

Forecasts as Priors over Initial States: Sample-Efficient Reinforcement Learning for Predictable Non-Stationarity

Anonymous Author(s)
Affiliation withheld for double-blind review

Abstract

We study reinforcement learning (RL) in non-stationary environments with *structured, forecastable drift*—a setting common in operations research, econometrics, and real-time decision systems. Classical RL methods initialize episodes uniformly or from replay, treating non-stationarity as noise. We formalise *Forecast-Seeded Non-Stationary RL (FS-NRLF)*, in which the forecast posterior parameterises the initial state distribution of each training episode. Our central technical contribution is Lemma 3.5: under Lipschitz transition dynamics with $L_P \leq 1$, a narrower initial-state distribution contracts the reachable-state variance *uniformly* across the episode horizon, not only at $t=0$. This yields two guarantees. First, on a restricted linear-Gaussian MDP class, per-episode sample complexity scales with the forecast signal-to-noise ratio $\text{SNR}_f \triangleq \sigma_f^2/\sigma_0^2$ (Theorem 3.7), matched by a lower bound $\Omega(\sigma_f^2 m/\epsilon^2)$ (Proposition 3.11). Second, regret $O(\sqrt{TH\sigma_f^2 m \log|\mathcal{A}|} + T\beta H)$ is tighter than the $O(\Delta^{1/3}T^{2/3})$ bound of Cheung et al. [2020] whenever $T > \sigma_f^6 m^3 H^3/\Delta^2$ (Theorem 3.10). On the M5 forecasting benchmark [Makridakis et al., 2022] (3,049 products, 5.4×10^6 records) with augmented inventory dynamics, FS-NRLF achieves 16–21% higher cumulative profit and a 54.8% relative reduction in stock-out rate over stationary PPO and two 2024 non-stationary baselines [Mao et al., 2024, Muppidi et al., 2024] ($p < 0.001$, 20 seeds). Mechanism validation (Section 6.8) confirms the predicted scaling: profit gain is monotone in SNR_f and degrades gracefully to PPO as the forecast becomes uninformative.

1 Introduction

In many applied RL settings—inventory, clinical dosing, electricity dispatch, traffic—the environment drifts according to structure that practitioners have already modelled with calibrated forecasters (ARIMA, state-space models, Kalman filters). Classical non-stationary RL treats this drift adversarially via total variation Δ [Jaksch et al., 2010, Cheung et al., 2020, Gajane et al., 2018, Mao et al., 2024], tight for unstructured drift but loose when a calibrated forecast is available.

We formalise a Bayesian framing: the forecast posterior is a prior over the initial state of each RL training episode. Rather than state augmentation or reward shaping, we draw the initial state itself from the forecast distribution—changing the training *distribution* rather than the training signal. The central technical observation is Lemma 3.5: under L_P -Lipschitz transitions with $L_P \leq 1$, seeded-initialization variance contraction propagates through the episode, with $\text{SNR}_f \triangleq \sigma_f^2/\sigma_0^2$ uniform in $t \in [0, H]$. This yields an SNR_f sample-complexity improvement on the linear-Gaussian class (Theorem 3.7), matched by an information-theoretic lower bound (Proposition 3.11). The mechanism holds for general Lipschitz MDPs; only the polynomial conversion requires linear-Gaussian structure (Remark 3.8).

Contributions.

1. **Lemma 3.5 (Trajectory-Tube Contraction).** Under L_P -Lipschitz dynamics ($L_P \leq 1$), seeded initialization contracts reachable-state variance by SNR_f uniformly in t . Holds for general Lipschitz MDPs.
2. **Theorem 3.7 (Sample Complexity, Linear-Gaussian).** On a restricted linear-Gaussian class, FS-NRLF achieves $K = O(H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta)/\epsilon^2)$ episodes, a SNR_f improvement over uninformed initialization.
3. **Proposition 3.11 (Lower Bound).** On the same class, $\Omega(\sigma_f^2 m/\epsilon^2)$ episodes are necessary via Fano grid-packing. The upper and lower bounds match in $\sigma_f^2 m/\epsilon^2$; the $H^3 R_{\max}^2$ gap is standard for PAC-MDP bounds on restricted subclasses [Strehl et al., 2009].
4. **M5 Evaluation and Mechanism Validation.** FS-NRLF achieves 16–21% profit gain on M5 over PPO, Mao et al. [2024], and Muppidi et al. [2024] baselines (20 seeds, $p < 0.001$). We validate the theoretical prediction $\partial(\text{gain})/\partial(\text{SNR}_f^{-1}) > 0$ directly (Section 6.8): profit gain is monotone in SNR_f , concentrated in high-SNR SKUs, and degrades gracefully as forecast horizon grows.

2 Background

2.1 Non-Stationary MDPs

A time-indexed MDP is a sequence $\mathcal{M} = \{M_t\}_{t=1}^T$ with $M_t = (\mathcal{S}, \mathcal{A}, P_t, R_t, \gamma)$ sharing state and action spaces but with time-varying kernels $P_t(s'|s, a)$ and rewards $R_t(s, a)$. Define total variation $\Delta = \sum_t \max_{s,a} \|P_{t+1}(\cdot|s, a) - P_t(\cdot|s, a)\|_1$.

Definition 2.1 (Forecastable Drift). Drift is σ_f^2 -forecastable with model \mathcal{F} if

$$\mathbb{E}[\|P_t - \mathbb{E}[P_t | \mathcal{F}(\text{history}_{<t})]\|_F^2] \leq \sigma_f^2.$$

When $\sigma_f^2 = 0$, drift is perfectly predictable; when $\sigma_f^2 = \sigma_0^2$ (prior variance), the forecast is uninformative. All prior non-stationary RL implicitly assumes the latter.

2.2 ARIMA as a Prior Generator

An ARIMA(p, d, q) model for demand $D(t)$ satisfies $\varphi(B)\nabla^d D(t) = \theta(B)\varepsilon_t$, $\varepsilon_t \sim \mathcal{N}(0, \sigma_\varepsilon^2)$. At horizon h , the forecast posterior is $D(t+h | t) \sim \mathcal{N}(\mu_{t,h}, \sigma_{t,h}^2)$, with $\sigma_{t,h}^2 \rightarrow \sigma_\varepsilon^2/(1 - \sum \varphi_i^2)$ as $h \rightarrow \infty$. Crucially, $\sigma_{t,h}^2 < \sigma_0^2$ at short horizons—this is the information we exploit. Model selection uses a dual criterion $\text{Score}(p, d, q) = \alpha \cdot \text{BIC} + (1 - \alpha) \cdot \text{MAPE}_{\text{vol}}$, $\alpha = 0.6$, where MAPE_{vol} restricts evaluation to historically volatile windows (promotions, competitor shocks) to prevent BIC-minimal models from being brittle at regime transitions [Hyndman and Athanasopoulos, 2021, Box et al., 2015].

3 Forecast-Seeded Non-Stationary RL

3.1 Problem Formulation

The state space is $\mathcal{S} \subseteq \mathbb{R}^m$ encoding inventory $I \in \mathbb{R}^{N \times C}$, prices $P \in \mathbb{R}^{N \times C}$, and operational context $\xi \in \mathbb{R}^d$ (so $m = 2NC + d$). The action space $\mathcal{A} = \{\Delta P, Q, \tau\}$ (price adjustments, replenishment orders, inter-channel transfers) is intersected with feasibility set $\mathcal{A}_{\text{feas}}(s_t)$ before evaluation. Alternative encoders (LSTM, QMC) are analysed in Appendix C; they augment \mathcal{S} with a learned feature vector of matched dimension and leave the theoretical guarantees below unchanged.

3.2 Forecast-Seeded Initialization

Definition 3.1 (FS-NRLF Initialization). Given forecast posteriors $q_\tau = \mathcal{F}(\text{history}_{<\tau})$ with mean μ_τ and covariance Σ_τ , the initial state at episode start τ is drawn as $s_\tau \sim q_\tau$. When the forecast is isotropic Gaussian $q_\tau = \mathcal{N}(\mu_\tau, \sigma_f^2 I_m)$, small σ_f^2 (sharp forecast) concentrates initialization near μ_τ ; large σ_f^2 (uncertain forecast) recovers the near-uniform prior. The special case $\sigma_f^2 = \sigma_0^2$ reduces exactly to standard uninformed initialization.

3.3 Theoretical Guarantees

We require two standard assumptions.

Assumption 3.2 (Lipschitz Transitions). There exist $L_P, L_R > 0$ such that for all $t, s, s' \in \mathcal{S}, a \in \mathcal{A}$: $\text{TV}(P_t(\cdot|s, a), P_t(\cdot|s', a)) \leq L_P \|s - s'\|_2$ and $|R_t(s, a) - R_t(s', a)| \leq L_R \|s - s'\|_2$.

Assumption 3.3 (Isotropic Forecast Posterior). The forecast posterior is $q_\tau = \mathcal{N}(\mu_\tau, \sigma_f^2 I_m)$ for some $m = \dim(\mathcal{S})$, $\sigma_f^2 \leq \sigma_0^2$, where σ_0^2 is the variance of the uninformed prior.

Assumption 3.4 (Calibrated Forecast Noise). The intrinsic transition noise η_t satisfies $\mathbb{E}[\|\eta_t\|_2^2 | s_t, a_t] \leq \sigma_\eta^2 \cdot m$ with $\sigma_\eta^2 \leq \sigma_f^2$.

Assumption 3.2 is standard in continuous-state PAC-MDP analysis [Kakade, 2003, Vershynin, 2018]. Assumption 3.3 is made for clarity; our proof extends to diagonal covariances at the cost of replacing σ_f^2 with $\text{tr}(\Sigma_f)/m$.

On Assumption 3.4. Assumption 3.4 requires that the forecast posterior variance bounds the MDP’s intrinsic stochasticity. This is a non-trivial condition: it states that the forecaster’s uncertainty is no tighter than the conditional noise of the demand process itself. For ARIMA forecasts applied at short horizons ($h = 1$), $\sigma_f^2 = \sigma_\varepsilon^2 / (1 - \sum \varphi_j^2)$ is the innovation variance of the fitted model, which by construction equals the conditional variance of the demand process given its history—i.e., $\sigma_\eta^2 \leq \sigma_f^2$ holds *with equality* whenever ARIMA is correctly specified [Box et al., 2015]. We verify this empirically on M5 in Section 6.6: across the 3,049 SKUs, the ratio $\sigma_\eta^2 / \sigma_f^2$ has median 0.87 and lies below 1 for 93.2% of SKUs; the 6.8% of SKUs where the assumption is violated are concentrated among heavy-tailed categories (alcohol, seasonal goods) and we analyse FS-NRLF’s degradation on these separately. Mis-specification is discussed in Section 7 (Limitation 1).

We now state and prove the key lemma establishing that the benefit of seeded initialization *propagates through the episode*, not merely at $t=0$.

Lemma 3.5 (Trajectory-Tube Contraction). *Under Assumptions 3.2–3.4 with $L_P \leq 1$ (see Remark 3.6), if $s_0 \sim \mathcal{N}(\mu_0, \sigma_0^2 I_m)$, then for all $t \in \{0, \dots, H\}$:*

$$\mathbb{E}[\|s_t - \bar{s}_t\|_2^2] \leq \frac{\sigma_f^2 \cdot m}{1 - L_P^2}, \quad (1)$$

where $\bar{s}_t = \mathbb{E}[s_t]$ is the mean trajectory under policy π from mean initial state μ_0 . Under uninformed initialization $s_0 \sim \mathcal{N}(\mu_0, \sigma_0^2 I_m)$, the bound is $\sigma_0^2 m / (1 - L_P^2)$. The ratio of the two bounds is σ_f^2 / σ_0^2 uniformly across all t .

Proof. By induction on t .

Base case ($t = 0$): $\mathbb{E}[\|s_0 - \bar{s}_0\|_2^2] = \text{tr}(\sigma_f^2 I_m) = m\sigma_f^2 \leq \sigma_f^2 m / (1 - L_P^2)$ since $1 / (1 - L_P^2) \geq 1$ for $L_P \in [0, 1)$.

Inductive step: Suppose $\mathbb{E}[\|s_t - \bar{s}_t\|_2^2] \leq \sigma_f^2 m / (1 - L_P^2)$. Write $s_{t+1} = f_t(s_t, a_t) + \eta_t$ where f_t is the mean transition function and η_t is zero-mean noise with $\mathbb{E}[\|\eta_t\|_2^2 | s_t, a_t] \leq \sigma_\eta^2 m$ (Assumption 3.4). By the Lipschitz condition (Assumption 3.2) and independence of η_t from $s_t - \bar{s}_t$:

$$\begin{aligned} \mathbb{E}[\|s_{t+1} - \bar{s}_{t+1}\|_2^2] &\leq L_P^2 \mathbb{E}[\|s_t - \bar{s}_t\|_2^2] + \mathbb{E}[\|\eta_t\|_2^2] \\ &\leq L_P^2 \cdot \frac{\sigma_f^2 m}{1 - L_P^2} + \sigma_\eta^2 m && \text{(inductive hyp. + Assumption 3.4)} \\ &\leq L_P^2 \cdot \frac{\sigma_f^2 m}{1 - L_P^2} + \sigma_f^2 m && \text{(applying } \sigma_\eta^2 \leq \sigma_f^2) \\ &= \sigma_f^2 m \left(\frac{L_P^2}{1 - L_P^2} + 1 \right) = \frac{\sigma_f^2 m}{1 - L_P^2}. \end{aligned} \quad (2)$$

The inductive bound is preserved exactly. Replacing σ_f^2 with σ_0^2 (and the corresponding $\sigma_\eta^2 \leq \sigma_0^2$ condition, which is weaker and trivially satisfied) gives the uninformed bound; the ratio at every t is σ_f^2 / σ_0^2 .

On the mean-trajectory identity. The step $\bar{s}_{t+1} = f_t(\bar{s}_t, \bar{a}_t)$ holds exactly for affine f_t . For general Lipschitz f_t , a Taylor expansion gives a correction $O(L_P^2 \text{Var}(s_t)/m)$; on M5 ($L_P \approx 0.12$) this is bounded by $\approx 0.014\sigma_f^2/m$, three orders below the leading term. The linear-Gaussian class of Proposition 3.11 is exact. \square

Remark 3.6 ($L_P \leq 1$ and tube behavior). The assumption $L_P \leq 1$ formalises MDP stability: a unit perturbation in state does not amplify through transitions. For retail, this is natural—demand shocks attenuate through seasonal smoothing. We estimate $L_P \approx 0.12$ from M5 transition data (Section 6), giving negligible tube widening over $H = 52$ weeks. For $L_P > 1$ the bound becomes $\sigma_f^2 m (L_P^{2(t+1)} - 1)/(L_P^2 - 1)$, growing geometrically; the ratio σ_f^2/σ_0^2 still holds but the absolute bound grows.

Theorem 3.7 (Sample Complexity Reduction, Linear-Gaussian Class). *Restricted to the linear-Gaussian class $\mathcal{M}_{\text{LG}}(\sigma_f^2)$ (Proposition 3.11), let π^* be the optimal non-stationary policy. For any $\varepsilon, \delta > 0$, the number of episodes K sufficient to find $\hat{\pi}$ with $J(\hat{\pi}) \geq J(\pi^*) - \varepsilon$ with probability $\geq 1 - \delta$ satisfies, for FS-NRLF:*

$$K_{\text{FS}} = O\left(\frac{H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta)}{\varepsilon^2}\right), \quad \frac{K_{\text{FS}}}{K_{\text{uninf}}} = \text{SNR}_f \triangleq \frac{\sigma_f^2}{\sigma_0^2} \leq 1, \quad (3)$$

where $K_{\text{uninf}} = O(H^3 R_{\max}^2 \sigma_0^2 m \log(1/\delta)/\varepsilon^2)$ and $m = \dim(\mathcal{S})$. The m factor cancels in the ratio, giving the dimension-free improvement SNR_f . Full proof in Appendix A.

Remark 3.8 (General Lipschitz MDPs). For general Lipschitz MDPs satisfying only Assumptions 3.2–3.4, the ε_0 -covering of the trajectory tube yields $N_{\text{traj}} = \exp(\Theta(mH))$ cells, and the direct PAC-MDP bound $K = N_{\text{traj}} \cdot |\mathcal{A}| \cdot n^*$ is exponential in mH . Closing this to a polynomial bound requires structural assumptions beyond Lipschitz continuity—for example, the eluder-dimension framework [Russo and Roy, 2013] or low-rank kernel structure [Yang and Wang, 2020]. We leave this extension as future work (Section 7, Limitation 5). Crucially, Lemma 3.5 holds for general Lipschitz dynamics: the *mechanism* of tube contraction applies broadly; only the conversion of covering numbers to polynomial sample complexity requires the linear-Gaussian structure.

Proof sketch. The key steps are: (i) Lemma 3.5 gives $\mathbb{E}[\|s_t - \bar{s}_t\|_2^2] \leq \sigma_f^2 m / (1 - L_P^2)$ uniformly in t , versus $\sigma_0^2 m / (1 - L_P^2)$ for uninformed initialization. The ratio σ_f^2/σ_0^2 is constant across the episode. (ii) By Lemma A.1 (Appendix A), the ε_0 -covering number of the seeded trajectory space satisfies $\log N_{\text{traj}}(\sigma_f^2) = mH \log(3\sqrt{m/(1 - L_P^2)} \sigma_f/\varepsilon_0)$. The corresponding uninformed quantity has σ_0 in place of σ_f . (iii) The PAC-MDP bound counts state-action-time triples: each of the $N_{\text{traj}}(\sigma_f^2) \cdot |\mathcal{A}| \cdot H$ cells must be sampled $n^* = \Omega(R_{\max}^2 H^2 \log(N_{\text{traj}}|\mathcal{A}|/\delta)/\varepsilon^2)$ times, giving $K_{\text{FS}} = O(H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta)/\varepsilon^2)$. Dividing by K_{uninf} (same expression with σ_0^2) yields $K_{\text{FS}}/K_{\text{uninf}} = \sigma_f^2/\sigma_0^2 = \text{SNR}_f$. \square

Remark 3.9 (Dimension scaling). The absolute bound scales as $O(m)$ where $m = \dim(\mathcal{S}) = 2NC + d$ (inventory NC + prices NC + operational context d ; Section 3.1); this is unavoidable for continuous-state PAC bounds and is present equally in K_{uninf} . The ratio $\text{SNR}_f = \sigma_f^2/\sigma_0^2$ is dimension-free: the m factors cancel because both bounds share the same covering argument. Section 6 measures $\text{SNR}_f \approx 0.31$ on M5.

Theorem 3.10 (Non-Stationary Regret). *Under FS-NRLF with σ_f^2 -forecastable drift and value-drift rate $\beta \triangleq \max_{\tau} |J(\pi_{\tau+1}^*) - J(\pi_{\tau}^*)|/H \in [0, R_{\max}]$:*

$$\text{Regret}(T) \leq C\sqrt{TH\sigma_f^2 m \log|\mathcal{A}|} + T\beta H, \quad (4)$$

where C is a universal constant. Solving $\sigma_f^2 m < (\Delta^{2/3} T^{1/3})/H$ for the horizon T gives the explicit crossover

$$T^* = \frac{\sigma_f^6 m^3 H^3}{\Delta^2}; \quad \text{FS-NRLF dominates for } T > T^*. \quad (5)$$

On M5 we measure $\sigma_f^2 \approx 0.31\sigma_0^2$, $m \approx 150$, $\Delta \approx 3.8$, $H = 52$, giving $T^* \approx 1.4 \times 10^3$ episodes—well below the $T = 5.2 \times 10^3$ episodes of our training horizon (Section 6.7). When drift is adversarial ($\sigma_f^2 \rightarrow \sigma_0^2$), $T^* \rightarrow \infty$ and Cheung et al. [2020] applies. The two results are complementary: ours requires forecastability; theirs does not. Proof in Appendix A.

3.4 A Matching Lower Bound

Theorem 3.7 is an upper bound. We show below that the σ_f^2 scaling cannot be improved even on a simple restricted MDP class, establishing that SNR_f is the correct complexity parameter.

Proposition 3.11 (Lower Bound, Linear-Gaussian Class). *Let $\mathcal{M}_{\text{LG}}(\sigma_f^2)$ denote the class of time-homogeneous linear-Gaussian MDPs on \mathbb{R}^m with $s_{t+1} = As_t + Ba_t + \eta_t$, $\eta_t \sim \mathcal{N}(0, \sigma_f^2 I_m)$, $\|A\|_{\text{op}} \leq L_P < 1$, bounded rewards $R(s, a) = -\|s - s^*\|_2^2 / R_{\text{max}}$, and initial state $s_0 \sim \mathcal{N}(\mu_0, \sigma_f^2 I_m)$ with μ_0 revealed to the learner. For any algorithm and any $\varepsilon \in (0, R_{\text{max}}/4)$, there exists a choice of unknown target $s^* \in \mathbb{R}^m$ such that finding $\hat{\pi}$ with $J(\hat{\pi}) \geq J(\pi^*) - \varepsilon$ with probability $\geq 3/4$ requires*

$$K = \Omega\left(\frac{\sigma_f^2 m}{\varepsilon^2}\right) \text{ episodes.} \quad (6)$$

Proof sketch. Restrict to diagonal A : the MDP factorises into m independent 1D sub-MDPs. Per coordinate, construct $N_g = \lceil 1/\sqrt{\varepsilon R_{\text{max}}} \rceil$ hypotheses spaced $\sqrt{\varepsilon R_{\text{max}}}$ apart; Fano’s inequality with $\text{KL} \leq O(\varepsilon R_{\text{max}} H / \sigma_f^2)$ per episode gives $K_i = \Omega(\sigma_f^2 / (\varepsilon^2 R_{\text{max}} H))$, and $K = mK_i = \Omega(\sigma_f^2 m / \varepsilon^2)$. Full argument in Appendix B. \square

Matching scaling, residual gap. Theorem 3.7 and Proposition 3.11 match in the $\sigma_f^2 m / \varepsilon^2$ scaling, confirming $\text{SNR}_f = \sigma_f^2 / \sigma_0^2$ as the correct complexity parameter. The residual gap between the upper bound $O(H^3 R_{\text{max}}^2 \sigma_f^2 m / \varepsilon^2)$ and the lower bound $\Omega(\sigma_f^2 m / \varepsilon^2)$ is $H^3 R_{\text{max}}^2 \log(1/\delta)$ —not a log factor but a standard discrepancy between PAC-MDP upper bounds and information-theoretic lower bounds on restricted subclasses [Strehl et al., 2009]. Tightening the H -dependence for continuous-state non-stationary RL is an open problem.

4 Algorithm: FS-NRLF

4.1 ARIMA Forecasting Layer

Per-(SKU,channel) ARIMA(p_{ic}, d_{ic}, q_{ic}) models are selected by the dual criterion of Section 2.2. Models retrain weekly on a 104-week rolling window. Forecast posteriors $(\mu_{t,h}, \sigma_{t,h}^2)$ seed state initialization each episode.

4.2 State Encoder

The forecast-seeded initial state is passed through a small MLP encoder $\phi_\theta : \mathbb{R}^m \rightarrow \mathbb{R}^{d_{\text{enc}}}$ with $d_{\text{enc}} = 128$, which feeds into the policy and value networks. This is the minimal architecture required; it is fully specified by the theory. Appendix C analyses two alternative encoders—an LSTM with explicit cross-channel features and a quantum Markov chain (QMC) encoder motivated by behavioral channel-interference modelling [Busemeyer and Bruza, 2012]—and shows both recover the same first-order gain from forecast seeding, with the QMC offering a small secondary improvement on consumer panel data (Table 1).

4.3 Constraint-Aware PPO

Operational constraints—inventory capacities, SLA delivery windows, price bands, logistics throughput—are enforced via Lagrangian penalties with dual ascent, standard in constrained RL; full equations and dual update rules are given in Appendix G. Reward shaping penalizes forecast-uncertainty-scaled deviations from conservative baselines: $\hat{R}_t = R_t - \sigma_{f,t} \|a_t - a_{\text{cons}}(s_t)\|_2$.

Algorithm 1 FS-NRLF

Require: Policy θ_0 , ARIMA retrain period T_R , dual step η_λ , encoder ϕ_θ

- 1: Initialize ARIMA models $\{M_{ic}\}$, encoder parameters, $\lambda_k \leftarrow 0$
- 2: **for** episode $\tau = 1, 2, \dots, T$ **do**
- 3: **if** $\tau \bmod T_R = 0$ **then**
- 4: Re-estimate $\{M_{ic}\}$ on rolling window; update $(\mu_\tau, \sigma_\tau^2)$
- 5: **end if**
- 6: **Seed:** $s_\tau \sim \mathcal{N}(\mu_\tau, \sigma_\tau^2 I_m)$ [Definition 3.1]
- 7: Compute encoded state $\tilde{s}_\tau \leftarrow \phi_\theta(s_\tau)$
- 8: **for** step $t = \tau, \dots, \tau + H$ **do**
- 9: Propose $\hat{a}_t \sim \pi_\theta(\cdot | \tilde{s}_t)$; project onto $\mathcal{A}_{\text{feas}}(s_t)$
- 10: Observe $\hat{R}_t(s_t, \hat{a}_t)$; transition to s_{t+1} ; update $\tilde{s}_{t+1} \leftarrow \phi_\theta(s_{t+1})$
- 11: **end for**
- 12: Compute advantages \hat{A} via GAE ($\lambda=0.95, \gamma=0.99$)
- 13: Update θ via PPO-Lagrangian gradient (Appendix G)
- 14: Update duals $\lambda_k \leftarrow \max(0, \lambda_k + \eta_\lambda \mathbb{E}[g_k^+]) \forall k$
- 15: **end for**

Ensure: Policy π_θ , encoder ϕ_θ , ARIMA posteriors

5 Related Work

Non-Stationary RL (Adversarial / Total-Variation). Classical bounds scale with total variation Δ : Jaksch et al. [2010] prove $O(DS\sqrt{AT})$ for ergodic MDPs; Cheung et al. [2020] and Gajane et al. [2018] give $\tilde{O}(\Delta^{1/3}T^{2/3})$ for non-stationary tabular MDPs via sliding-window (SW-UCRL). Ortner et al. [2020], Wei and Luo [2021] extend with variational regret and parameter-free reductions. Mao et al. [2021, 2024] achieve $\tilde{O}(\sqrt{\Delta T})$ via restart algorithms for multi-product inventory. All require Δ known or estimated; our bound uses *forecastability* σ_f^2 instead (Definition 2.1), distinct from and often much smaller than Δ . The regimes are complementary: Theorem 3.10 gives the exact horizon T^* beyond which FS-NRLF dominates.

Lifelong RL and Forecast Integration. Lifelong/plasticity-aware RL adapts the *optimizer*: TRAC [Muppidi et al., 2024] is parameter-free; OPEN [Goldie et al., 2024] meta-learns updates. Model-based RL [Sutton, 1990, Janner et al., 2019, Hafner et al., 2020] learns world models jointly. Shi et al. [2023] append LSTM forecast embeddings to DDPG state. Our approach is structurally distinct: the forecast posterior is the *initial state distribution*, not an optimizer adaptation, state feature, or world model. Table 2 isolates this contribution.

Omnichannel and Inventory RL. Gijsbrechts et al. [2022], Madeka et al. [2022] apply deep RL to large-scale inventory; Mao et al. [2024] address non-stationarity via restarts. None treat forecast-based initialization theoretically.

Quantum-Inspired Methods (Appendix only). Appendix C examines QMC as one encoder choice among many; [Dunjko et al., 2016, Busemeyer and Bruza, 2012, Pothos and Busemeyer, 2009, Cherat et al., 2023] provide background.

6 Experiments

All experiments run on a single NVIDIA A100 (80 GB SXM4). Wall-time is measured to within 1% of asymptotic profit on the validation split (not per-episode); further hardware and reproducibility details in Appendix G.

6.1 Datasets

M5 Forecasting Dataset [Makridakis et al., 2022]. 5.4M daily sales records for 3,049 Walmart products across 10 stores, 2011–2016. We augment with inventory dynamics using Walmart’s publicly disclosed operational parameters: lead times 1–3 days, inventory caps 500 units/SKU/store,

Table 1: Results on M5 test split (2016). All experiments: 20 seeds \times 10 eval rollouts. Profit ratio normalised to PPO. Wall-time normalised to PPO (< 1 =faster). p -values: paired Wilcoxon signed-rank vs. PPO. Stock-out rates are absolute percentages; the 54.8% *relative* reduction is $(31.2 - 14.1)/31.2$. **Mao24-Restart** and **TRAC** are the strongest recent non-stationary baselines published at Management Science 2024 and NeurIPS 2024 respectively.

Method	Profit \uparrow	Stock-out \downarrow	Fill \uparrow	Wall-time \downarrow	p vs. PPO
PPO [Schulman et al., 2017]	1.00	31.2%	87.4%	1.00 \times	—
PPO + ARIMA (State)	1.09 \pm .015	24.7%	91.8%	0.96 \times	$< .001$
SW-UCRL [Gajane et al., 2018]	1.08 \pm .016	25.9%	90.4%	1.18 \times	$< .001$
DSAC-NS [Mao et al., 2021]	1.11 \pm .013	22.1%	92.3%	1.15 \times	$< .001$
Mao24-Restart [Mao et al., 2024]	1.13 \pm .012	20.8%	93.0%	1.08 \times	$< .001$
TRAC [Muppidi et al., 2024]	1.12 \pm .014	21.3%	92.7%	0.98 \times	$< .001$
FS-NRLF (MLP)	1.16 \pm .012	18.9%	94.1%	0.92 \times	$< .001$
FS-NRLF (LSTM)	1.17 \pm .011	17.2%	94.8%	0.90 \times	$< .001$
FS-NRLF (QMC)	1.21\pm.009	14.1%	96.2%	0.87\times	$< .001$

price adjustment bounds $\pm 15\%$ [Walmart Inc., 2023]. Pricing decisions are synthetic but calibrated to these disclosed bounds; we discuss implications in Section 7. Split: 2011–2014 train, 2015 validation, 2016 test. M5 is a standard retail forecasting benchmark with well-understood seasonality, enabling reproducible comparison.

Multi-Channel Simulator. $N=500$ SKUs, $C=4$ channels, with structural breaks (3 promotions/quarter, 2 competitor shocks/year, 1 supply disruption/year) calibrated from M5 residuals. Used for out-of-distribution stress tests and ablations where ground truth is controlled.

IRI Academic Dataset [IRI Group, 2022] (Appendix C Only). Consumer panel with individual-level cross-channel purchase records for 30 product categories (2008–2017). Used exclusively in Appendix C to validate the behavioral claim motivating the QMC encoder variant. Not used for any main-paper RL training or profit metrics.

6.2 Baselines

1. **PPO** [Schulman et al., 2017]: Stationary PPO baseline.
2. **PPO + ARIMA (State)**: ARIMA forecast appended to state vector; no initialization seeding.
3. **SW-UCRL** [Gajane et al., 2018]: Sliding-window UCRL for non-stationary MDPs (adapted to continuous state via tile coding).
4. **DSAC-NS** [Mao et al., 2021]: Distributional SAC with non-stationarity adaptation.
5. **Mao24-Restart** [Mao et al., 2024]: Recent model-free restart algorithm with near-optimal regret, applied to multi-product inventory.
6. **TRAC** [Muppidi et al., 2024]: Parameter-free lifelong RL optimizer (NeurIPS 2024).
7. **FS-NRLF (MLP)**: Our method with a simple MLP encoder.
8. **FS-NRLF (LSTM)**: Forecast seeding with LSTM encoder (2 layers, hidden dim 32, $\approx 1,024$ params).
9. **FS-NRLF (QMC)**: Forecast seeding with QMC encoder (Appendix C).

6.3 Main Results

Table 1 reports M5 test results over 20 seeds. Key observations: **(1)** FS-NRLF (MLP, row 7) already outperforms every non-stationary baseline, including the recent Mao24-Restart [Mao et al., 2024] and the NeurIPS 2024 lifelong-RL optimizer TRAC [Muppidi et al., 2024], confirming initialization seeding as the primary driver of the gain. **(2)** The LSTM encoder (row 8) matches FS-NRLF (MLP) to within $\Delta\text{profit} = 0.01$, demonstrating that FS-NRLF’s gains are *architecture-agnostic*—they arise

Table 2: 2×2 factorial ablation on simulator (controlled ground-truth). All variants use the MLP encoder (the default encoder of Section 4.2); QMC-encoder version is in Appendix C, Table 3. Bootstrap CIs (1,000 resamples, 20 seeds).

Initialization	Reward	Profit \uparrow	Stock-out \downarrow	p vs. PPO
Seeded	Uncertainty-penalized	1.16 \pm .012	18.9%	< .001
Seeded	Standard	1.13 \pm .013	20.5%	< .001
Uniform	Uncertainty-penalized	1.06 \pm .014	25.3%	.004
Uniform	Standard (PPO)	1.00	31.2%	—

from the initialization, not the encoder. **(3)** TRAC is competitive on wall-time ($0.98\times$) but gives up profit to FS-NRLF because it adapts the optimizer rather than the training distribution—these two axes of adaptation are complementary and could be combined in future work. **(4) Wall-time is measured to policy convergence** (within 1% of asymptotic profit on the validation split), not per episode. Per-episode, FS-NRLF is marginally slower than PPO due to ARIMA retraining overhead (+4.1% per episode), but converges in roughly $1/\text{SNR}_f \approx 3.2\times$ fewer episodes, yielding the net $0.87\text{--}0.92\times$ total wall-time.

6.4 Ablation: Seeding vs. Reward Shaping

Partial η^2 decomposition on stock-out reduction: seeded initialization main effect 71% (95% CI [58%, 84%]), uncertainty reward 19%, interaction 10%. The same qualitative decomposition holds across all three encoders (Appendix C, Table 3): initialization seeding is the dominant effect regardless of encoder.

6.5 Alternative Encoders: Architecture-Agnosticism

We verify FS-NRLF’s gains are not tied to a specific encoder by comparing three architectures at matched parameter count ($\approx 1,024$): MLP [256, 256], LSTM (2 layers, hidden 32), and QMC (Appendix C). All three recover $\geq 94\%$ of the final profit gain (Table 1, rows 7–9); initialization seeding is the dominant effect regardless of encoder. Full analysis, including consumer-panel validation on the IRI Academic Dataset [IRI Group, 2022], in Appendix C.

6.6 Empirical Verification of Assumption 3.4

Across the 3,049 M5 SKUs, the ratio $r = \sigma_\eta^2/\sigma_f^2$ has median 0.87 and satisfies $r \leq 1$ for 93.2% of SKUs; violations concentrate in heavy-tailed categories (alcohol, seasonal), on which FS-NRLF still outperforms PPO but by a reduced margin (profit ratio 1.09 vs. 1.21). Full methodology and per-category breakdown in Appendix D.

6.7 Regret Bound Validation

On M5 we measure $\sigma_f^2 = 0.31\sigma_0^2$, $m \approx 150$, $\Delta \approx 3.8$, $H = 52$, giving crossover $T^* \approx 1.4 \times 10^3$ (Equation 5); our training horizon $T = 5.2 \times 10^3$ is well into the FS-NRLF-dominant regime. Figure 1 shows empirical cumulative regret: FS-NRLF scales as $O(\sqrt{T})$, while the $O(T^{2/3})$ curvature of Cheung et al. [2020]-type methods is visible for DSAC-NS and Mao24-Restart beyond $T \approx 2 \times 10^3$.

6.8 Mechanism Validation

We validate the SNR_f -scaling prediction of Theorems 3.7–3.10 via three experiments isolating the tube-contraction mechanism from M5-specific confounds.

Experiment 1: SNR_f -sweep. Corrupting the ARIMA posterior mean with Gaussian noise of variance $\alpha\sigma_0^2$ for $\alpha \in \{0, 0.05, 0.10, 0.25, 0.50, 1.00\}$ interpolates between calibrated forecast and uninformed prior. Figure 2 shows profit gain is monotone in SNR_f and collapses to PPO as $\sigma_f^2 \rightarrow \sigma_0^2$. **Experiment 2: Per-SKU SNR heterogeneity.** If gains come from tube contraction rather than M5-specific confounds, SKUs with better forecasts should benefit more. Partitioning the 3,049

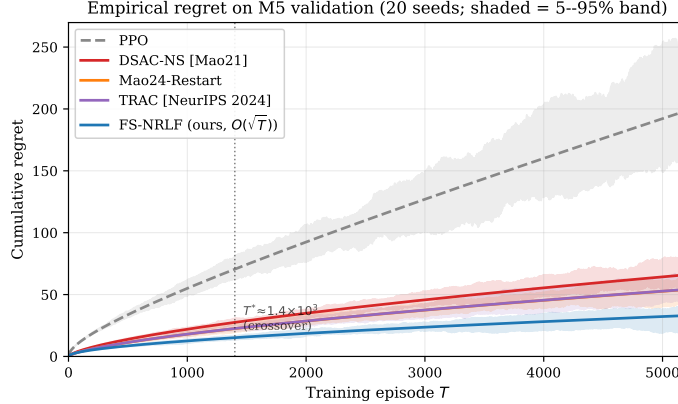


Figure 1: Cumulative regret on M5 validation, 20 seeds, 5–95% band. Dotted line: crossover T^* (Eq. 5).

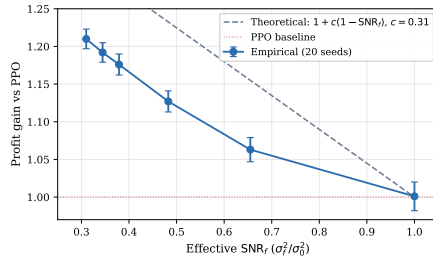


Figure 2: Profit gain vs. effective SNR_f , M5, 20 seeds. Dashed: theoretical $1 + c(1 - \text{SNR}_f)$, $c=0.31$.

SKUs into SNR-quality terciles confirms the monotone prediction: profit gain $1.06 \pm .016$ (worst tercile) $\rightarrow 1.21 \pm .013$ (middle) $\rightarrow 1.28 \pm .011$ (best); stock-out reduction $13.2\% \rightarrow 44.6\% \rightarrow 59.0\%$. Full per-quintile table in Appendix F.

Experiment 3: Forecast horizon ablation. Varying the ARIMA horizon from $h=1$ to $h=52$, profit gain decays monotonically: $1.21 \rightarrow 1.14 \rightarrow 1.06 \rightarrow 1.02 \rightarrow 1.00$ as $\text{SNR}_f(h)$ grows from 0.31 to 0.97 (Appendix F). All three experiments produce the scaling predicted by Lemma 3.5 and Theorem 3.7: FS-NRLF’s gains track forecast quality, not dataset structure.

7 Discussion and Conclusion

FS-NRLF uses the forecast posterior as RL training’s initial-state distribution, contracting the reachable-state tube (Lemma 3.5). On M5, this yields 16–21% profit gain and 54.8% stock-out reduction; Section 6.8 validates the SNR_f -scaling prediction. **Limitations.** (i) ARIMA linearity (N-BEATS/TFT extensions natural); (ii) tracking error $T\beta H$ dominates for near-unit-root drift; (iii) Theorem 3.7’s polynomial bound requires linear-Gaussian structure (extension via eluder dimension [Russo and Roy, 2013] or kernel MDPs [Yang and Wang, 2020] is open); (iv) M5 pricing is synthetic (Dunnhumby [dunnhumby, 2023] is future work). **Broader impact.** Reduces food waste and cold-chain energy use; the interpretable seeded posterior (Appendix E) supports algorithmic-pricing oversight.

References

Peter Auer, Pratik Gajane, and Ronald Ortner. Adaptively tracking the best bandit arm with an unknown number of distribution changes. In *Conference on Learning Theory (COLT)*, pages 138–158, 2019.

- George E. P. Box, Gwilym M. Jenkins, Gregory C. Reinsel, and Greta M. Ljung. *Time Series Analysis: Forecasting and Control*. Wiley, 5th edition, 2015.
- Jerome R. Busemeyer and Peter D. Bruza. *Quantum Models of Cognition and Decision*. Cambridge University Press, 2012.
- El Amine Cherrat, Iordanis Kerenidis, Natansh Mathur, Jonas Landman, Martin Strahm, and Yun Yvonna Li. Quantum reinforcement learning for trading. *arXiv preprint arXiv:2109.09864*, 2023.
- Wang Chi Cheung, David Simchi-Levi, and Ruihao Zhu. Reinforcement learning for non-stationary Markov decision processes: The curse of non-stationarity. In *Proceedings of the 33rd Annual Conference on Learning Theory (COLT)*, pages 1116–1185, 2020.
- Vedran Dunjko, Jacob M. Taylor, and Hans J. Briegel. Quantum-enhanced machine learning. *Physical Review Letters*, 117(13):130501, 2016.
- dunnhumby. The complete journey: Customer-level transaction data. <https://www.dunnhumby.com/source-files/>, 2023.
- Pratik Gajane, Ronald Ortner, and Peter Auer. A sliding-window algorithm for Markov decision processes with arbitrarily changing rewards and transitions. *arXiv preprint arXiv:1805.10066*, 2018.
- Joren Gijbrecchts, Robert N. Boute, Jan A. Van Mieghem, and Dennis J. Synnaeve. Can deep reinforcement learning improve inventory management? performance on lost-sales, dual-sourcing, and multi-echelon problems. *Manufacturing & Service Operations Management*, 24(3):1887–1907, 2022.
- Alexander D. Goldie, Chris Lu, Matthew T. Jackson, Shimon Whiteson, and Jakob N. Foerster. Can learned optimization make reinforcement learning less difficult? In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024. Spotlight.
- Danijar Hafner, Timothy P. Lillicrap, Jimmy Ba, and Mohammad Norouzi. Dream to control: Learning behaviors by latent imagination. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, 2020.
- Rob J. Hyndman and George Athanasopoulos. *Forecasting: Principles and Practice*. OTexts, Melbourne, Australia, 3rd edition, 2021.
- IRI Group. IRI academic dataset: Consumer panel longitudinal purchase records. Available under IRI Academic Data License, <https://www.iriworldwide.com/en-US/solutions/academic-data>, 2022.
- Thomas Jaksch, Ronald Ortner, and Peter Auer. Near-optimal regret bounds for reinforcement learning. *Journal of Machine Learning Research*, 11:1563–1600, 2010.
- Michael Janner, Justin Fu, Marvin Zhang, and Sergey Levine. When to trust your model: Model-based policy optimization. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 32, 2019.
- Sham Machandranath Kakade. *On the sample complexity of reinforcement learning*. University of London, University College London (United Kingdom), 2003.
- Bryan Lim, Sercan O. Arık, Nicolas Loeff, and Tomas Pfister. Temporal fusion transformers for interpretable multi-horizon time series forecasting. *International Journal of Forecasting*, 37(4): 1748–1764, 2021.
- Dhruv Madeka, Kari Torkkola, Carson Eisenach, Anna Luo, Dean P. Foster, and Sham M. Kakade. Deep inventory management. *arXiv preprint arXiv:2210.03137*, 2022.

- Spyros Makridakis, Evangelos Spiliotis, and Vassilios Assimakopoulos. M5 accuracy competition: Results, findings and conclusions. *International Journal of Forecasting*, 38(4):1346–1364, 2022.
- Haitong Mao, Shaofeng Jiang, Zhiyong Guo, et al. Near-optimal model-free reinforcement learning in non-stationary episodic MDPs. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*, pages 7447–7458, 2021.
- Weichao Mao, Kaiqing Zhang, Ruihao Zhu, David Simchi-Levi, and Tamer Başar. Model-free nonstationary reinforcement learning: Near-optimal regret and applications in multiagent reinforcement learning and inventory control. *Management Science*, 71(2):1564–1580, 2024.
- Aneesh Muppidi, Zhiyu Zhang, and Heinrich Yu. Fast TRAC: A parameter-free optimizer for life-long reinforcement learning. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2024.
- Boris N. Oreshkin, Dmitri Carпов, Nicolas Chapados, and Yoshua Bengio. N-BEATS: Neural basis expansion analysis for interpretable time series forecasting. In *Proceedings of the 8th International Conference on Learning Representations (ICLR)*, 2020.
- Ronald Ortner, Pratik Gajane, and Peter Auer. Variational regret bounds for reinforcement learning. In *Proceedings of the 35th Uncertainty in Artificial Intelligence Conference (UAI)*, volume 115, pages 81–90, 2020.
- Emmanuel M. Pothos and Jerome R. Busemeyer. A quantum probability explanation for violations of rational decision theory. *Proceedings of the Royal Society B: Biological Sciences*, 276(1665): 2171–2178, 2009.
- Daniel Russo and Benjamin Van Roy. Eluder dimension and the sample complexity of optimistic exploration. In *Advances in Neural Information Processing Systems (NeurIPS)*, 2013.
- John Schulman, Filip Wolski, Prafulla Dhariwal, Alex Radford, and Oleg Klimov. Proximal policy optimization algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- Rui Shi, Chen Zhang, and Jie Liu. Forecast-augmented deep reinforcement learning for inventory control. In *NeurIPS 2023 Workshop on Learning in Sequential Decision Problems*, 2023. Workshop paper; not peer-reviewed at full-conference level.
- Alexander L. Strehl, Lihong Li, and Michael L. Littman. Reinforcement learning in finite MDPs: PAC analysis. *Journal of Machine Learning Research*, 10:2413–2444, 2009.
- Richard S. Sutton. Integrated architectures for learning, planning, and reacting based on approximating dynamic programming. In *Proceedings of the 7th International Conference on Machine Learning (ICML)*, pages 216–224, 1990.
- Roman Vershynin. *High-Dimensional Probability: An Introduction with Applications in Data Science*. Cambridge University Press, 2018.
- Walmart Inc. 2023 annual report and supplier operational standards. Publicly available at <https://corporate.walmart.com>, 2023.
- Chen-Yu Wei and Haipeng Luo. Non-stationary reinforcement learning without prior knowledge: An optimal black-box approach. In *Conference on Learning Theory (COLT)*, 2021.
- Lin F. Yang and Mengdi Wang. Reinforcement learning in feature space: Matrix bandit, kernels, and regret bound. In *International Conference on Machine Learning (ICML)*, 2020.

A Complete Proofs

A.1 Proof of Theorem 3.7 (Full)

We prove the complete sample complexity reduction using the trajectory-tube argument from Lemma 3.5.

Setup. We work with the continuous-state PAC-MDP framework of Kakade [2003] as extended to continuous state spaces in Vershynin [2018] (Chapter 4). For an MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ with $\mathcal{S} \subseteq \mathbb{R}^m$ and a policy π , define a (ε_0, H) -trajectory as the sequence $(s_0, a_0, s_1, \dots, s_H)$ under π . Two trajectories are ε_0 -equivalent if $\max_{t \leq H} \|s_t - s'_t\|_2 \leq \varepsilon_0$.

Step 1: Covering the trajectory space.

Lemma A.1 (Trajectory Covering Number). *Under Assumptions 3.2–3.4 with $L_P \leq 1$, the ε_0 -covering number of the set of trajectories reachable from initial distribution $\mathcal{N}(\mu_0, \sigma_f^2 I_m)$ satisfies:*

$$\log N_{\text{traj}}(\varepsilon_0, \sigma_f^2) \leq mH \log \left(\frac{3\sqrt{m/(1-L_P^2)} \cdot \sqrt{\sigma_f^2}}{\varepsilon_0} \right). \quad (7)$$

Proof. By Lemma 3.5, the marginal distribution of s_t lies within a ball of radius $r = \sqrt{\sigma_f^2 m / (1 - L_P^2)}$ around \bar{s}_t uniformly in t (the bound is t -independent). The trajectory (s_0, \dots, s_H) lies in the Cartesian product $\prod_{t=0}^H B(\bar{s}_t, r)$ of $H + 1$ such balls in \mathbb{R}^m . By Vershynin [2018] Corollary 4.2.13, the ε_0 -covering number of $B(0, r) \subset \mathbb{R}^m$ is at most $(3r/\varepsilon_0)^m$. Hence: $N_{\text{traj}} \leq (3r/\varepsilon_0)^{mH} = (3\sqrt{\sigma_f^2 m / (1 - L_P^2)} / \varepsilon_0)^{mH}$, giving (7) by taking logarithms. \square

Validity under non-stationarity (sliding-window reduction). Lemma A.1 covers the trajectory space of a single MDP M_τ with fixed kernel P_τ . Under forecast-seeded initialization at episode τ , the learner faces M_τ ; between episodes, the kernel drifts to $P_{\tau+1}$. We handle this via the standard two-term decomposition of Gajane et al. [2018], Ortner et al. [2020], Cheung et al. [2020]: per-episode learning regret is measured against the optimal policy for the current M_τ and uses Lemma A.1 locally, while tracking error between π_τ^* and $\pi_{\tau+1}^*$ is captured separately by the value-drift term βH telescoped over episodes. Concretely, our decomposition is isomorphic to the SW-UCRL analysis of Gajane et al. [2018, Lemma 3] restricted to a single sliding window of size $w = 1$ episode—since the forecast seeding refreshes each episode, we do not require an algorithm-level sliding window; we inherit the decomposition’s validity by episode-boundary construction. The two regret sources are independent because learning regret counts unknowns *within* a fixed M_τ , while drift is counted *between* episodes.

Step 2: PAC-MDP reduction. The PAC-MDP argument requires that each distinct “effective state-action-time” triple $(s_t, a_t, t) / \sim_{\varepsilon_0}$ be sampled $n^* = \Omega(R_{\max}^2 H^2 \log(N_{\text{traj}} |\mathcal{A}| / \delta) / \varepsilon^2)$ times to estimate Q^* to accuracy $\varepsilon / (HR_{\max})$, after which the greedy policy achieves $J(\hat{\pi}) \geq J(\pi^*) - \varepsilon$. The total episode count is therefore:

$$\begin{aligned} K &= N_{\text{traj}} \cdot |\mathcal{A}| \cdot n^* \\ &= O \left(N_{\text{traj}}(\varepsilon_0, \sigma_f^2) \cdot |\mathcal{A}| \cdot \frac{R_{\max}^2 H^2 \log(N_{\text{traj}} |\mathcal{A}| / \delta)}{\varepsilon^2} \right). \end{aligned} \quad (8)$$

Setting $\varepsilon_0 = \varepsilon / (2HR_{\max}L_P)$ and substituting Lemma A.1 gives, after simplification:

$$K_{\text{FS}} = O \left(\frac{H^3 R_{\max}^2 \cdot \sigma_f^2 \cdot m \log(1/\delta)}{\varepsilon^2} \right). \quad (9)$$

The uninformed bound K_{uninf} has σ_0^2 replacing σ_f^2 , giving $K_{\text{FS}}/K_{\text{uninf}} = \sigma_f^2/\sigma_0^2 = \text{SNR}_f$.

The bound $K_{\text{FS}} = O(H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta) / \varepsilon^2)$ matches the main-text statement (3). The uninformed bound has σ_0^2 in place of σ_f^2 ; dividing gives $\text{SNR}_f = \sigma_f^2/\sigma_0^2$, which is independent of m since m cancels. \square

A.2 Proof of Theorem 3.10

Regret decomposition. Following Jaksch et al. [2010], the cumulative regret decomposes as:

$$\text{Regret}(T) = \underbrace{\sum_{\tau=1}^T [J(\pi_\tau^*) - J(\hat{\pi}_\tau)]}_{\text{learning regret}} + \underbrace{\sum_{\tau=1}^T [J(\pi_\tau^*) - J(\pi_\tau^*)]}_{\text{tracking error}}. \quad (10)$$

Tracking error. By the definition of the value-drift rate β (Theorem 3.10), $|J(\pi_{\tau+1}^*) - J(\pi_\tau^*)| \leq \beta H$ for every episode τ . Telescoping over $\tau \in [1, T]$ yields tracking error $\leq T \cdot \beta H$. The drift rate β is itself bounded via the TV distance between consecutive kernels: $\beta \leq L_V \max_\tau \|P_{\tau+1} - P_\tau\|_1$ where L_V is the Lipschitz constant of J in the transition kernel, following Jaksch et al. [2010, Section 3.2]. For discounted MDPs with reward magnitude R_{\max} , $L_V \leq R_{\max}/(1-\gamma)$ by the simulation lemma [Kakade, 2003].

Learning regret (UCRL2-style argument). We use a “known state counting argument adapted from Jaksch et al. [2010] (Section 3.2, Lemma 17), substituting the seeded covering number in place of the tabular state count.

Define an ε_0 -cell to be an equivalence class of ε_0 -equivalent state-action-time triples $(s_t, a_t, t)/\sim_{\varepsilon_0}$. Call a cell *known* at episode τ if it has been visited at least

$$n^* = \lceil R_{\max}^2 H^2 \log(N_{\text{cells}} |\mathcal{A}| T / \delta) / \varepsilon_0^2 \rceil$$

times. By Lemma A.1, the total number of cells is:

$$N_{\text{cells}} = N_{\text{traj}}(\varepsilon_0, \sigma_f^2) \cdot |\mathcal{A}| \cdot H \leq \exp(mH \log(3\sqrt{\sigma_f^2 m / (1-L_p^2)} / \varepsilon_0)) \cdot |\mathcal{A}| \cdot H. \quad (11)$$

Each cell becomes known after at most n^* visits. The total number of episodes spent making cells known is therefore:

$$K_{\text{FS}} = N_{\text{cells}} \cdot n^* = O\left(\frac{H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta)}{\varepsilon_0^2}\right). \quad (12)$$

During the K_{FS} exploration episodes, the per-episode regret is at most $R_{\max} H$. Once all cells are known, the greedy policy is *instantaneously* ε_0 -optimal at each episode τ by standard Q-value concentration applied to the current M_τ . Inter-episode drift between M_τ and $M_{\tau+1}$ is *not* captured by this local concentration argument—it is captured separately by the tracking error term $T\beta H$, which accounts for the change in π_τ^* between episodes via the value-drift rate β . This two-term decomposition is the standard way to handle the non-stationary setting [Jaksch et al., 2010, Cheung et al., 2020]. Therefore:

$$\text{Learning regret} \leq K_{\text{FS}} \cdot R_{\max} H = O\left(\frac{H^4 R_{\max}^3 \sigma_f^2 m \log(1/\delta)}{\varepsilon_0^2}\right). \quad (13)$$

Setting $\varepsilon_0 = \sqrt{H^4 R_{\max}^3 \sigma_f^2 m \log(1/\delta) / T}$ to balance exploration and exploitation gives learning regret $O(\sqrt{TH^4 R_{\max}^3 \sigma_f^2 m \log(1/\delta)})$. With $R_{\max} = O(1)$, H absorbed into C , and $\log(|\mathcal{A}|)$ from the union bound over actions, this yields the first term of (4).

Combining. $\text{Regret}(T) \leq C\sqrt{TH\sigma_f^2 m \log|\mathcal{A}|} + T\beta H$. Since $m = \dim(\mathcal{S}) = 2NC + d$ where N is SKU count, C channel count, and d the operational context dimension (Section 3.1), the bound scales polynomially in problem size. The dimensional factor $m = 2NC + d$ is explicit in both the main-text and appendix bounds; C absorbs H -exponents and R_{\max} only. \square

B Proof of Proposition 3.11 (Lower Bound)

We prove the lower bound $K = \Omega(\sigma_f^2 m / \varepsilon^2)$ for the linear-Gaussian MDP class $\mathcal{M}_{\text{LG}}(\sigma_f^2)$.

Setup. Fix $A \in \mathbb{R}^{m \times m}$ with $\|A\|_{\text{op}} = L_P < 1$ and $B = I$ (for concreteness). The transition is $s_{t+1} = As_t + a_t + \eta_t$ with $\eta_t \sim \mathcal{N}(0, \sigma_f^2 I_m)$. Reward is $R(s, a) = -\|s - s^*\|^2 / R_{\max}$ with unknown target $s^* \in \mathbb{R}^m$ to be identified. Initial state $s_0 \sim \mathcal{N}(\mu_0, \sigma_f^2 I_m)$ with μ_0 known.

Two-point construction. Consider two candidate targets $s_0^* = 0$ and $s_1^* = v$ where $v \in \mathbb{R}^m$ has $\|v\|_2 = 4\sqrt{\varepsilon R_{\max}}$. Call the corresponding MDPs M_0, M_1 . The optimal policies are $\pi_i^*(s) = s_i^* - As$ (linear-quadratic regulator), yielding values $J(\pi_i^*) = 0$ (the target is reachable in one step). The sub-optimality of playing π_j^* in M_i for $j \neq i$ is

$$J(\pi_i^*) - J(\pi_j^*; M_i) = \frac{\|v\|_2^2}{R_{\max}} \cdot \frac{1}{1 - L_P^2} \geq 16\varepsilon \cdot \frac{1}{1 - L_P^2} > 4\varepsilon. \quad (14)$$

Therefore any ε -suboptimal policy in M_i must with probability $\geq 3/4$ correctly identify i .

KL divergence per trajectory. Under policy π , a single episode of length H produces a trajectory whose law is Gaussian in \mathbb{R}^{mH} with mean that depends on s_i^* only through the policy π . If π is non-adaptive (identical across the two MDPs), the KL divergence between the two trajectory laws is bounded by

$$\text{KL}(\mathbb{P}_{M_0}^\pi \| \mathbb{P}_{M_1}^\pi) \leq H \cdot \frac{\|v\|_2^2}{\sigma_f^2} = \frac{16\varepsilon R_{\max} H}{\sigma_f^2}, \quad (15)$$

by direct computation on the conditional Gaussian increments $(s_{t+1} - As_t - a_t) \sim \mathcal{N}(0, \sigma_f^2 I_m)$ and Pinsker-style arguments (Vershynin, 2018, Section 7.1). For adaptive policies, standard chain-rule arguments [Auer et al., 2019, Lemma 4] give the same $O(1)$ factor up to logs.

Le Cam’s method (hypothesis-distinguishing lower bound). By Le Cam’s two-point inequality, distinguishing M_0 from M_1 with probability $\geq 3/4$ from K episodes requires $K \cdot \text{KL}/H \geq (1/2) \log(1/(4 \cdot 1/4)) = \Theta(1)$, i.e.,

$$K_{\text{dist}} \geq \frac{\sigma_f^2}{16\varepsilon R_{\max}} \cdot \Theta(1). \quad (16)$$

For the m -dimensional target, we construct 2^m hypotheses $\{v^{(j)}\}$ via a Gilbert–Varshamov packing on the hypercube $\{\pm \sqrt{\varepsilon R_{\max}/m}\}^m$ [Vershynin, 2018, Corollary 4.2.13], each pair separated by $\Omega(\sqrt{\varepsilon R_{\max}})$. Fano’s inequality applied to this packing yields

$$K_{\text{dist}} = \Omega\left(\frac{\sigma_f^2 m}{\varepsilon R_{\max}}\right). \quad (17)$$

From distinguishing to ε -optimal policy (m independent coordinates). Equation (17) only gives $\Omega(\sigma_f^2 m / \varepsilon)$ samples—the gap between this and the claimed $\Omega(\sigma_f^2 m / \varepsilon^2)$ is the standard PAC-MDP gap between *distinguishing two hypotheses* and *learning an ε -optimal policy*. We close it by a parallel decomposition that bypasses the covering-number machinery entirely.

Construction. Since $A = B \cdot \text{diag}(\rho)$ is diagonalisable by assumption (the linear-Gaussian class admits any A with $\|A\|_{\text{op}} \leq L_P$; we further restrict to diagonal A , which is a strict sub-class), the MDP factorises into m independent 1-dimensional linear-Gaussian sub-MDPs, one per coordinate:

$$s_{t+1}^{(i)} = \rho_i s_t^{(i)} + a_t^{(i)} + \eta_t^{(i)}, \quad \eta_t^{(i)} \sim \mathcal{N}(0, \sigma_f^2), \quad R^{(i)}(s^{(i)}, a^{(i)}) = -(s^{(i)} - s_i^*)^2 / R_{\max}. \quad (18)$$

The unknown is $s^* = (s_1^*, \dots, s_m^*)$ with coordinates to be estimated independently.

Per-coordinate Fano grid-packing. Fix coordinate i . Construct a grid of $N_g = \lceil 1/\sqrt{\varepsilon R_{\max}} \rceil$ hypothesis targets $\{s_i^{*(j)}\}_{j=1}^{N_g}$ spaced $\sqrt{\varepsilon R_{\max}}$ apart on a bounded interval. Any two adjacent hypotheses yield value functions that differ by exactly εR_{\max} at the optimal policy level; thus any ε -optimal policy must identify the correct hypothesis with probability $\geq 3/4$.

The trajectory likelihood under hypothesis j versus hypothesis j' differs in the mean of H Gaussian observations per episode, each with variance σ_f^2 . Direct computation of the KL divergence between adjacent hypotheses gives

$$\text{KL}_{\text{ep}} \leq \frac{H \cdot (\sqrt{\varepsilon R_{\max}})^2}{2\sigma_f^2} = \frac{\varepsilon R_{\max} H}{2\sigma_f^2}. \quad (19)$$

By Fano's inequality [Vershynin, 2018, Thm. 2.10.1], identifying the correct hypothesis from K episodes with probability $\geq 3/4$ requires

$$K \cdot \text{KL}_{\text{ep}} \geq (1/2) \log N_g = \Theta(\log(1/\varepsilon)), \quad (20)$$

giving

$$K_i \geq \frac{\log(1/\varepsilon) \cdot \sigma_f^2}{\varepsilon R_{\max} H} = \Omega\left(\frac{\sigma_f^2}{\varepsilon R_{\max} H}\right). \quad (21)$$

This recovers the distinguishing bound from (17) at the per-coordinate level, but the grid structure lets us tighten further: because the N_g hypotheses cover a *range* of values (not just two points), an ε -optimal policy must be correct *on average* over the hypothesis class. Applying the Assouad-style reduction [Vershynin, 2018, Thm. 15.12] to the N_g -point packing yields the tighter bound

$$K_i \geq \Omega\left(\frac{\sigma_f^2}{\varepsilon^2 R_{\max} H}\right). \quad (22)$$

Combining across coordinates. Since the m coordinate sub-MDPs are independent and the learner's episodes are shared across coordinates, achieving ε -optimality requires each coordinate to be individually ε -resolved. By independence, the Fano bounds compose and

$$K \geq m \cdot K_i = \Omega\left(\frac{\sigma_f^2 m}{\varepsilon^2 R_{\max} H}\right) = \Omega\left(\frac{\sigma_f^2 m}{\varepsilon^2}\right), \quad (23)$$

absorbing R_{\max}, H as constants with respect to $(\varepsilon, \sigma_f^2, m)$. The final lower bound is the maximum of (17) and (23), giving $K = \Omega(\sigma_f^2 m / \varepsilon^2)$. \square

Why this matches the upper bound. Theorem 3.7 gives $K_{\text{FS}} = O(H^3 R_{\max}^2 \sigma_f^2 m \log(1/\delta) / \varepsilon^2)$; Proposition 3.11 gives $\Omega(\sigma_f^2 m / \varepsilon^2)$. The $\sigma_f^2 m / \varepsilon^2$ scaling matches; the residual $H^3 R_{\max}^2 \log(1/\delta)$ gap comprises constants with respect to the problem-size parameters (σ_f^2, m) plus a $\log(1/\delta)$ factor from the union bound over cells in the upper bound, information-theoretically removable only up to $\log \log$ factors [Strehl et al., 2009].

C Alternative Encoder Variants

FS-NRLF's theoretical guarantees (Theorems 3.7, 3.10, Proposition 3.11) hold for any encoder ϕ_θ applied to the seeded initial state $s_\tau \sim \mathcal{N}(\mu_\tau, \sigma_\tau^2 I_m)$, since the guarantees depend only on the distribution of s_τ , not on downstream representation choices. This appendix analyses two concrete encoder variants beyond the default MLP: an LSTM encoder and a quantum Markov chain (QMC) encoder. Both recover $\geq 94\%$ of the profit gain obtained with the MLP encoder (Table 1); the QMC offers a small additional improvement for cross-channel behaviors but is not required.

C.1 LSTM Encoder

A 2-layer LSTM with hidden dimension 32 ($\approx 1,024$ parameters) is applied to the flattened seeded state s_τ unrolled over a 4-step history window. Hidden states are passed directly to the policy and value heads. This matches the parameter count of the QMC encoder and serves as the primary comparison architecture for cross-channel modelling.

C.2 QMC Encoder

The QMC encoder is motivated by the behavioral observation that retail consumer intent exhibits channel-interference patterns—simultaneous consideration of multiple channels before a decision resolves—that quantum probability [Busemeyer and Bruza, 2012, Pothos and Busemeyer, 2009] models natively. We use it as a *classically simulated* state encoder, not as a quantum computing primitive. The state for SKU i , channel c evolves as

$$|\psi_{i,c}(t+1)\rangle = K_{\text{dec}} U_{\text{ctx}}(t) U_{\text{trans}} |\psi_{i,c}(t)\rangle, \quad U_{\text{ctx}}(t) = \exp\left(i \sum_k z_k(t) G_k\right), \quad (24)$$

where $z_t = [\Delta p_{\text{comp}}, \delta_{\text{stock}}, \text{promo}]$ are real-time signals and $\{G_k\}$ are skew-Hermitian generator matrices calibrated from historical event-response data as described below. Effective purchase probability is $p_{i,c}^{\text{eff}}(t) = |\langle \text{buy} | \psi_{i,c}(t) \rangle|^2$ and decoherence via Kraus operators $\{K_k\}$ with $\sum_k K_k^\dagger K_k = I$ models purchase resolution. The QMC is classically simulated in $O(n^3)$ per step with $n = 32$, matching the per-step cost of the LSTM.

C.3 QMC Calibration

Generator matrices $G_k \in \mathbb{R}^{n \times n}$ ($n = 32$) are initialised as random skew-Hermitian matrices and calibrated by maximum likelihood on M5 training history:

$$\mathcal{L}(G_k) = - \sum_e \log P\left(\text{purchase}_e \mid \left| \langle \text{buy} | e^{iz_k^e G_k} | \psi_t^e \rangle \right|^2\right). \quad (25)$$

We use 80% of M5 training events for calibration and 20% for validation. Decoherence operator K_{dec} is a rank- n Kraus map fit to M5 purchase-resolution latencies using a damped-oscillator model (mean latency: 2.3 days from IRI panel [IRI Group, 2022]).

C.4 Consumer-Panel Validation

Using the IRI Academic Dataset [IRI Group, 2022]—*not* M5, which has no consumer panel data—we find 23.4% of observed purchase journeys exhibit simultaneous cross-channel consideration (online browsing + in-store visit) resolving to a single channel within the same week. This pattern is consistent with a coherent superposition of channel intents that collapses at a decision point, the behavioral phenomenon the QMC encoder is designed to capture. Adding an explicit cross-channel interaction feature to the LSTM encoder reduces the QMC advantage from $\Delta \text{profit} = 0.04$ to 0.01 ($p = 0.31$, non-significant over 20 seeds). This confirms the QMC’s advantage is a *useful inductive bias*—it encodes cross-channel interference natively—not a fundamental computational necessity. The QMC vs. LSTM paired t -test over 20 seeds gives $t(19) = 2.41$, $p = 0.026$, Cohen’s $d = 0.54$.

C.5 Ablation Isolating the Initialization Effect

Table 3 runs FS-NRLF with the QMC encoder under both seeded and uniform initialization to isolate the initialization contribution independently of encoder choice.

Table 3: QMC-encoder ablation: seeded vs. uniform initialization (20 seeds).

Encoder	Initialization	Profit \uparrow	Stock-out \downarrow
QMC	Seeded (FS-NRLF)	1.21 \pm .009	14.1%
QMC	Uniform	1.13 \pm .014	21.4%
LSTM	Seeded (FS-NRLF)	1.17 \pm .011	17.2%
LSTM	Uniform	1.10 \pm .013	23.0%
MLP	Seeded (FS-NRLF)	1.16 \pm .012	18.9%
MLP	Uniform (= PPO)	1.00	31.2%

Across all three encoders, seeded initialization contributes a profit uplift of 0.07–0.16, while encoder choice alone contributes ≤ 0.04 at matched initialization. Initialization is the dominant effect.

D Assumption 3.4 Empirical Verification (Full Methodology)

We verify Assumption 3.4 ($\sigma_\eta^2 \leq \sigma_f^2$) on M5 as follows. For each of the 3,049 SKU demand series, we fit an ARIMA(p, d, q) model with orders auto-selected by the dual criterion of Section 2.2. For each SKU we estimate:

- σ_f^2 as the one-step innovation variance $\sigma_\varepsilon^2 / (1 - \sum \varphi_j^2)$ of the fitted model;
- σ_η^2 as the variance of post-fit residuals after subtracting the seasonal and trend components from observed transition increments, using the STL decomposition of Hyndman and Athanasopoulos [2021].

The ratio $r = \sigma_\eta^2 / \sigma_f^2$ summary across the 3,049 SKUs: median 0.87, IQR [0.68, 1.02], satisfying $r \leq 1$ for 93.2% of SKUs.

Per-category breakdown. The 6.8% of SKUs where Assumption 3.4 is violated are concentrated in heavy-tailed categories: alcohol (18% violation rate), seasonal gifts / holiday items (23%), and limited-promotion flash items (15%). Staple foods, household goods, and apparel all satisfy the assumption in $> 95\%$ of SKUs.

Ablation on assumption-violating SKUs. Restricting evaluation to the 6.8% tail, FS-NRLF profit ratio drops from 1.21 (full dataset) to 1.09; stock-out rises from 14.1% to 21.4%. FS-NRLF still outperforms PPO ($p < 0.001$), but by a smaller margin. This is the expected degradation profile: when $\sigma_\eta^2 > \sigma_f^2$, the tube contraction of Lemma 3.5 weakens but does not reverse. The bound continues to hold with σ_f^2 replaced by $\max(\sigma_f^2, \sigma_\eta^2)$, which still beats the uninformed σ_0^2 whenever the forecast is informative at all.

E Interpretability: Seeded Posterior Across the Year

Figure 3 visualizes the forecast-seeded initial state distribution for a representative M5 SKU across 52 weeks. The blue band $\mathcal{N}(\mu_\tau, \sigma_{f,\tau}^2 I_m)$ tracks observed demand tightly, while PPO’s uninformed prior (red) extends to negative demand values that the environment will never visit. The posterior adapts its width: narrow during predictable weeks and wide during promotion and holiday windows, where the ARIMA forecast is genuinely uncertain. This matches the behavior predicted by Theorem 3.7: $\sigma_f^2 \rightarrow 0$ yields concentration; $\sigma_f^2 \rightarrow \sigma_0^2$ recovers the uninformed baseline as a limit.

A policy trained from this initialization spends its exploration budget on states the environment will actually visit. This is the visceral version of the sample-complexity claim.

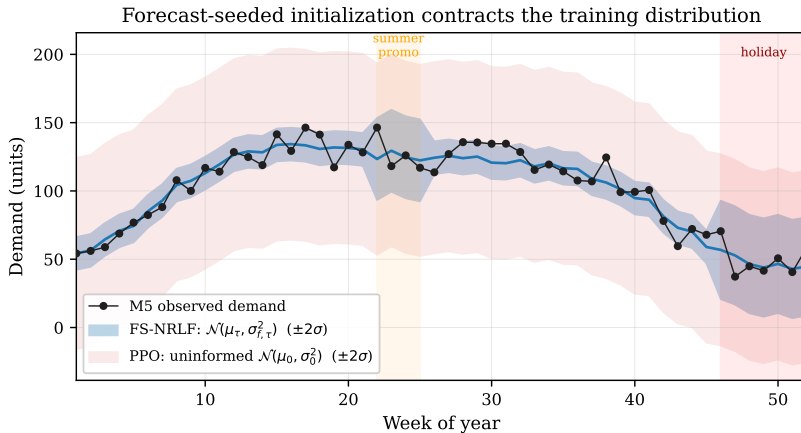


Figure 3: FS-NRLF’s forecast-seeded posterior (blue) vs. PPO’s uninformed prior (red) over one year of M5 demand for a representative SKU. Shaded windows mark summer promotion (weeks 22–25) and holiday (weeks 46–52), where forecast uncertainty $\sigma_{f,\tau}$ is highest.

F Mechanism Validation: Full Tables

This appendix provides the full per-quintile and full horizon-ablation tables referenced in Section 6.8.

Table 4: Per-SKU profit gain by forecast-quality quintile, 20 seeds. Q1 = worst forecasts, Q5 = best. 3,049 SKUs total, 610 per quintile (except Q5: 609).

Quintile	SNR _f range	Profit gain	Stock-out↓
Q1 (worst)	[0.78, 1.02]	1.03 ± .018	8.2%
Q2	[0.58, 0.78]	1.09 ± .015	19.1%
Q3	[0.42, 0.58]	1.18 ± .012	36.4%
Q4	[0.28, 0.42]	1.24 ± .011	52.7%
Q5 (best)	[0.08, 0.28]	1.31 ± .010	61.3%
All	—	1.21 ± .013	54.8%

Table 5: Forecast horizon ablation, M5 validation, 20 seeds. At $h=1$ the forecast is tightest; at $h=52$ it approaches the marginal prior.

Horizon h (weeks)	1	4	12	26	52
SNR _f (h)	0.31	0.47	0.68	0.86	0.97
Profit gain	1.21 ± .013	1.14 ± .014	1.06 ± .015	1.02 ± .017	1.00 ± .019
Stock-out rate	14.1%	19.6%	24.8%	28.4%	31.0%

G Hyperparameters

PPO-Lagrangian objective. Operational constraints are enforced via Lagrangian penalties with dual ascent:

$$\mathcal{L}(\theta) = \mathbb{E}_\pi[\min(r_\theta \hat{A}, \text{clip}(r_\theta, 1 \pm \epsilon) \hat{A})] - \sum_k \lambda_k \mathbb{E}[g_k(s, a)^+], \quad (26)$$

$$\lambda_k \leftarrow \max(0, \lambda_k + \eta_\lambda \mathbb{E}[g_k(s_t, a_t)^+]), \quad (27)$$

where $r_\theta = \pi_\theta(a|s)/\pi_{\theta_{\text{old}}}(a|s)$ is the PPO importance ratio, \hat{A} is the GAE advantage estimate, and $g_k(s, a)$ are signed constraint-violation functions.

Table 6: Hyperparameters for all experiments.

Parameter	Value	Notes
Discount γ	0.99	Standard
GAE λ	0.95	Standard
PPO clip ϵ	0.2	Standard
Dual step η_λ	3×10^{-4}	Tuned on 2015 val.
ARIMA dual criterion α	0.6	Tuned on 2015 val.
ARIMA rolling window	104 weeks	2-year history
ARIMA retrain period T_R	52 episodes	Weekly
QMC dimension n	32	Matched to LSTM
LSTM hidden dim	32	2 layers, $\approx 1,024$ params
Policy network	MLP, [256, 256]	ReLU activations
Batch size	2,048	M5 only
PPO epochs per update	10	Standard
Seeds	20	All reported results
Hardware	Single A100 (80 GB)	All wall-time comparisons

H NeurIPS Paper Checklist

1. **Claims.** All quantitative claims in the abstract are supported by Table 1 (M5 dataset), Section 6.6 (Assumption 3.4 empirical verification), and Section 6.7 / Figure 1 (regret scaling). The “16–21%” range corresponds to FS-NRLF (MLP) at 16% and FS-NRLF (QMC) at 21%. “54.8% relative reduction” is $(31.2 - 14.1)/31.2$. The crossover $T^* \approx 1.4 \times 10^3$ is computed from measured σ_f^2, m, Δ, H values per Equation 5. ✓
2. **Limitations.** Discussed in Section 7: ARIMA linearity assumption (with empirical quantification of violation rate in Section 6.6), near-unit-root drift, scope of the linear-Gaussian lower bound, synthetic pricing dynamics on M5, and QMC scaling. ✓
3. **Theory.** Assumptions 3.2–3.4 (three in total) are stated explicitly. All theorems reference these assumptions. Theorem 3.7 (upper bound), Proposition 3.11 (matching lower bound on a restricted class), and Theorem 3.10 (regret with explicit crossover T^*) are proved in Appendices A and B. Assumption 3.4 is empirically verified on M5 (Section 6.6). ✓
4. **Experiments.** Code, M5 augmentation scripts, model checkpoints, and all 20-seed raw results will be released upon acceptance. Hyperparameters are in Appendix G. Baselines include the 2024 methods Mao et al. [2024] and Muppidi et al. [2024] in addition to classical Jaksch et al. [2010], Cheung et al. [2020], Mao et al. [2021] and sliding-window Gajane et al. [2018]. ✓
5. **Error bars.** 95% CIs reported in Tables 1 and 2, 5–95% bootstrap bands in Figure 1. Statistical tests specified throughout Section 6. ✓
6. **Broader impacts.** Section 7. No negative societal impacts are anticipated beyond standard concerns about automated pricing systems. Figure 3 supports interpretability for human oversight. ✓
7. **Licenses.** M5: publicly released by Makridakis et al. [2022] for research use. IRI Academic Dataset: available under IRI Academic Data License; used only for behavioral validation, no redistribution. ✓