

Deep RL for Multi-Echelon Supply Chains

Anonymous authors
Paper under double-blind review

Keywords: Deep RL in supply chain, PPO, optimization heuristics, imitation learning

Summary

The present article studies RL methods for supply chain optimization, one of the most natural real-world applications of RL. A lot of attempts appeared during the past years in the operations research community. We approach the problem more from the RL point of view. To this end we design an abstraction that covers features of real-life supply chains typical in the process industry. Our abstraction can be implemented as a gymnasium environment to be trained with standard algorithms. It is proposed to combine optimization heuristics from operations research in combination with imitation learning to pre-train the algorithm. We compare experimentally PPO with and without pre-training to the optimization heuristic. In particular, we give a zero-shot comparison to show that deep RL agents generalize better to disruptions in the supply chain.

Contribution(s)

1. This paper proposes an abstraction of typical real-world multi-echelon supply chain problems from the process industry (e.g. the chemical or pharmaceutical industries) in the form of an RL environment. Our supply chain MDP uses order-based actions to be realistic and reduce the complexity of the action-space. We explain how to use action-shaping and -masking in different ways and apply PPO to the problem.
Context: Without getting the attention it deserves, supply chain optimization has always been one of the prime examples for the use of reinforcement learning in real-world. Many recent articles have tackled the problem, mostly from the operations research perspective. A main caveat is often a simplified view on supply chain planning, not capturing real-world restrictions. We try to give an implementable and more realistic MDP formulation that tries to stay close to today's real-life supply chain planning.
2. The paper shows how to use classical optimization heuristics in combination with imitation learning to pre-train deep RL agents. The numerical advantage is shown on typical multi-echelon supply chain problems.
Context: It is always desirable to start deep RL training in reasonably well-trained policies. For the supply chain problem, we do not have access to known pre-trained agents but can use optimization heuristics to create reasonable rollouts that can be fed into imitation learning algorithms to obtain reasonable policies.
3. The paper compares experimentally the deep RL to a classical planning heuristic. We show that deep RL agents can be more robust, in our experiments, they improve heuristics in zero-shot learning on changing demand.
Context: Since the pandemic, the question of supply chain robustness gets a lot of attention. A question of large practical importance is the understanding of robustness of planning algorithms used in suddenly changed environments, such as suddenly increased or decreased demand. We show experimentally that deep RL agents are more robust than typical planning heuristics that are run in standard supply chain planning software.

Deep RL for Multi-Echelon Supply Chains

Anonymous authors

Paper under double-blind review

Abstract

1 The present article studies RL methods for multi-echelon inventory optimization, one
 2 of the most natural real-world applications of RL. A lot of attempts appeared during
 3 the past years in the operations research community; we approach the problem from the
 4 RL point of view. To this end, we design an abstraction that covers features of real-
 5 life supply chains typical in the process industry. Our abstraction can be implemented
 6 as a gymnasium environment to be trained with standard algorithms. We propose to
 7 combine MRP optimization heuristics from operations research in combination with
 8 imitation learning to pre-train the RL algorithms. We compare experimentally PPO
 9 with and without pre-training to the MRP heuristic. In particular, we give a zero-shot
 10 comparison to show that deep RL agents generalize better to disruptions in the supply
 11 chain.

12 1 Introduction

13 The field of deep RL has seen remarkable success in various fields. Breakthroughs range from
 14 games [Silver et al. \(2016\)](#), over the optimization of fusion reactors [Degraeve et al. \(2022\)](#), to various
 15 topics in the training of LLMs. A field that has seen relatively little progress given its tremendous
 16 industry importance, is multi-echelon (also called multi-level) optimization, see e.g. [Barbosa-Povoa
 17 et al. \(2017\)](#). In supply chain optimization, decisions are made (in real life by teams of supply chain
 18 planners) that range from the procurement of materials, over manufacturing steps, to the distribution
 19 in a logistic network. In this article, we are interested in problems that are typical for the process
 20 manufacturing industry (e.g. the chemical or pharmaceutical industries). Process industry have the
 21 challenge of manufacturing ingredients via multiple processing, with long lead times, shared and
 22 constraint resources, with typically global supply via 1-2 plants.

23 A typical real-world situation is displayed in Figure 1.
 24 The example has a single factory and a number of ware-
 25 houses distributed over several continents. Customer de-
 26 mand is fulfilled only by gray warehouses, while white
 27 warehouses are inter-company distribution centers. The
 28 lead times in replenishment vary a lot depending on inter-
 29 continental or local replenishment.

30 The topic we are interested in is *multi-echelon inventory
 31 optimization (MEIO)*. How to make optimal manufactur-
 32 ing/replenishment decisions that fulfill customer demand
 33 at minimal cost (storage, shipment, and backorder costs).
 34 In this article, we introduce an order-based RL framework
 35 that imitates real-world supply chain planning by placing
 36 orders into an order book. Our framework is implemented
 37 as a gymnasium environment to be trained with standard
 38 deep RL algorithms. It turns out that training the environments is delicate, fast PPO training requires
 39 a reasonable initial policy to avoid metastability effects. We introduce a new trick to deep RL train-
 40 ing in supply chain optimization by combining deep RL and standard optimization heuristics. Using

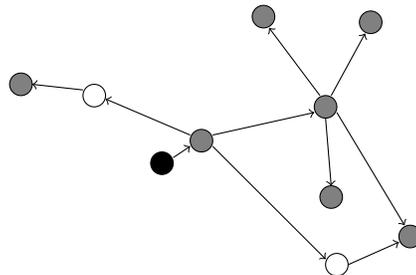


Figure 1: Real-world logistics network over continents. Black is a factory, gray warehouses have customer demand, and white intercompany distribution centers.

41 some heuristic (good or bad), we use imitation learning to imitate a policy network that is used to
42 initiate the PPO training. We show experimentally that pre-training through imitation improves the
43 training a lot. Previous work on deep RL in supply chain optimization has shown that beating known
44 rule-based optimization heuristics is not easy, see e.g. [Gijsbrechts et al. \(2022\)](#). This is similar for
45 our order-based MDP problem. Still, learning (near) optimal policy networks has the big advantage
46 that the evaluation is much cheaper compared to rule-based heuristics and additionally one can ex-
47 pect generalization effects of the neural networks. The latter might turn out powerful for future deep
48 RL applications in supply chain problems as recent years have shown that disruptions in the supply
49 chain cause major challenges, more robust supply chain planning is key.

50 **Related work:** Despite the challenges, industry solutions and a huge body of operations research
51 literature give approximate solution methods to various optimization problems in inventory man-
52 agement, production planning, and logistics operations. Problems often fall into the category of
53 NP-hard problems, where finding optimal solutions becomes computationally infeasible as the scale
54 and intricacy of the supply chain increase ([Gayon et al., 2017](#)). Consequently, traditional systems
55 rely on heuristic algorithms to provide approximate solutions ([Marklund & Rosling, 2012](#); [Zhao &
56 Zhang, 2020](#)). While many methods offer computational efficiency and simplicity, they are limited
57 by their simple nature, they tend to struggle to solve the optimization problems if the complexity
58 is high, and they do not generalize to changing conditions. There is a growing body of research
59 towards deep learning and RL for supply chain optimization. Q -learning was used, for instance, for
60 the beer game ([Oroojlooyjadid et al., 2022](#)), PPO was used, for instance, on a small toy problem in
61 [Vanvuchelen et al. \(2020\)](#). Inventory management for a single-node problem has been studied in [Qi
62 et al. \(2023\)](#). [Hubbs et al. \(2020\)](#) provide a gym environment for simpler supply chain problems.
63 [Perez et al. \(2021\)](#) trained PPO and compared it to perfect information oracles. A comparison of
64 A3C with classical operations research methods has been provided in [Gijsbrechts et al. \(2022\)](#).

65 In contrast to the present articles, all these articles have in common that their focus is less on the
66 specifics of RL training on supply chain problems but much more on the comparison to different
67 heuristics from operations research. For a recent article that contributed to the RL specifics of
68 (rather different) supply chain problems we refer to [Madeka et al. \(2022\)](#).

69 **Main contributions:**

- 70 • **Modeling:** An order-based MDP formalism is set up with main features of industry supply chains.
- 71 • **RL techniques:** Action-shaping and -masking are used to deal with huge action-spaces. Imitation
72 learning is used to leverage simple optimization heuristics for PPO training.
- 73 • **Experiments:** Different learning strategies are compared experimentally for learning efficiency,
74 cost vs. backorder performance, and zero-shot generalization to demand changes.

75 **2 Modeling order-based supply chain environments**

76 As a first approximation for a multi-echelon supply chain problem, the reader might want to think
77 of a multi-graph G in which nodes represent material-storage-locations (different materials at the
78 same location are different nodes) and edges their abstract relations. In a logistics network for one
79 material, the nodes might just represent warehouses, in manufacturing, nodes might also be different
80 materials at the same location. Our modeling allows both. A relation could be a shipment from one
81 location to another, but also a production step involving multiple materials manufactured into a new
82 material at the same location. With regards to real-world supply chain planning we distinguish
83 between *procurement*, buying materials from suppliers, *manufacturing*, producing new materials
84 from other materials, and *replenishment*, transporting materials between storage locations. To have
85 a joint modeling of actions we introduce the concept of *value creation* objects. These are supply
86 chain planning steps that add value to the supply chain. A value creation is based on a labeled finite
87 weighted sub-graph of the supply chain graph G and performs a typical supply chain step. This can
88 be a shipment that only involves two nodes connected by an edge, or a production step that involves
89 a sub-tree of G , where several materials are manufactured into a new material. Value creations v
90 are executed through orders. If the order is feasible, the value creation specifies the quantities of

91 materials to be shipped or produced. Quantities are measured in integer multiples of a fixed lot size
 92 L_v . Each value creation has an associated (random) lead time τ_v , the time between placing the order
 93 and receiving the material. We assume that materials are consumed immediately at the source nodes,
 94 become available at the destination node after the lead time τ_v at a cost c_v per lot size. Required real-
 95 world information needed for this setup is available in ERP (enterprise resource planning) systems,
 96 more precisely from master data such as BOM (bill-of-materials) files or from historic data. Each
 97 node n maintains an inventory level $I_{n,t}$ at time t , representing the quantity of stored materials. For
 98 demand nodes (e.g. customers can buy at these locations), random customer demand $d_{n,t}$ specifies
 99 the demand at time t . When $I_{n,t} \geq d_{n,t}$, demand is fully met and inventory decreases accordingly.
 100 If $d_{n,t} > I_{n,t}$, the unmet amount $B_{n,t}$ is called a backorder. Both inventory levels and backorders
 101 are dynamically updated based on sequential order placements that trigger value creations.

102 **Supply chain environment:** To create an MDP that models the order-based decision making in
 103 supply chain planning a number of quantities need to be defined. In the reality, most data from the
 104 table are either known from master data or can be estimated from historic data.

Symbol	Description
N	Set of all material-storage-locations.
V	Set of value creations, a set of weighted subgraphs of the complete graph with N nodes. The edge weights determine the portions of each ingredient.
d_v	Destination of value creation v , the node that receives the material
s_v	Set of sources of value creation v , nodes that contributed to the value creation v
c_n^h	Storage cost per time unit and unit of material stored at node n .
c_n^b	Backorder cost per time unit and unit of material stored at node n .
w_n	Maximum inventory capacity at node n .
\tilde{w}_n	Maximum allowed backorder at node n .
c_v	Cost per lot for orders of a value creation v .
q_v	Lot size for value creation v .
$g_{v,n}$	Input material quantity from node n per one lot size of value creation v .
k_v^{\max}	Maximum number of lots per order for a value creation v .
\mathcal{D}_n	Demand distribution for material at node n .
\mathcal{Q}_v	Lead time distribution on value creation v .
$I_{n,t}$	Inventory level at node n at time t .
$B_{n,t}$	Backorder level at node n at time t .
$D_{n,t}$	Demand for material at node n at time t (distributed according to \mathcal{D}_n).
L_v	Number of lots created in value creation v .
\mathcal{O}	Set of orders $o = (v, t, \tau, L)$, where v is a value creation, t the placement time, τ the lead time (distributed according to \mathcal{Q}_v), and L the number of lot sizes.

Table 1: Definitions for constants, distributions and variables.



Figure 2: Schematic description of value creations. Left: Replenishment v with lot size q_v , cost c_v , maximal replenishment capacity k_v^{\max} . Right: Manufacturing step using materials s_1, s_2 requiring quantities q_{v,s_1}, q_{v,s_2} to create a lot of size q_v at cost c_v . k_v^{\max} lots can be processed at once.

105 We define the problem as a discounted Markov Decision Process (MDP) with states, actions, and
 106 rewards. An MDP is defined by a tuple $\mathcal{M} = \langle \mathcal{S}, \mathcal{A}, P, R, \gamma \rangle$, where \mathcal{S} is the set of states, \mathcal{A} is the
 107 set of actions, P is the transition probability function, R is the reward function, and γ is the discount

108 factor. The goal is to maximize the expected discounted total reward $\max_{\pi} \mathbb{E}_{\pi} [\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$
 109 over all policies (kernels on $\mathcal{A} \times \mathcal{S}$). Defining an MDP for supply chain problems is non-trivial, in
 110 particular caused by the complexity of allowed actions. An example to take into account is a planner
 111 that can only ship as much stock as is available.

112 **States:** The state-space is $\mathcal{S} = \mathbb{R}^{|N|} \times \mathbb{R}^{|N|} \times O^{\infty}$, where we denote by O^{∞} vectors of orders
 113 (triplets) of arbitrary finite length. States represent inventory levels and backorder levels for each
 114 node n and the current order book. The order book contains a finite number of tuples with start time,
 115 lead time, and number of lots demanded for a value creation.

116 **All actions:** The formal action-space is $\mathcal{A} = \mathbb{N}^V$. A finite action vector a represents the lot sizes of
 117 all action creations asked by the RL planner to be placed. If an action a is performed at time t , then
 118 new orders (v, t, τ, L) will be placed in the order book and used to update the environment.

119 **Feasible actions:** Not all actions are allowed, actions are only allowed when they are possible. For
 120 example, a shipment order can only be placed if material is available at the source. The feasible set
 121 of actions is defined by three constraints: resource capacity, order capacity and material availability.

- 122 • *Available manufacturing capacity:* Manufacturing orders can only be placed if the previous order
 123 o of the same value creation is processed, i.e. $t_o + \tau_o < t$. There is no such constraint for
 124 procurement and replenishment, new shipment orders are always possible.
- 125 • *Order capacity:* There is a maximal number k_v^{\max} of lots that can be manufact-
 126 ured/replenished/procured. For replenishment e.g. the maximal number of containers on a vessel.
- *Material availability:* The sum over all value creations placed by the action a must satisfy

$$\sum g_{v,n} L_v \leq I_n$$

127 for all nodes n . The constraint asks that all orders added to the orderbook can be satisfied imme-
 128 diately from the current inventory.

129 **Environment dynamics:** The environment runs in discrete time-steps by fulfilling orders from the
 130 order book O that is filled by the actions. To see how states transition we need to identify new
 131 inventory and backorder values, as well as the change in the order book. Changes in the order book
 132 are simple. New orders are included, finished orders (i.e. $t_o + \tau_o < t$) are removed. To describe the
 133 inventory and backorder update, we keep track of a single inventory variable (instead of inventory
 134 and backorder separately) that can have negative values, called the generalized inventory level. The
 135 generalized inventory level at time t at node n is denoted by $G_{n,t} \in \mathbb{R}$. This is equal to either the
 136 inventory level $I_{n,t}$ when it is positive or to the negative backorder $-B_{n,t}$ when it is negative. The
 137 update rule for generalized inventory values is as follows:

$$G_{n,t} = G_{n,t-1} + \underbrace{\sum_{o \in O: d_{v_o} = n} q_{v_o} L_{v_o} \mathbf{1}_{t_o + \tau_o = t}}_{\text{goods received by orders finished at time } t} - \underbrace{\sum_{o \in O: n \in s_{v_o}} g_{v_o, n} L_{v_o, t_o} \mathbf{1}_{t_o = t}}_{\text{goods issued to satisfy orders}} - \underbrace{D_{n,t}}_{\text{demand}}, \quad (1)$$

138 where the random demand $d_{n,t}$ is sampled from an iid distribution $d_{n,t} \stackrel{\text{i.i.d.}}{\sim} \mathcal{D}_n$ at each time t . In the
 139 update formula 1 denotes the indicator function. $G_{n,t}$ is then clipped to the range $G_{n,t} \in [-\tilde{w}_n, w_n]$
 140 to enforce the storage capacity and backorder capacity conditions. While the maximum inventory
 141 constraint w_n is relevant to model maximum storage capacities, the maximal backorder constraint
 142 \tilde{w}_n is mostly technical and will be set to a large number. Changes in generalized inventory level
 143 result from three factors: goods received from finished orders, goods issued due to started orders,
 144 and goods issued due to sales towards customers. In this model, unfulfilled demand is not erased,
 145 but is instead stored as backorder, eventually converting into fulfilled sales once inventory becomes
 146 available to clear the backlog. The assumption is unrealistic in B2C (business-to-customer) models
 147 but common in B2B models where customers cannot change easily their suppliers for regulatory
 148 reasons. Inventory level and backorder are deduced from $G_{n,t}$:

$$I_{n,t} = \max(0, G_{n,t}) \quad \text{and} \quad B_{n,t} = \max(0, -G_{n,t}). \quad (2)$$

149 **Reward function:** For state-action pairs (s, a) described above, the reward is defined as

$$R(s, a) = - \sum_{n \in \mathcal{N}} \left(\underbrace{c_n^h I_n}_{\text{inventory holding cost}} + w \underbrace{c_n^b B_n}_{\text{backorder cost}} \right) - \sum_{v \in \mathcal{V}} \underbrace{c_v L_v}_{\text{value creation execution costs}}, \quad (3)$$

150 where inventory/backorder values are from the state and the lot sizes from the action. We introduce
 151 a *backorder weight* $w > 0$ that is not part of the model. This non-trivial weight allows the decision
 152 maker to adjust preferences between inventory/backorder (see Figure 5). We use $w = 30$ in our
 153 experiments (justified by the Pareto curve in Figure 5).

154 3 RL training

155 Training multi-echelon supply chain environments is challenging for a number of reasons. Other
 156 than the challenge of creating an efficient simulator, the action-spaces are too big for naive explora-
 157 tion to succeed. The action-space of a multi-echelon supply chain can grow exponentially with
 158 the size of the supply chain network, making it extremely hard for an RL agent to find a good policy
 159 through unguided exploration. To enable deep RL training with PPO, we built a gymnasium envi-
 160 ronment for the order-based supply chain management described above. To get hold of the large
 161 action-space, we introduce an action-shaping and -masking component. Finally, we propose a new
 162 pre-training concept for RL problems with known optimization heuristics.

163 3.1 Action-shaping and -masking

164 While our order-based modeling already gives
 165 a lot of structure and reduction in size to the
 166 action-space, two more ingredients are used.

167 **Multi-discrete action shaping:** In the order-
 168 based modeling above, we introduced the concept
 169 of value creations to structure dependent
 170 sub-actions in the supply chain. Value crea-
 171 tions allow us to simplify the action-space

172 enormously, as not all combinations of actions at all nodes must be taken into account when placing
 173 an order. Additionally, we use multi-discrete action shaping, see [Kanervisto et al. \(2020\)](#). The il-
 174 lustration shows that we do not need to consider all combinations of value creation orders possible,
 175 but only those with dependencies through their sources. Structuring the actions into independent
 176 components and utilizing multi-discrete action-spaces, we improve computational efficiency.

177 **Invalid action masking:** A substantial portion of actions is typically invalid due to constraints such
 178 as limited material availability. If no products exist in the supply chain, production must precede
 179 any other action, rendering nearly all actions initially invalid. We utilize invalid action masking, see
 180 e.g. [Huang & Ontañón \(2020\)](#), in the policy network that outputs probabilities for all actions. To
 181 implement masking, we adjust the logits before applying a final softmax function. Invalid actions
 182 are assigned a large negative number to ensure invalid actions have near zero probabilities in the
 183 softmax output.

184 3.2 Reward optimized MRP heuristic

185 Material Requirements Planning (MRP) is a dynamic programming inspired rule-based heuristic
 186 used to determine the quantity and timing of production and replenishment orders, more details in the
 187 Appendix 7 for a very simple variant of MRP. Based on a number of checks, MRP suggests times and
 188 quantities for value creations. It can be interpreted as a policy which is inefficient to evaluate and is
 189 *not* designed to maximize a particular reward function. What makes the MRP algorithm unpleasant
 190 is not only the effort in computing the action but also that the algorithm requires so-called safety
 191 stock levels S_n as a hyperparameter. For our experiments, we use a novel safety stock approach that
 192 is promising in its own rights, as it can be used for complex supply chains. This is in contrast to
 193 safety stock policies (such as based on Gaussian tail bounds for estimated one-step workloads) used

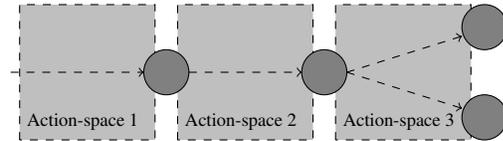


Figure 3: Independent action-spaces, dotted lines represent parts of the supply chain.

194 in real-life systems that are optimal in very small examples but can fail badly for complex systems.
 195 Interpreting safety stocks as a hyperparameter to the "MRP policy" we use Bayesian optimization
 196 (scikit-optimize library) to maximize the value function $V(S) = \mathbb{E}_{\text{MRP}(S)}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$. The
 197 approach can be seen as a combination of standard MRP and reward optimization. The maximal
 198 safety stock vector is denoted by S^* , the *reward optimized MRP policy* by $\text{MRP}(S^*)$.

199 3.3 Online imitation learning

200 We will show below that standard RL training using PPO is delicate for multi-echelon supply chain
 201 problems and propose two fixes. The most effective one is to imitate an MRP heuristic and use
 202 the imitated policy network as a pre-training for PPO. For the imitation we use DAgger (dataset
 203 aggregation), a popular imitation learning algorithm for online learning (Ross et al., 2011). DAgger
 204 is an improvement on Behavior Cloning (BC) (Pomerleau, 1991) that replaces the static dataset of
 205 BC with a dynamic dataset, containing trajectories generated by the training agent and the expert
 206 together and annotated by the expert. DAgger training consists of N iterations, generating new
 207 training data each iteration while increasing the influence of the learned agent on the generated
 208 trajectories. At iteration t , a mixed policy π_t is used to generate environment interactions:

$$\pi_t = \begin{cases} \pi_* & : \text{with probability } \beta_t \\ \pi_{\theta_t} & : \text{with probability } 1 - \beta_t \end{cases},$$

209 where $\beta_t \in [0, 1]$ is a mixing parameter that controls the probability of following the expert policy
 210 π_* versus the learned policy π_{θ_t} . Here π_{θ} is a policy parametrized by a neural network with fixed
 211 architecture and weight vector θ . The collected dataset $\mathcal{D}_i = \{(s, \pi_t(s))\}$ of visited states and
 212 corresponding actions is added to an aggregated dataset: $\mathcal{D} = \bigcup_{i=1}^N \mathcal{D}_i$ and a new policy $\pi_{\theta_{i+1}}$ is
 213 then trained on \mathcal{D} . We follow the typical method of linearly annealing β_i from 1 to 0.

214 3.4 Training

215 **Standard PPO:** Training is performed using PPO with random (Glorot) initialization. We use the
 216 maskable-PPO implementation of Stable Baselines3 available in SB3-Contrib (Raffin et al., 2021),
 217 see appendix for details. It turns out that the training is extremely slow, multi-echelon supply chain
 218 training is very much prone to metastability effects (see appendix). What happens is the following:
 219 the agent learns quickly to optimize all nodes except the last, where all material exceeding the
 220 storage capacity is discarded. Leaving this local maximum is delicate for gradient-based algorithms.
 221 To learn to optimize the terminal node, earlier sub-decisions become suboptimal so the algorithm
 222 must first worsen before improving. Orange learning curves in the experiments below show the
 223 metastability, the agent is stuck in suboptimal strategies for millions of interactions and suddenly
 224 improves (when inventory of the last node is decreased).

225 **Modified reward function:** A first mitigation of the metastability problem is to modify the reward
 226 function to penalize the waste of excess material. We change the reward function R into

$$\bar{R}(s, a) = R(s, a) - c \sum_{n \in N} c_n^h \cdot e_{n,t}, \quad (4)$$

227 where $e_{n,t}$ is the amount of excess material discarded from node n at time t when it exceeds the
 228 capacity w_n . c is a penalty weight, set to 100 in the experiments. We note that the change of R
 229 to \bar{R} does not change the optimal solution, since the capacities w_n are typically much higher than
 230 the optimal inventory levels. However, our experiments show that the use of \bar{R} speeds up training
 231 significantly and allows the agent to escape suboptimal local maxima.

232 **MRP pre-trained PPO:** We now propose a novel approach to drastically improve PPO training.
 233 The idea is to make the PPO agent start in a favorable parameter valley, a valley that represents
 234 policies that do not sacrifice optimality at single nodes (typically the last). For the pre-training
 235 we imitate the $\text{MRP}(S)$ policy. We learn a neural network that mimics the same action making as
 236 $\text{MRP}(S)$. Every evaluation of the $\text{MRP}(S)$ policy is costly, the entire rule-based algorithm needs

237 to be performed. We run DAgger for ten iterations only, generating a single trajectory of 3600 (10
 238 years) time-steps in each iteration. For the i th iteration of DAgger the expert policy $\pi_* := \text{MRP}(S)$
 239 is chosen in every step with decreasing probability β_i . This is enough to bring the policy above 80%
 240 accuracy predicting MRP’s actions. Even though the evaluation of the MRP policy is costly, a few
 241 thousand evaluations are negligible compared to millions of iterations needed for PPO. The MRP
 242 pre-trained policy is then used to initialize the PPO training. It is important to note that the choice
 of safety-stock vector S is not crucial, every reasonable choice improves the PPO training a lot.

Algorithm 1 MRP pre-trained PPO for multi-echelon inventory optimization

Step 1 (Safety stocks): Determine good safety stock vector S , e.g. using Bayesian opt.

Step 2 (Imitation learning): Perform DAgger using expert policy $\pi_* = \text{MRP}(S)$:

for $i = 1, \dots, N$ **do**

 Obtain rollout using in every step $\text{MRP}(S)$ (resp. π_{θ_i}) with probability β_i (resp. $(1 - \beta_i)$).

 Add rollout state-action pairs to buffer \mathcal{D} , use supervised imitation on \mathcal{D} to obtain θ_{i+1} .

end for

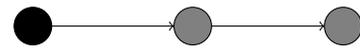
Step 3 (PPO): Run PPO with initial policy network obtained in Step 2. =0

243

244 **4 Experiments**

245 The algorithms are compared on a simple three-nodes example and the more complex ten-nodes
 246 example from Figure 1.

247 The simple problem consists of a factory (black) and two
 248 warehouses (gray) with random customer demands (negative
 249 binomial with means 367 (middle) and 172 (right)). Replen-
 250 ishment from factory to warehouse and from warehouse to
 251 warehouse have random lead times (Poisson auto-regressive
 252 with means 6 (resp. 15)). Variants of the simple toy example are often used in the supply chain
 253 literature. They also appear in industry, parameters for our experiment (see Appendix) stem from a
 254 real-world problem in the process manufacturing industry.



Example, gray nodes have demand

255 To show that the choice of S is not crucial to have a well-performing MRP pre-trained PPO agent
 256 we used a decent safety stock vector S instead of the optimal S^* .

257 **Learning curves:** MRP value functions are plotted as horizontal lines. RL agents are evaluated
 258 every $2.5 \cdot 10^4$ environment interactions, value functions are evaluated on 80 episodes of length
 259 3600 (representing ten years). The training curves are smoothed with a Gaussian filter with $\sigma = 2$.
 260 Standard PPO struggles to learn the optimal policy. Further investigation of single node inventories
 261 shows that the algorithm is stuck in local maxima of strategies that sacrifice the final warehouse
 262 which is kept at maximal capacity. The PPO agent takes millions of interactions to find a param-
 263 eter region that reduces the inventory at the final warehouse. PPO with modified rewards is much
 264 worse at the beginning (due to the penalization) and then improves much quicker, as expected the
 265 metastable behavior disappeared. MRP pre-trained PPO clearly beats PPO without pre-training.

266 **Performance - cost vs. backorder:** In Figure 5 we provide experiments to see the effects on cost vs.
 267 backorder when varying inventory/backorder preferences. The graphs show that MRP pre-trained
 268 PPO agents slightly outperform the reward-optimized MRP policy. It is natural to ask why it might
 269 be interesting to prefer RL policies to the simple (reward optimized) MRP policy. First, it is timely
 270 to evaluate for every decision the rule-based MRP policy, see Appendix 7. In contrast to evaluating
 271 a neural network for RL policies, it requires a number of algorithmic steps that in real-life large
 272 supply chains is an important problem. Secondly, as our next result shows, the RL agent generalizes
 273 better to changes in the environment, one of the key topics in current supply chain organizations.

274 **Robustness:** In Figure 6 we provide zero-shot experiments, where the trained policies are compared
 275 on unseen demand distributions with changed customer demand. For decreased demand the RL
 276 agent generalizes better than reward-optimized MRP, the inventory/backorder situation improves.

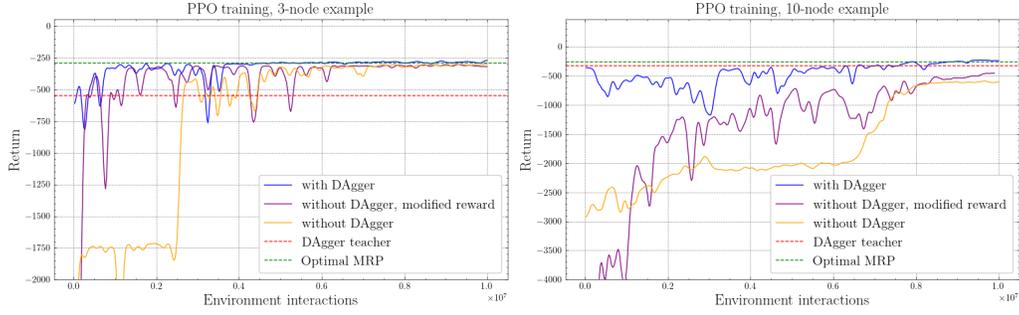


Figure 4: Learning curves for three PPO trained agents compared to the (constant) value function of MRP agents (green with reward optimized S^* , red with suboptimal S used for DAgger training). DAgger pre-training (blue) enables better performance than optimized MRP. Without pretraining, adding the capacity penalty of Eq. (4) helps to avoid overstocking material (purple vs. orange).

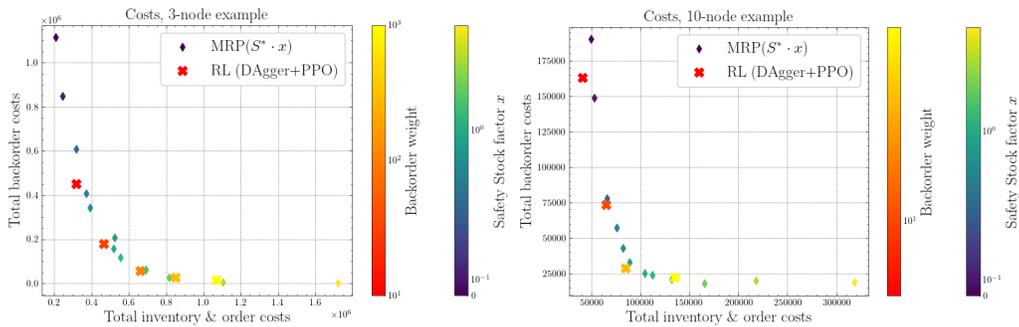


Figure 5: Pareto curves (backorder vs. inventory costs) of RL agents and MRP. We fix all RL training parameters but scale the backorder weight w of Eq. (3). Similarly, for MRP we change inventory/backorder preferences by scaling all safety stock values a factor x .

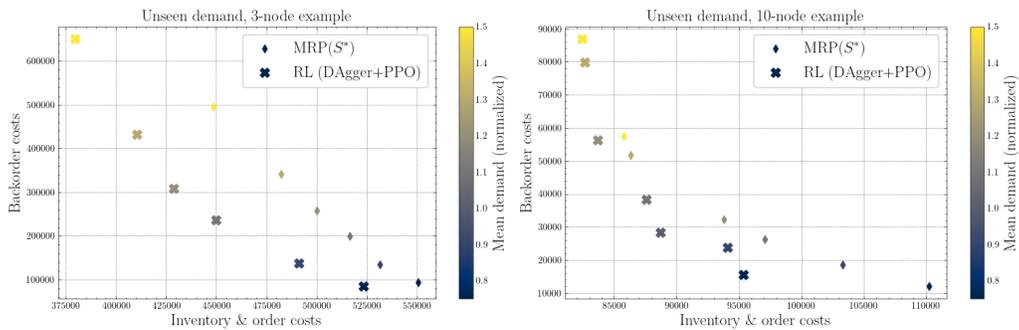


Figure 6: Zero-shot performance for new demand. Changing magnitude of demand, for lower demand RL agent generalizes clearly better.

277 5 Conclusion

278 In the present article we introduced an MDP modeling for order-based supply chain management.
 279 Combining action-shaping and -masking to reduce the action-space we created a gymnasium envi-
 280 ronment to run RL algorithms. Since plain vanilla PPO is slowed down by metastability effects we
 281 introduced supply chain specific patches, in particular using imitation learning on a rule-based MRP
 282 heuristic. The approach has potential to be used in other RL approaches to classical OR problems.

283 Due to its enormous importance in supply chain management it would be very beneficial in future
 284 work to understand how policy networks can be improved on zero-shot learning by training the
 285 agent on more and different extreme scenarios. The generalization ability of neural networks has a
 286 potential huge benefit to supply chain robustness for sudden changes in real-life.

287 **References**

- 288 Ana Barbosa-Povoa, Albert Corominas, and João Miranda. *Optimization and Decision Support Sys-*
289 *tems for Supply Chains*. 01 2017. ISBN 978-3-319-42419-4. DOI: 10.1007/978-3-319-42421-7.
- 290 Jonas Degraeve, Federico Felici, Jonas Buchli, Michael Neunert, Brendan Tracey, Francesco
291 Carpanese, Timo Ewalds, Roland Hafner, Abbas Abdolmaleki, Diego Casas, Craig Donner,
292 Leslie Fritz, Cristian Galperti, Andrea Huber, James Keeling, Maria Tsimpoukelli, Jackie
293 Kay, Antoine Merle, Jean-Marc Moret, and Martin Riedmiller. Magnetic control of toka-
294 mak plasmas through deep reinforcement learning. *Nature*, 602:414–419, 02 2022. DOI:
295 10.1038/s41586-021-04301-9.
- 296 Jean-Philippe Gayon, Guillaume Massonnet, Christophe Rapine, and Guy Stauffer. Fast approxi-
297 mation algorithms for the one-warehouse multi-retailer problem under general cost structures and
298 capacity constraints. *Mathematics of Operations Research*, 42(3):854–875, 2017.
- 299 Joren Gijsbrechts, Robert N Boute, Jan A Van Mieghem, and Dennis J Zhang. Can deep reinforce-
300 ment learning improve inventory management? performance on lost sales, dual-sourcing, and
301 multi-echelon problems. *Manufacturing & Service Operations Management*, 24(3):1349–1368,
302 2022.
- 303 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recog-
304 nition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp.
305 770–778, 2016.
- 306 Shengyi Huang and Santiago Ontañón. A closer look at invalid action masking in policy gradient
307 algorithms. *arXiv preprint arXiv:2006.14171*, 2020.
- 308 Christian D Hubbs, Hector D Perez, Owais Sarwar, Nikolaos V Sahinidis, Ignacio E Grossmann,
309 and John M Wassick. Or-gym: A reinforcement learning library for operations research problems.
310 *arXiv preprint arXiv:2008.06319*, 2020.
- 311 Anssi Kanervisto, Christian Scheller, and Ville Hautamäki. Action space shaping in deep reinforce-
312 ment learning. In *2020 IEEE conference on games (CoG)*, pp. 479–486. IEEE, 2020.
- 313 Dhruv Madeka, Kari Torkkola, Carson Eisenach, Anna Luo, Dean P Foster, and Sham M Kakade.
314 Deep inventory management. *arXiv preprint arXiv:2210.03137*, 2022.
- 315 Johan Marklund and Kristoffer Rosling. Lower bounds and heuristics for supply chain stock alloca-
316 tion. *Operations Research*, 60(1):92–105, 2012.
- 317 J. Orlicky. *Material requirements planning: the new way of life in production and inventory man-*
318 *agement*. New York: McGraw-Hill, 1975.
- 319 Afshin Oroojlooyjadid, Mohammad Nazari, Lawrence V. Snyder, and Martin Takáč. A deep q-
320 network for the beer game: Deep reinforcement learning for inventory optimization. *Manufactur-*
321 *ing & Service Operations Management*, 24(1):285–304, 2022.
- 322 Hector D Perez, Christian D Hubbs, Can Li, and Ignacio E Grossmann. Algorithmic approaches to
323 inventory management optimization. *Processes*, 9(1):102, 2021.
- 324 Dean A Pomerleau. Efficient training of artificial neural networks for autonomous navigation. *Neu-*
325 *ral computation*, 3(1):88–97, 1991.
- 326 Meng Qi, Yuanyuan Shi, Yongzhi Qi, Chenxin Ma, Rong Yuan, Di Wu, and Zuo-Jun Shen. A
327 practical end-to-end inventory management model with deep learning. *Management Science*, 69
328 (2):759–773, 2023.

- 329 Antonin Raffin, Ashley Hill, Adam Gleave, Anssi Kanervisto, Maximilian Ernestus, and Noah
330 Dormann. Stable-baselines3: Reliable reinforcement learning implementations. *Journal of*
331 *Machine Learning Research*, 22(268):1–8, 2021. URL [http://jmlr.org/papers/v22/](http://jmlr.org/papers/v22/20-1364.html)
332 [20-1364.html](http://jmlr.org/papers/v22/20-1364.html).
- 333 Stéphane Ross, Geoffrey Gordon, and Drew Bagnell. A reduction of imitation learning and struc-
334 tured prediction to no-regret online learning. In *Proceedings of the fourteenth international con-*
335 *ference on artificial intelligence and statistics*, pp. 627–635. JMLR Workshop and Conference
336 Proceedings, 2011.
- 337 David Silver, Aja Huang, Chris J Maddison, Arthur Guez, Laurent Sifre, George Van Den Driessche,
338 Julian Schrittwieser, Ioannis Antonoglou, Veda Panneershelvam, Marc Lanctot, et al. Mastering
339 the game of go with deep neural networks and tree search. *Nature*, 529(7587):484–489, 2016.
- 340 Nathalie Vanvuchelen, Joren Gijsbrechts, and Robert Boute. Use of proximal policy optimization
341 for the joint replenishment problem. *Computers in Industry*, 119:103239, 2020.
- 342 Min Zhao and Mingyu Zhang. Multiechelon lot sizing: New complexities and inequalities. *Opera-*
343 *tions Research*, 68(2):534–551, 2020.

344
345
346

Supplementary Materials

The following content was not necessarily subject to peer review.

347 6 Experimental details

348 6.1 Environment

349 Here are the paramters for the 3-node example, the exact parameters for the 10-node example must
350 be kept confidential.

node	c_n^h	c_n^b	w_n	\tilde{w}_n	demand
middle	26.7	34	170093	12465	NB(0.11, 0.0003)
right	29.5	42	127646	8025	NB(0.081, 0.00047)
replenishment	q_v	k_v^{\max}	k_v^{\max}		lead time
left	4180	2	0		ARPois(6, 20, 0.98)
right	4100	2	1500		ARPois(15, 20, 0.98)

351

352 Both demand distributions are Negative Binomial, NB(r, p). To reflect realistic fluctuating lead
353 times, we use an autoregressive variant of the Poisson distribution, ARPois(λ_0, h, ϕ). This distribu-
354 tion produces each time-step t a lead time τ_t which is drawn from the distribution Pois(λ_t), where
355 λ_t depends on the previously h drawn lead time values τ according to the formula:

$$\lambda_t = \max\left(0, \phi \cdot \frac{\sum_{t' \in [t-h, t]} \tau_{t'}}{h} + (1 - \phi)\lambda_0\right) \quad (5)$$

356 This produces a time-correlated Poisson distribution which retains an expected value of λ_0 .

357 6.2 Details on PPO training

358 We train feedforward neural nets with PPO. We add skip connections (He et al., 2016) every two
359 layers to enable training deep networks, effectively using 2 residual blocks in both the value and
360 policy networks.

Hyper-parameter	parameter value
no. of environment interactions	10^7
policy network	width 256, depth 4
actor network	width 256, depth 4
activation function	ReLU
discount factor γ	0.99
GAE paramter	0.95
Adam learning rate	$2.5 \cdot 10^{-4}$
batch size	64

361 6.3 Metastability

362 The learning curves in Figure 4 highlight the difficulty of deep RL training for multi-echelon supply
363 chain optimization if the reward function is not chosen carefully. Metastability effects occur as there

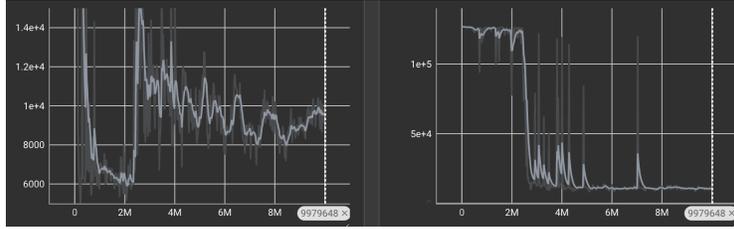


Figure 7: 3-node example, inventory during training at middle node (left) and right node (right)

364 are suboptimal strategies that manifest local maxima for the parameter vectors of the policy network.
 365 As a consequence, gradient based algorithms struggle to improve the suboptimal strategy. To explain
 366 what happens we plotted the inventory levels of the 3-node example during training (inventory plots
 367 look similar for the 10-node example). The learning curve of Figure 4 shows a sudden improvement
 368 after 2.5M environment interactions, preceded by a return drop. This is reflected in Figure 7. The
 369 agent finds quickly the suboptimal strategy to optimize the inventory at the middle node while com-
 370 pletely sacrificing the terminal right node, running at maximal inventory 127646. At 2.5M iterations
 371 the RL agent deviates from the local maximum and explores a new strategy, reducing inventory at
 372 the terminal node and increasing inventory at the middle node.

373 7 Material requirements planning (MRP)

374 The main algorithmic invention of this article is to combine off-the-shelf RL training (PPO) through
 375 imitation learning with rule-based heuristics from operations research (OR). There are several OR
 376 algorithms to solve approximately different supply chain problems. For the multi-echelon inventory
 377 optimization (MEIO) problem studied in the present article we use a dynamic programming inspired
 378 rule-based algorithm that is (with various modifications) implemented in many industry supply chain
 379 solutions. We now give a quick overview for the interested RL researcher.

380 The rule-based algorithm implements a time-phased Material Requirements Planning (MRP I) sys-
 381 tem to maintain inventory levels above safety stock thresholds across all nodes in the supply chain.
 382 Rooted in the foundational work [Orlicky \(1975\)](#), the process begins by exploding dependencies from
 383 downstream nodes (e.g., retailers or finished goods) to upstream suppliers, following the hierarchical
 384 structure of the multi-echelon supply chain. Inventory projections are calculated in daily time buck-
 385 ets over a fixed $H = 150$ planning horizon. Starting from the current day t , the system computes the
 386 projected available balance (PAB) for each subsequent day $s \in [t, t + H]$, accounting for scheduled
 387 receipts, planned orders, and demand forecasts. If the PAB is projected to fall below the safety stock
 388 level at time T , a planned order is generated to replenish the deficit. Orders are offset by lead times
 389 using backward scheduling: for an order requiring τ days of lead time, the release date is set to
 390 $T - \tau$. If this calculated release date precedes the current day t , the order is flagged as overdue and
 391 scheduled for immediate release. This daily recalibration ensures alignment with the core principle
 392 of time-phased net requirement calculation, where material plans are dynamically adjusted to reflect
 393 real-time demand and supply conditions. Rule-based MRP algorithms are dynamic-programming,
 394 heuristic-based algorithms. It implements a safety-stock approach to managing the supply chain,
 395 meaning it predicts the future inventory levels of all nodes in the chain and tries to ensure inventory
 396 never falls below the "safety stock" that must be given to the algorithm.

397 Since we use the MRP algorithm in our examples without multi-material manufacturing steps, we
 398 give pseudo-code for a simplified version of MRP. It should be noted that the algorithm is a very
 399 simple MRP variant that does not estimate demand expectations and lead times on the run. We do
 400 this for a fair comparison to the RL agents, otherwise demand distributions should also be included
 401 in the MDP state-space and not be given as part of the model.

402 There are two novel ideas we add on the standard OR literature.

Algorithm 2 MRP Algorithm (without multi-material manufacturing steps)

Input: expected demands $\mathbb{E}(d)$ and lead times $\mathbb{E}(l)$, safety stocks S_n for all nodes, current generalized inventories $G_{n,t}$ and running orders set O .

for each node n in topological order **do**

for each time-step $s \in [t, t + H]$ **do**

$gen(n, \tau) \leftarrow$ amount of additional material by finished orders

$G_{n,s} \leftarrow G_{n,s-1} + gen(n, s) - \mathbb{E}(d(n))$

if $G_{n,s+1} < S_n$ **then**

 Set number of lots L to minimal number containing at least amount $S_n - G_{n,s+1}$.

if procurement is possible **then**

 Add to O an order from a source node. L lots, start time: $\max(t, s - \mathbb{E}(l))$

else

 Choose source node n' that maximizes $\frac{G_{n',s+1}}{\text{num_outgoing}(n')}$, where $\text{num_outgoing}(n')$ denotes the number of nodes supplied by node n' .

 Add to O an order from node n' to node n . L lots and start time $\max(t, s - \mathbb{E}(l))$

end if

end if

end for

end for

Output: all orders in O that start at time t . =0

- 403 • We interpret $\text{MRP}(S) =: \pi$ as a policy. The action (orders) in the state S (inventory level and
404 current order book) are the output orders of the algorithm given above (the orders suggested by
405 the algorithm to be placed at initial time t).
- 406 • The safety stock vector S is a required input to the algorithm. We define the reward-based optimal
407 safety stock vector S^* by maximizing the expected reward R under the MRP run defined by
408 the MRP algorithm: $V(S) = \mathbb{E}_{\text{MRP}(S)}[\sum_{t=0}^{\infty} \gamma^t R(s_t, a_t)]$. Since S is a hyperparameter to the
409 algorithm, it is natural to use a Bayesian optimization algorithm to do so.