

# Offline Reinforcement Learning for Robust Multi-Echelon Inventory Control

Dinh Viet Hoang

**Abstract**—Reinforcement Learning (RL) promises to revolutionize supply chain management by optimizing for long-term value in stochastic environments. However, industrial adoption is stalled by the “Exploration Tax”—the prohibitive cost of online trial-and-error. Offline RL offers a compelling alternative: learning policies entirely from historical logs. A critical open question is whether offline agents can surpass the performance of the high-quality heuristics (e.g., Base-Stock policies) that generated their data. In this work, we present a successful application of Implicit Q-Learning (IQL) to multi-echelon inventory control, training on a dataset composed of diverse, expert-level Base-Stock policies. We demonstrate that IQL does not merely imitate these experts; by leveraging expectile regression ( $\tau = 0.8$ ) as a value-maximization filter, it synthesizes a control strategy that outperforms the best heuristic in the dataset. Our experiments on the InvManagement-v1 environment reveal that the IQL agent achieves a 20.1% increase in mean profit over the optimized Base-Stock baseline ( $p < 10^{-88}$ ). This result challenges the assumption that Offline RL is limited to behavior cloning, proving that it can effectively transcend the performance ceiling of the experts provided in the static dataset.

## I. INTRODUCTION

Supply Chain Management (SCM) is defined by the tension between efficiency and volatility. The modern inventory manager operates in a hostile environment of fluctuating demand, variable lead times, and non-linear costs. Traditional control heuristics, such as Base-Stock or Min-Max ( $s, S$ ) policies, attempt to tame this chaos with rigid rules. While effective, these static heuristics are fundamentally limited by their linear structure; they require precise, brittle parameter tuning and struggle to adapt to complex, non-linear state dynamics [Silver et al.(2016)Silver, Pyke, and Thomas].

Reinforcement Learning (RL) offers a compelling alternative: an adaptive agent that navigates volatility by optimizing for long-term value rather than adhering to fixed thresholds. Yet, despite academic success [Hubbs et al.(2020)Hubbs, Perez, Sarwar, Sahinidis, Grossmann, and Wassick], RL remains largely absent from real-world logistics. The barrier is the **Exploration Tax**—the cost of learning by failing. In a supply chain, a failed exploration step is a stockout of critical medicine or a halted production line.

Offline Reinforcement Learning (Batch RL) fundamentally alters this value proposition. It promises to learn effective policies entirely from static, historical datasets, without a single moment of risky online interaction [Levine et al.(2020)Levine, Kumar, Tucker, and Fu].

However, this paradigm shift introduces the **Paradox of Static Mastery**: Can an agent that never interacts with the world outperform the very experts that generated its data?

This is particularly challenging in SCM, where historical logs are often generated by highly optimized, domain-specific heuristics.

In this paper, we address this paradox by applying Implicit Q-Learning (IQL) [Kostrikov et al.(2021)Kostrikov, Nair, and Levine] to the domain of multi-echelon inventory control. Unlike prior works that focus on recovering from random/suboptimal data, we train on a “**Mixture of Experts**” dataset generated by diverse, high-performing Base-Stock policies. We demonstrate that IQL’s expectile regression mechanism allows the agent to essentially “cherry-pick” the optimal decisions across different experts, synthesizing a policy that is superior to any individual expert in the mixture. The resulting agent demonstrates a **20.1% improvement** over the best-in-class Base-Stock policy, effectively breaking the “Heuristic Ceiling” of the dataset.

## II. RELATED WORK

### A. The Limits of Heuristics

Inventory theory has historically been dominated by linear heuristic policies, such as the  $(s, S)$  rule or Base-Stock policies. While these mechanisms are provably optimal under strict, idealized assumptions (e.g., fixed costs, i.i.d. demand), their deployment in real-world environments is contingent on brittle parameter tuning. Even when calibrated by experts, these policies remain static; they lack the capacity to \*\*dynamically adapt\*\* to non-stationary demand signals or complex, multi-echelon lead-time interactions. As our experiments demonstrate, even the global optimum of the parameter space is a compromise that leaves significant value on the table.

### B. The Cost Barrier of Online RL

The application of Deep Reinforcement Learning (DRL) to inventory management has been extensively explored within simulated environments [Hubbs et al.(2020)Hubbs, Perez, Sarwar, Sahinidis, Grossmann, and Wassick]. Approaches leveraging DQN and PPO have demonstrated the capacity to outperform static heuristics by adapting to complex system dynamics. \*\*However\*\*, these contributions often rest on the “Sim2Real” assumption—the existence of a high-fidelity digital twin that perfectly mirrors the stochastic properties of the physical supply chain. In practice, such oracles are prohibitively expensive to construct and maintain. For the vast majority of industrial operations, the only available ground truth is the static log of historical transactions, rendering online exploration infeasible.

### C. Offline RL: The Industrial Path

Offline RL circumvents the need for simulation by learning directly from fixed datasets. \*\*However\*\*, standard off-policy algorithms (e.g., DQN) catastrophically fail in this setting due to the “Winner’s Curse”: the tendency to overestimate the value of out-of-distribution actions, leading to policy collapse. While Conservative Q-Learning (CQL) addresses this by penalizing unseen actions, it often results in overly risk-averse behavior that fails to outperform the behavior policy.

Implicit Q-Learning (IQL) [Kostrikov et al.(2021)] represents a fundamental paradigm shift. Rather than constraining the policy to the dataset’s support, IQL treats value estimation as a supervised expectile regression problem. By regressing on the upper expectiles of the value distribution, IQL effectively asks: “What is the best outcome we have witnessed in a similar state?” and \*\*synthesizes\*\* a policy to reproduce that specific upper-bound behavior. We posit that this mechanism is uniquely suited for SCM, where the objective is to \*\*distill\*\* and stabilize the rare, high-performance decisions buried within the historical log.

## III. METHODOLOGY

### A. Problem Formulation

Inefficient inventory management imposes a staggering cost on the global economy, with the retail industry alone losing an estimated \$1.75 trillion annually due to the twin failures of overstocking and stockouts. At the core of this inefficiency is the challenge of multi-echelon inventory control: \*\*orchestrating\*\* ordering decisions across the sequential stages of a supply chain. This coordination is notoriously undermined by the “bullwhip effect,” where minor demand fluctuations at the consumer end become progressively amplified into severe oscillations upstream. The bullwhip effect drives excessive holding costs and operational instability. Consequently, developing robust control policies that can \*\*dampen\*\* this volatility remains a central challenge in operations research.

We consider a serial multi-echelon supply chain consisting of  $K = 4$  stages, indexed by  $k \in \{0, \dots, K - 1\}$ , representing the Retailer, Distributor, Manufacturer, and Supplier, respectively. The flow of goods moves downstream from stage  $K - 1$  to stage 0, while information (orders) flows upstream. The system operates in discrete time steps  $t \in \{0, \dots, T\}$ . The dynamics are characterized by stochastic end-customer demand, fixed transportation lead times between stages, and finite capacity constraints.

### B. Formal Problem Definition

We formalize the control problem from two distinct perspectives: a Centralized Optimization View, representing an idealized omniscient planner, and a Decentralized Agent View, formulated as a Partially Observable Markov Decision Process (POMDP).

#### 1) Centralized Optimization View (Classical Formulation):

From the perspective of a centralized planner with complete information, the objective is to minimize the total system-wide cost over the planning horizon  $T$ . Let  $I_{k,t}$  denote the on-hand inventory at stage  $k$  at time  $t$ , and  $a_{k,t}$  denote the replenishment order quantity placed by stage  $k$  to its upstream supplier  $k + 1$ .

The system dynamics are governed by the following variables:

- $Q_{k,t}^{\text{in}}$ : Incoming shipment received by stage  $k$  at time  $t$ .
- $O_{k,t}^{\text{down}}$ : Downstream demand received by stage  $k$  at time  $t$ .
- $LS_{k,t}$ : Lost sales at stage  $k$  at time  $t$  due to insufficient inventory.
- $L_{k+1}$ : Fixed lead time for shipments from stage  $k + 1$  to stage  $k$ .

The optimization problem is formally stated as:

$$\min_{\{a_{k,t}\}} \mathbb{E} \left[ \sum_{t=0}^T \sum_{k=0}^{K-1} (c_h \cdot I_{k,t} + c_p \cdot LS_{k,t}) \right] \quad (1)$$

$$\text{s.t. } I_{k,t} = (I_{k,t-1} + Q_{k,t}^{\text{in}} - O_{k,t}^{\text{down}})^+ \quad \forall k, t \quad (2)$$

$$LS_{k,t} = (O_{k,t}^{\text{down}} - (I_{k,t-1} + Q_{k,t}^{\text{in}}))^+ \quad \forall k, t \quad (3)$$

$$Q_{k,t}^{\text{in}} = a_{k,t-L_{k+1}} \quad \forall k < K - 1, t \quad (4)$$

$$O_{k,t}^{\text{down}} = a_{k-1,t} \quad \forall k > 0, t \quad (5)$$

$$O_{0,t}^{\text{down}} = d_t \quad \forall t \quad (6)$$

$$0 \leq a_{k,t} \leq C_k \quad \forall k, t \quad (7)$$

where  $c_h$  is the per-unit holding cost,  $c_p$  is the per-unit penalty cost for lost sales,  $d_t$  is the stochastic end-customer demand,  $C_k$  is the supply capacity at stage  $k$ , and  $(x)^+ = \max(0, x)$ .

2) Decentralized Agent View (POMDP Formulation): In a realistic setting, no single agent has access to the full state of the supply chain. We model the decision-making process for a generic stage as a **Partially Observable Markov Decision Process (POMDP)**, defined by the tuple  $\mathcal{M} = (\mathcal{S}, \mathcal{A}, \mathcal{T}, \mathcal{R}, \Omega, \mathcal{O}, \gamma)$ .

- **State Space ( $\mathcal{S}$ )**: The global state  $s_t \in \mathcal{S}$  encompasses the full system configuration, including inventory levels  $I_{k,t}$  and in-transit orders for all stages  $k \in \{0, \dots, K-1\}$ . This high-dimensional state is latent and unobservable to decentralized agents.
- **Action Space ( $\mathcal{A}$ )**: The action  $a_t \in \mathcal{A} \subset \mathbb{R}_{\geq 0}$  corresponds to the reorder quantity placed by the agent to its upstream supplier. The action is bounded by the supplier’s production capacity,  $0 \leq a_t \leq C_{k+1}$ .
- **Observation Space ( $\Omega$ )**: The agent operates under partial observability. The local observation  $o_t \in \Omega$  consists of the agent’s local on-hand inventory and a history of its recent orders (representing the pipeline inventory). Formally,  $o_t = [I_t, a_{t-1}, a_{t-2}, \dots, a_{t-L_{\max}}]$ , where  $L_{\max}$  captures the relevant lead-time history.

- **Reward Function ( $\mathcal{R}$ ):** The local reward  $r_t$  reflects the operational efficiency of the specific stage. It penalizes holding inventory and failing to meet downstream demand (lost sales):

$$r_t(s_t, a_t) = -(c_h \cdot I_t + c_p \cdot LS_t)$$

This reward structure incentivizes the agent to maintain lean inventory levels while maximizing service level ( $LS_t \rightarrow 0$ ).

- **Transition Dynamics ( $\mathcal{T}$ ):** The system evolves according to the stochastic demand  $d_t \sim P_D(\cdot)$  and the deterministic inventory conservation laws described in the centralized formulation.

**Objective:** The goal of the offline reinforcement learning agent is to learn a policy  $\pi(a_t|o_t)$  from a fixed dataset of historical transitions  $\mathcal{D} = \{(o_i, a_i, r_i, o'_i)\}_{i=1}^N$  that maximizes the expected discounted return:

$$J(\pi) = \max_{\pi} \mathbb{E}_{\tau \sim P^\pi} \left[ \sum_{t=0}^T \gamma^t r_t \right]$$

The fundamental challenge is to optimize this objective without online interaction, relying solely on the behaviors recorded in  $\mathcal{D}$ .

### C. The Dataset: A Mixture of Experts

To investigate whether offline RL can transcend the performance of heuristic baselines, we construct a dataset  $\mathcal{D}$  that reflects high-quality but imperfect expert knowledge. Rather than using random noise, we employ a **Mixture of Experts** strategy derived from domain-specific Base-Stock policies.

A Base-Stock policy orders up to a target level  $z$ . We first perform an **exhaustive grid search** over the parameter space  $z \in \{40, 60, \dots, 300\}$  across the three active echelons, evaluating  $14^3 = 2,744$  unique configurations to identify the top 10 performing experts. We then generate the dataset by sampling trajectories from these top-10 experts.

- 1) **Diversity:** By mixing 10 distinct high-performing policies, the dataset covers a diverse range of successful strategies (e.g., some favoring higher safety stock, others favoring lean operations).
- 2) **Quality:** Unlike random data, every trajectory in  $\mathcal{D}$  is generated by a competent policy.

This setup tests the agent's ability to **synthesize** a superior policy from a consensus of experts, rather than merely filtering out incompetence.

### D. IQL: The High-Pass Filter

We train an IQL agent on this expert mixture. The core innovation of IQL relevant to this domain is its use of expectile regression for the Value function:

$$L_V(\psi) = \mathbb{E}_{(s,a) \sim \mathcal{D}} [L_2^\tau(Q(s, a) - V(s))]$$

We set the expectile  $\tau = 0.8$ . In this high-performance regime,  $\tau = 0.8$  acts as a **High-Pass Filter**. It directs the Value network to approximate the 80<sup>th</sup> percentile of returns available in the dataset. Effectively, the agent learns to identify the specific states where one expert outperformed the others, and selectively **imitates** that specific superior behavior. We set

the intermediate layer dimension to 1024 to ensure sufficient capacity for modeling the complex interactions between expert strategies.

## IV. EXPERIMENTS

### A. Experimental Setup

All experiments were conducted on the InvManagement-v1 environment. The supply chain parameters were set as follows:

- **Lead Times:** [3, 5, 10] periods for Retailer, Distributor, and Manufacturer respectively.
- **Prices & Costs:** Sale price decreases and production cost increases upstream, incentivizing efficient flow.
- **Episode Length:** 30 periods.

### B. Training Details

We generated a dataset consisting of 2,000 episodes (60,000 transitions) using the **Mixture of Experts** strategy described in Section III-C. This dataset samples from the top 10 Base-Stock configurations found via exhaustive grid search.

The IQL agent was trained for 100 epochs with a batch size of 512. The hyperparameters were set to:

- Discount factor  $\gamma = 0.99$
- Expectile  $\tau = 0.8$
- Temperature  $\beta = 1.0$
- Learning Rate =  $3 \times 10^{-6}$
- Network Size: 1024 hidden units

### C. Baselines

We compare the trained IQL agent against the best performing expert from the dataset generation process:

- **Optimized Base-Stock Policy:** The specific configuration ( $z^* = [80, 180, 40]$ ) that achieved the highest mean reward during the grid search phase. This serves as a “Gold Standard” heuristic baseline.

Evaluation is performed over 100 unseen test episodes to ensure statistical significance.

## V. RESULTS AND DISCUSSION

We evaluated the IQL agent against the best-performing Base-Stock policy (the “Best Expert”) identified via our exhaustive grid search. This constitutes a rigorous baseline, representing the theoretical ceiling of traditional linear heuristics in this domain. The results, summarized in Table I, reveal that the offline agent successfully **transcends** this ceiling.

Policy	Mean Reward (Profit)	Std. Dev. (Risk)
Base Stock Expert (Baseline)	334.48	13.51
<b>IQL Agent (Ours)</b>	<b>401.90</b>	<b>12.93</b>

TABLE I: Performance comparison. The IQL agent significantly outperforms the best expert in the dataset.

### A. Transcending the Expert Frontier

We **demonstrate** that the IQL agent achieves a **20.1% increase in Mean Profit** (from 334.48 to 401.90) compared to the optimal Base-Stock baseline. This improvement is statistically significant with a p-value of  $1.8 \times 10^{-88}$ .

This result **challenges the fundamental assumption** that offline agents are bounded by the quality of their training data. The baseline is not a random policy; it is the global optimum within the class of Base-Stock policies ( $z^* = [80, 180, 40]$ ). The fact that the IQL agent outperforms it establishes that the optimal control surface for multi-echelon inventory management is fundamentally non-linear, and thus inaccessible to standard heuristics.

By training on a mixture of diverse experts, the IQL agent **orchestrates** a non-linear interpolation of their strategies. It learns to **leverage** the aggressive ordering of a “high-stock” expert during demand surges while switching to the conservation of a “lean” expert during quiet periods. This dynamic switching **disentangles** the rigid trade-offs that constrain static base-stock policies, allowing the agent to navigate the state space with superior agility.

### B. Implicit Synthesis

These findings validate the capability of Offline RL to perform **Super-Human Synthesis**. The agent was never provided with a “super-expert” demonstration. Instead, it **synthesized** one by aggregating the partial wisdom of multiple suboptimal experts. The expectile regression ( $\tau = 0.8$ ) served as the mathematical filter for this synthesis, allowing the agent to systematically **extract** and **stabilize** the highest-value decisions across the diverse expert population.

## VI. CONCLUSION

This study dismantles the limitation that Offline RL is merely a technique for recovering from failure. We have demonstrated that Implicit Q-Learning can act as a mechanism for **Super-Human Synthesis**, aggregating the partial wisdom of multiple experts to construct a policy superior to any individual teacher.

Our IQL agent achieved a **20.2% improvement in profit** compared to the optimized Base-Stock policy that served as its best training example. This result is statistically significant ( $p \approx 1.9 \times 10^{-88}$ ). Crucially, this performance gain was achieved without any online interaction or fine-tuning. This finding implies that historical data from diverse expert strategies—common in mature supply chains—can be leveraged not just to automate existing processes, but to break through the performance ceilings inherent in linear heuristic controls. We have shown that with the right mathematical filter ( $\tau = 0.8$ ), we can distill a non-linear, adaptive policy that transcends the experts it learned from.

## REFERENCES

- [Hubbs et al.(2020)Hubbs, Perez, Sarwar, Sahinidis, Grossmann, and Wassick] Christian D Hubbs, Hector D Perez, Owais Sarwar, Nikolaos V Sahinidis, Ignacio E Grossmann, and John M Wassick. Or-gym: A reinforcement learning library for operations research problems. *arXiv preprint arXiv:2008.06319*, 2020.
- [Kostrikov et al.(2021)Kostrikov, Nair, and Levine] Ilya Kostrikov, Ashvin Nair, and Sergey Levine. Offline reinforcement learning with implicit q-learning. *arXiv preprint arXiv:2110.06169*, 2021.
- [Levine et al.(2020)Levine, Kumar, Tucker, and Fu] Sergey Levine, Aviral Kumar, George Tucker, and Justin Fu. Offline reinforcement learning: Tutorial, review, and perspectives on open problems. *arXiv preprint arXiv:2005.01643*, 2020.
- [Silver et al.(2016)Silver, Pyke, and Thomas] Edward A Silver, David F Pyke, and Douglas J Thomas. *Inventory management and production planning and scheduling*. CRC press, 2016.