies). Thus, we can look at the novel Conf-Semi-POMDP framework as either (i) an extension of the Conf-MDP setting to semi-Markov, multi-agent, partially-observable environments or (ii) an extension of the Semi-POMDP to configurable environments in which we have no transition model p , that can indeed be altered as an effect of the environment configuration activity. We focus on the case where \mathcal{P} is a parametric space of transition probability functions. This assumption, which is usual in Configurable MDPs [7], nicely fits the ECR+FD domain, in which each $p \in \mathcal{P}$ encodes assignments of vessels to routes. More specifically, each configuration $p \in \mathcal{P}$ corresponds to $\{(v_i, e_i, h_i)\}_{i=1}^{|\mathcal{V}|}$, where each element (v, e, h) encodes the fact that vessel v follows route e and starting from port h . We also enforce the constraint that h must belong to e .

The performance of a model-policy pair $(p, \pi) \in \mathcal{P} \times \Pi$ is defined via the expected return, as for Semi-POMDPs:

$$\mathcal{J}^{p, \pi} := \mathbb{E} \left[\sum_{i=0}^{+\infty} \gamma^i r(s_{t_i}, \mathbf{a}_{t_i}) \mid s_0 \sim \mu, \pi, p \right], \quad (1)$$

where $t_0 = 0$ and $t_i = t_{i-1} + k_i$ for $i \geq 1$. Thus, the goal, as for Conf-MDPs, consists of finding the optimal model-policy pair $(p^*, \pi^*) \in \mathcal{P} \times \Pi$ such that $\mathcal{J}^{p, \pi}$ is maximized. We denote with π_p^* an optimal policy for a generic configuration $p \in \mathcal{P}$. For ease of notation, whenever it is clear from the context, we drop the specification of the environment in which a policy is run. For instance, $\mathcal{J}^{\pi_p^*}$ measures the performance of pair (p, π_p^*) . In our work, we consider the case in which a central entity (i.e., the shipping company) is responsible for optimizing/taking decisions on the adopted configuration p .

V. CONFIGURE & CONQUER

We now introduce Configure & Conquer, our method to solve Conf-Semi-POMDPs. To appreciate its generality, we first present CC to solve a generic Conf-Semi-POMDP, and then discuss how it works in the ECR+FD domain.

Imagine having an oracle that, given a model $p \in \mathcal{P}$, provides the performance index $\mathcal{J}^{\pi_p^*}$ of the optimal policy π_p^* for that specific environment p . In this case, the original joint optimization problem described in Section IV reduces to:

$$p^* \in \operatorname{argmax}_{p \in \mathcal{P}} \mathcal{J}^{\pi_p^*}. \quad (2)$$

In practice, however, we do not have access to such an oracle. Nevertheless, given a configuration p , it is possible to train an algorithm \mathcal{A}_c of our choice to learn an approximation of the optimal policy $\tilde{\pi}_p^*$, and, consequently, $\mathcal{J}^{\tilde{\pi}_p^*}$. In particular, there might exist sample efficient (yet suboptimal) algorithms that can be used to obtain such approximations. In that case, we can leverage these methods to optimize the empirical version of the objective function $\mathcal{J}^{\tilde{\pi}_p^*}$, where the expectation in $\mathcal{J}^{p, \pi}$ is estimated with trajectories collected within model p using $\tilde{\pi}_p^*$. We can notice that, with this new formulation, the contribution of the configurator to the optimal solution is completely decoupled from the problem of finding the optimal agents' policy. Indeed, when the approximation of p^* is found (i.e., *configure step*), CC optimizes the agents' policy in the tuned environment with a more complex and expensive algorithm \mathcal{A}_* (i.e., *conquer step*), which aims at obtaining better approximations of π_p^* . For this reason, CC is a two-stage optimization algorithm. The general pseudo-code is reported in Algorithm 1.

With respect to the ECR+FD setting, Algorithm 1 evaluates a given assignment p of vessels to routes exploiting a cheaper algorithm \mathcal{A}_c to train a cooperative ECR policy in p . Once this is done, \mathcal{A}_* is used to train the final ECR policy that will be deployed in the approximation of optimal fleet \tilde{p}^* . We now discuss how to choose \mathcal{A}_* and \mathcal{A}_c .

A. Choosing \mathcal{A}_c and \mathcal{A}_* .

The choice of algorithms for computing approximations of the optimal policy in some configuration p depends on

Algorithm 1 Configure & Conquer (CC).

Require: Algorithms \mathcal{A}_* , \mathcal{A}_c , model and policy spaces \mathcal{P} , Π

- 1: Solve $\tilde{p}^* \in \operatorname{argmax}_{p \in \mathcal{P}} \mathcal{J}^{\tilde{\pi}_p^*}$ using \mathcal{A}_c to estimate $\tilde{\pi}_p^*$
 - 2: Solve $\tilde{\pi}_{\tilde{p}^*}^* \in \operatorname{argmax}_{\pi \in \Pi} \mathcal{J}^{\tilde{\pi}_{\tilde{p}^*}^*}$ using \mathcal{A}_*
-

the specific problem at hand. For what concerns \mathcal{A}_c , there is a trade-off between computational/sample efficiency and performance. Indeed, since \mathcal{A}_c is used to train the configurator, the ideal method should be fast to compute and provide good approximations of π_p^* . The main issue is that the faster the method, the more configurations we can evaluate in a reasonable amount of time so to find \tilde{p}^* . However, this usually comes at the cost of precision, which might impact the optimization landscape of $\mathcal{J}^{\tilde{\pi}_p^*}$. For these reasons, depending on the problem, one might use heuristics, mathematical programming, experts, reinforcement learning agents trained on limited data and so on. On the other hand, the choice of \mathcal{A}_* is more critical for the final performance of the system since it is responsible for computing agents' policy that will be actually deployed in \tilde{p}^* . In this sense, the optimal choice for \mathcal{A}_* is the state-of-the-art for the considered industrial setting.

In our experiments, we analyze the performance of CC varying \mathcal{A}_c and \mathcal{A}_* among traditional methods usually employed in ECR domains. More specifically, we consider algorithms that ranges from simple heuristics, to the more complex OR [2] and MARL approaches [2], [6].

B. Configurator Optimization

We have seen how CC decouples the joint optimization process into two subsequent stages: configure and conquer. More specifically, once the configuration step is over, the conquer step consists in simply applying any algorithm of choice in the tuned environmental setting. For this reason, the crucial step of CC consists in finding an approximation of p^* . We now provide an in-depth description of how one can use Reinforcement Learning (RL) to solve this issue.

Suppose w.l.o.g. that D is the dimension of the parametric space \mathcal{P} (i.e., the number of parameters required to define every $p \in \mathcal{P}$). At this point, it is possible to construct an MDP whose sequence of actions will define a configuration $p \in \mathcal{P}$. More specifically, at each timestep the agent chooses a dimension-value pair (d, z) that will assign the value z for dimension d . At time t , the state contains information about the previously selected pairs $\{(d_i, z_i)\}_{i=1}^{t-1}$, and dummy values for the missing configuration parameters. After D timestamps, an entire configuration p is produced and the agent will receive a reward according to $\mathcal{J}^{\tilde{\pi}_p^*}$, where $\tilde{\pi}_p^*$ is the approximation of the optimal policy of algorithm \mathcal{A}_c in model p . For all previous timestamps, the reward is fixed at 0 for any action. Given this formulation, it is easy to show that the expected discounted reward of the fixed initial state (i.e., no value-dimension assignment done) is proportional to $\mathbb{E}_{p \sim \nu_c} [\mathcal{J}^{\tilde{\pi}_p^*}]$, where ν_c is the single-agent configurator policy. At this point, one can use any RL algorithm to train the configurator to solve this MDP. In our experiments, we model ν_c with neural networks, and rely on PPO [19] for its optimization.

We notice that, in the ECR+FD setting, the configurator, at each step, needs to assign a vessel to a route and an initial port; i.e., it selects a triplet (v, e, h) . It follows that, with the most naive implementation of ν_c , the action space would be given by the cartesian product of $H \times V \times E$, which, in practice, can be quite large, thus making the learning process hard, unstable, and inefficient.⁴ Moreover, we expect a good configurator to be able to exploit the inner structure within the parametric space \mathcal{P} , such as similarities between configurations. Imagine two distinct configurations $p_1, p_2 \in \mathcal{P}$ in which the only difference is that a given vessel $\bar{v} \in V$ is assigned to the same route $\bar{e} \in E$ but two different initial ports h_1 and h_2 in \bar{e} . Clearly, we can expect the rewards $\mathcal{J}^{\bar{p}}$ that the configurator will obtain for the two configurations to be very similar. Therefore, to make training efficient and effective, we rely on the following more complex architecture for ν_c . First of all, we reduce significantly the action space size using *autoregressive* policies [20]. More specifically, each action (i.e., triplet (v, e, h)), is split into three components (i.e., v , e and h) that will be sequentially picked one before the next. We first pick the route e , then given the route, we pick the port h , and given the route and the port we pick the vessel v . By doing so the size of the action space is reduced from $|H||V||E|$ to $|H| + |V| + |E|$. Moreover, to exploit similarity between configurations, we enlarge the state of the agent at timestep t with features of the current uncompleted configuration (e.g., number of vessels assigned to each route, total capacity of the vessels assigned to each route). These features are processed by a neural network f to create an embedding of the current state of the agent. This embedding is used to select the first sub-action (i.e., the route), and it is concatenated to previous selected sub-actions to compute the next sub-actions. Moreover, to further exploit structure the parametric space \mathcal{P} , any unfeasible action (e.g., a port that does not exist in a route) is masked.

VI. EXPERIMENTS

A. Experimental Setup

Similar to previous studies [2], [6], our experiments aims at testing our approach in scenarios that mimic dimensions and behaviors of international transportation companies. To this end, we rely on a patched version of the MARO simulator [21]. More specifically, we test our method on two different topologies WWT_1 and WWT_2 . WWT_1 is composed of 46 vessels, 22 ports, and 13 routes (as in [6]); in WWT_2 , the number of ports and vessels is the same, but the number of routes is 6. Moreover, in WWT_1 , we consider an optimization horizon of 400 days, while in WWT_2 , 200 days are considered. For both problems, as in previous works [6], order distributions have complex trigonometric shapes with multiple periods.

Given this real-world inspired setup, our experiments aim at answering the following questions: (i) Can CC find better configurations in which the agents operate? (ii) How does the performance of \mathcal{A}_c and \mathcal{A}_* impact the final results? (iii) How does CC compare to existing algorithms? To this end, we

conduct an extensive empirical study of CC in WWT_1 and WWT_2 , picking as \mathcal{A}_c and \mathcal{A}_* the following algorithms:

- **Random policy (Rand)**. A random repositioning action is taken every-time a vessel arrives at a certain port.
- **Heuristic policy (Heur)**. This is a stochastic intuitive heuristic that we propose to solve ECR problems. If a port is an exporting one (i.e., it exports much more goods w.r.t. the ones that it imports), then, whenever a vessel arrives, we randomly discharge at least parts of the empty containers that it carries. If a port is an importing one (i.e., it imports much more goods than the ones it exports), instead, we randomly load empty containers on the vessel.
- **Operation Research (OR)**. Noisy estimates of future orders and vessel arrivals are used at the beginning of the interaction with the environment to compute a plan by solving the mathematical formulation of the problem (see Appendix in [2]).
- **Operation Research methods with iterative plan (OR(I))**. Noisy estimates of future orders and vessel arrivals are used to solve the mathematical formulation of the ECR problem. The plan is computed for a long horizon but executed only for a short window. Once the window expires, a new plan is recomputed using the current state of the environment, so to prevent the plan to diverge from reality [1], [2].
- **MARL system**. The application of MARL techniques to ECR problems has been studied in several works [2], [6], [16]. In our experiments, we use a variant of [6].

We select a subset of combinations of \mathcal{A}_c and \mathcal{A}_* that highlight our contributions and that successfully answer to the previous questions. The notation that will be used for CC is $CC-\mathcal{A}_c-\mathcal{A}_*$.

Moreover, we compare CC with the following baselines:

- **LS-NET [17]**. LS-NET [17] was originally proposed to tackle the joint problem of *network design* and ECR, however, the extension to the ECR+FD is direct. We define elements of a GA population so that each element describes a configuration $p \in \mathcal{P}$ (i.e., assignments of vessels to routes and initial ports), and we evaluate each of the elements in the population using as fitness function the value of the objective function of the OR formulation of the problem. Once the method has reached convergence, OR(I) is used to evaluate the performance of the approximation of the optimal configuration.
- **Genetic Algorithm - Joint (GA joint)**. GAs are used to jointly optimize the configuration and the agents policy. The agents' policy is represented as a matrix in which cell (i, j) specifies the repositioning j -th action of the i -th vessel. Configurations are represented as assignments of vessels to routes. An element in the GA population is, thus, a concatenation of a policy with a configuration.
- **Random Configuration and OR (RandomConf-OR(I))**. Configurations are generated at random; OR(I) computes the policy on these sampled configurations.

B. Results

Table I reports mean and 95% confidence intervals (5 runs) of the percentage of satisfied demand of different algorithms

⁴In our experiments, we consider real-world instances of 46 ships, 22 ports and 13 routes. This means that the action space would have dimension 13156.

TABLE I
ECR+FD RESULTS (5 RUNS, MEAN \pm 95% C.I.).

Algorithm	WTT_1	WTT_2
CC-Heur-Heur	86.21 \pm 0.15	78.64 \pm 0.27
CC-Heur-OR(I)	90.93 \pm 0.34	90.58 \pm 0.39
CC-Rand-Rand	77.72 \pm 0.26	42.69 \pm 0.11
CC-Rand-OR(I)	87.93 \pm 0.12	87.15 \pm 0.35
CC-OR-OR	88.38 \pm 0.23	86.15 \pm 0.23
CC-OR-MARL	88.77 \pm 0.41	94.88 \pm 0.60
CC-OR-OR(I)	90.95 \pm 0.12	93.03 \pm 0.60
CC-OR(I)-OR(I)	92.28 \pm 0.35	94.52 \pm 0.47
CC-OR(I)-MARL	88.85 \pm 0.21	95.35 \pm 0.88
CC-OR(I)-Rand	74.16 \pm 1.66	39.15 \pm 5.45
GA joint	86.32 \pm 0.18	83.73 \pm 0.90
LS-NET	88.21 \pm 0.54	84.66 \pm 1.70
RandomConf-OR(I)	77.42 \pm 0.18	68.49 \pm 0.56

on WTT_1 and WTT_2 . We highlight in bold the highest performance reached in each domain.

First of all, as long as we choose good algorithms for \mathcal{A}_* (i.e., OR, OR(I), MARL), we can notice that CC is able to find configurations in which the agents operate that are better than random. This is confirmed by the fact that a good method on random configurations (i.e., RandomConf-OR(I)) reaches lower performance w.r.t. cases in which configurations have been tuned using CC. In this sense, the choice of \mathcal{A}_* is the one that most impact the final performance of our two-stage optimization algorithm. This is expected; indeed, even though we find the optimal configuration, if our agents ignore how to behave (see CC-OR(I)-Rand), the performance will be highly sub-optimal (even worse than using a good method on random configurations such as RandConfig-OR(I) does). The more sophisticated \mathcal{A}_* is (e.g., MARL and OR(I)) the better the performance we can obtain. As for the choice of \mathcal{A}_c , instead, most surprisingly we notice that its choice does not affect performance significantly. Indeed, even when the reward for the configurator is computed using random policies, one can still significantly improve the performance. In particular, we notice that in WTT_2 , even if the random policy is highly suboptimal (see CC-Rand-Rand; 42.69%), CC is still able to find configurations that lead to significantly good performance (see CC-Rand-OR(I); 87.15%). Finally, we see that GA joint and LS-NET underperform w.r.t. all the non-ablation versions of CC. For what concerns GA joint, we conjecture that its sub-optimality arises from the difficulty of the joint optimization process. LS-NET, on the other hand, computes the plan beforehand and optimizes the configuration relying on the value of an objective that might not be respected in practice. Indeed, if the value of the objective function of the mathematical formulation (i.e., the fitness function used in the GA optimization steps) diverges from reality, then LS-NET optimizes toward sub-optimal configurations.

VII. CONCLUSIONS

In this work, we have studied the ECR+FD setting under the novel perspective of Configurable Semi-POMDPs. We have modeled the fleet deployment problem as configurations of the environment that can be tuned to increase the overall performance of the multi-agent system. In particular, we focused on the case in which decisions on the fleet of vessels

are taken by a central entity (i.e., the shipping company) that cooperates with the agents to minimize the total shortage of containers. We proposed a novel two-stage optimization algorithm (CC) that, as our experiments show, successfully solve the joint ECR+FD optimization problem in large and real-world inspired problem instances, outperforming competitive baselines.

We remark the generality of the proposed approach. Indeed, a broad number of multi-agent environments have the possibility to be configured to maximize the performance of the system. In this sense, CC represents a viable option that can successfully scale to large and complex domains. Our work represents a further step in the literature of Configurable MDPs, that, so far, have managed to solve problem of much smaller dimensions only. We also notice that CC, can be directly applied as-is in single-agent problems as well.

As future works, we plan to apply the Conf-Semi-POMDPs framework to the more complex ECR+network design setting. In this case, the set of routes is not pre-determined but paths for each of the vessels are computed by the configurator.

REFERENCES

- [1] Y. Long, L. H. Lee, and E. P. Chew, "The sample average approximation method for empty container repositioning with uncertainties," *EJOR*, 2012.
- [2] X. Li, J. Zhang, J. Bian, Y. Tong, and T.-Y. Liu, "A cooperative multi-agent reinforcement learning framework for resource balancing in complex logistics network," in *AAMAS*, 2019, pp. 980–988.
- [3] D.-P. Song and J.-X. Dong, "Empty container repositioning," *Handbook of ocean container transport logistics*, pp. 163–208, 2015.
- [4] N. K. Tran and H.-D. Haasis, "Literature survey of network optimization in container liner shipping," *Flexible Services and Manufacturing Journal*, vol. 27, no. 2, pp. 139–179, 2015.
- [5] S. Wang and Q. Meng, "Container liner fleet deployment: a systematic overview," *TR_C*, 2017.
- [6] W. Shi, X. Wei, J. Zhang, X. Ni, A. Jiang, J. Bian, and T.-Y. Liu, "Cooperative policy learning with pre-trained heterogeneous observation representations," in *AAMAS*, 2021, pp. 1191–1199.
- [7] A. M. Metelli, M. Mutti, and M. Restelli, "Configurable markov decision processes," in *ICML*, 2018, pp. 3491–3500.
- [8] A. M. Metelli, G. Manneschi, and M. Restelli, "Policy space identification in configurable environments," *Machine Learning*, pp. 1–53, 2021.
- [9] A. M. Metelli, E. Ghelfi, and M. Restelli, "Reinforcement learning in configurable continuous environments," in *ICML*, 2019.
- [10] M. L. Puterman, *Markov Decision Processes: Discrete Stochastic Dynamic Programming*, 1st ed., USA, 1994.
- [11] G. Ramponi, A. M. Metelli, A. Concetti, and M. Restelli, "Learning in non-cooperative configurable Markov decision processes," *NeurIPS*, vol. 34, 2021.
- [12] K. Zhang, Z. Yang, and T. Başar, "Multi-agent reinforcement learning: A selective overview of theories and algorithms," *Handbook of Reinforcement Learning and Control*, pp. 321–384, 2021.
- [13] N. Vithayathil Varghese and Q. H. Mahmoud, "A survey of multi-task deep reinforcement learning," *Electronics*, 2020.
- [14] H. Zhang, Y. Chen, and D. C. Parkes, "A general approach to environment design with one agent," in *IJCAI*, 2009.
- [15] C.-J. Ho, Y.-L. Kuo, and J. Y.-j. Hsu, "Multiagent environment design in human computation," in *AAMAS*, 2011, pp. 1279–1280.
- [16] Q. Luo and X. Huang, "Multi-agent reinforcement learning for empty container repositioning," in *ICSESS*. IEEE, 2018.
- [17] K. Takano and M. Arai, "Study on a liner shipping network design considering empty container repositioning," *JASNAOE*, 2011.
- [18] S. Mirjalili, "Genetic algorithm," in *Evolutionary algorithms and neural networks*. Springer, 2019, pp. 43–55.
- [19] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [20] L. Metz, J. Ibarz, N. Jaitly, and J. Davidson, "Discrete sequential prediction of continuous actions for deep rl," *arXiv preprint arXiv:1705.05035*, 2017.

- [21] "Maro: A multi-agent resource optimization platform," 2020. [Online]. Available: <https://github.com/microsoft/maro>