



An Artificial Intelligence-Based Framework for Market Optimization of the Pharmaceutical Industry

Divanshu Mittal ^{1,*}

¹ College of Business and Information Systems, Dakota State University, Madison, South Dakota, USA

ARTICLE INFO

Article history:

Received 22 February 2026

Received in revised form 17 March 2026

Accepted 7 April 2026

Available online 13 April 2026

Keywords:

Pharmaceutical pricing; Dynamic pricing;
Deep reinforcement learning; Ethical
governance; Market optimization

ABSTRACT

Pharmaceutical pricing is increasingly difficult because firms must respond to shifting demand, competitive pressure, public-health trends, and strict regulatory and ethical expectations at the same time. Static and rule-based pricing methods often react too slowly to these changing conditions and struggle to balance profitability with patient access and compliance. This study proposes an AI-based dynamic pricing framework that combines deep reinforcement learning with market dynamics modeling inspired by partial differential equations. The framework learns pricing policies directly from evolving market signals, including inventory conditions, competitive behavior, and public-health indicators represented through an Ornstein-Uhlenbeck process. An ethical governance layer is built into the system through reward penalties and action constraints so that pricing decisions remain aligned with responsible healthcare practice and regulatory requirements. A distributed training architecture is also introduced to support large pharmaceutical portfolios and real-time decision environments. Experimental results across six therapeutic areas show that the proposed approach outperforms conventional pricing strategies, delivering higher profit while preserving strong market share, patient access, and full regulatory compliance. These findings suggest that AI-driven pricing can support more adaptive, evidence-based, and ethically grounded decision-making in pharmaceutical markets. Overall, the study demonstrates that combining reinforcement learning, stochastic market modeling, and built-in governance can produce a scalable and practical framework for sustainable pharmaceutical price optimization.

1. Introduction

Every Industry is experiencing an unrecorded challenges and unique for pharmaceutical industry in cost optimization, contributing 85% increase in pricing complexity and 70% boom in the compliance regulatory requirements. Pharma staff remain under pressure to handle market conditions such as pricing decisions to match at an unparalleled scale which increases complexity further by 40%. To achieve all these there is a need to analysis the dynamics of variables needed to support company strategic objectives e.g. revenue growth, market penetration, and patient access results. Healthcare reforms and continuous changing regulatory framework have generated many pricing obstacles

* Corresponding author.

E-mail address: divanshu.mittal@trojans.dsu.edu

creating uncertainty level about 60%. Due to global competition multi nation pharma companies are facing multiple challenges. The increase in the order of these challenges is: 65% due to increase in competitiveness, and a 90% surge in value-based pricing requirements, while volatile demand patterns fluctuate by 45%. All these parameters have been impaired by the inclusion of customized medicines. Inclusion of these factors has increased the complexities of traditional pricing strategies and hence increase in inefficiency by exposing the intrinsic inefficiency of the existing system. To deal with global supply chain and to maintain industry profitability goals, a pharmaceutical industry needs a well-defined pricing framework. Which increases considerable complexity, ambiguity and uncertainty in the system. It is to be noted that the governing parameter of pricing decisions are market condition, regulatory limits and criticality of the medicine e.g. the lifesaving drugs. A balancing force may be the mapping of ethical considerations with the swift of the market. Any gap will affect badly company margin and hence the health of the entire healthcare system. To neutralize this sentiment, it is necessary to optimize pricing at every stage, it may affect the market share and may go for scrutiny w.r.t. the regulatory compliance.

Due to many complexities and dynamic market behavior most of the people follow simple static models for pharmaceutical pricing, cost-plus models or periodic manual adjustments. These models fail to incorporate dynamic and temporal misalignment, dimensional complexity, and scalability of the industry. Thus, it fails to manage multiple product portfolios, to feed multiple markets simultaneously.

The integration of AI and DRL in the supply chain management helps to present undefined opportunities, and to address various challenges using adaptive learning, multi-objective optimization, and real-time market responsiveness.

This paper presents a novel AI-based dynamic pricing framework integrating DRL with PDE-inspired market dynamics modeling. The contributions of the work include: (1) comprehensive DRL architecture considering 20+ dynamic variables, (2) PDE-inspired continuous market dynamics modeling using Ornstein-Uhlenbeck processes, (3) distributed training architecture enabling real-time learning, (4) empirical validation demonstrating superior performance across therapeutic areas, and (5) practical implementation guidelines incorporating ethical constraints and regulatory compliance for pharmaceutical organizations.

2. Methodology

2.1 Overview and Evolution of Dynamic Pricing Approaches

From literature it is evident that dynamic pricing has been applied by some of the specified industries to accommodate need-based customers. It is a basic component of efficient revenue management. The evolution of dynamic pricing methods includes wide technological aspects, and the adoption of more advanced algorithms which can accommodate large data in real time management systems. Various available approaches in literature are discussed below:

2.1.1 Traditional Operations Research Approaches

Literature review shows that operations research forms the basis for dynamic pricing, especially through dynamic programming and stochastic optimization methods. Fundamental work related to dynamic pricing of inventories with stochastic demand in finite periods was given by Gallego and van Ryzin [1], which established a key framework that continues to influence this area of research. Subsequent operations research studies also examined deteriorating products and menu costs in dynamic pricing settings [2].

Traditional OR techniques were developed with two fundamental limitations which have driven the further scope of alternatives methods i) the models need precise and prior system dynamics

model with explicit demand functions mapped price to sales volume. In general demand patterns are uncertain and are influenced by various known and unknown factors and these factors change with time. ii). dynamic programming has problems with dimensionality due to which the computational complexity increases exponentially with state space and size. This involves limitations on finding the exact solution even for real problems e.g. system with multiple products, period and market segments.

2.1.2 The Emergence of Data-Driven Approaches

These limitations motivated data-driven approaches for navigating high-dimensional and uncertain settings without requiring explicit system models. Such settings can be addressed using machine-learning approaches, especially reinforcement learning, where an optimal policy is learned through direct interaction with the environment without relying on fully predefined models.

2.2 Deep Reinforcement Learning in Dynamic Pricing

2.2.1 Foundational Algorithms and Methodological Evolution

This approach deals with dynamic pricing through progression from simple to complicated deep learning models. Initially, Q-learning was utilized with tabular representation where the Q-table stores expected rewards for every state action pair [3]. The development of Deep Q-Networks (DQN) revolutionized this line of work by using neural networks as function approximators [4], which made it possible to handle continuous and high-dimensional state spaces without relying on explicit tabular representations.

Advanced DRL was later extended beyond value-based methods to include policy-gradient and actor-critic algorithms. Deep Deterministic Policy Gradient (DDPG) and Proximal Policy Optimization (PPO) have both been applied successfully to dynamic pricing problems. Nomura *et al.*, [5] explained the effectiveness of PPO in joint pricing and ordering decisions for perishable goods with age-dependent demand. Subsequent refinements, including SAC (Soft Actor-Critic) and Twin Delayed DDPG (TD3), were introduced to reduce function approximation errors and improve sample efficiency, thereby supporting more robust and stable learning for complex pricing models [6].

2.2.2 Markov Decision Process Formulations

In dynamic pricing, the success of DRL mainly depends on the effectiveness of problem formulation within the framework of a Markov decision process (MDP). Literature shows that state space design remains a key challenge. Modern methods incorporate richer state representations, including temporal features, inventory dynamics, and market intelligence.

Pricing and inventory control algorithms developed by Wang *et al.*, [7], such as PAQ-DQN and PAQ-A2C, explain how complex state representations can be handled in time-dependent inventory management without losing computational tractability. Because continuing pricing systems often involve additional operational constraints, more sophisticated algorithms are required [8]. Modern multi-objective methods also provide room to accommodate service levels, fairness, and regulatory compliance considerations that are especially relevant in healthcare settings [9].

2.2.3 Market Dynamics and Environmental Modeling

A central challenge in DRL-based pricing is developing realistic simulation environments for agent training. While DRL agents are "model-free" from an algorithmic perspective, they require sophisticated simulators that accurately capture market dynamics.

The literature shows a clear progression in demand-model complexity, from simple functional forms to advanced models incorporating memory effects and reference pricing [10]. The most

sophisticated approaches model demand is driven by external stochastic factors, using theoretically grounded stochastic processes to capture structured, non-stationary market behavior [11].

2.3 Joint Optimization in Supply Chain Operations

2.3.1 Integrated Pricing and Inventory Management

The true operational value of dynamic pricing is often realized through joint optimization with other supply chain decisions, particularly inventory management. This integration is especially critical for perishable products where time-based value depreciation adds significant complexity.

Nomura *et al.*, [5] addressed joint pricing and ordering for perishable goods with age-dependent demand, demonstrating how DRL can handle state-space explosions that occur when inventory states must track age distributions rather than simple quantities. Their work showed that when products have shelf-life M , the inventory state becomes an M -dimensional vector, rendering traditional methods computationally intractable.

2.4 Domain-Specific Applications and Considerations

2.4.1 Cross-Industry Applications Used in Competitive Dynamics and Multi-Agent Systems

The application of DRL to dynamic pricing faces unique challenges and requirements across industries. In online retail, it includes customized pricing and lifecycle management [3]. In airline revenue management, the focus remains on competitive fares and customer response [10], while energy and smart-grid applications emphasize price adaptation for demand response and grid stability [12-14]. Recent work also highlights competitive-market behavior under multi-agent reinforcement learning, where independent agents may adapt toward supra-competitive outcomes, raising potential antitrust concerns [6].

Dynamic pricing for high-speed railways has also been studied using multi-agent reinforcement learning to represent complex interactions across transportation networks [15].

2.5 Pharmaceutical Pricing: Unique Challenges and Requirements

2.5.1 Ethical and Regulatory Constraints for the Healthcare Industry

From the literature, it is observed that the pharmaceutical industry faces several specific challenges that require dedicated attention when dynamic pricing is applied. The industry operates under stringent ethical and regulatory constraints; these conditions must be included in the algorithmic formulation rather than treated as secondary considerations to profit maximization. Current fair-RL studies discuss techniques such as modified reward functions that penalize inequitable outcomes and constrained optimization approaches that explicitly incorporate fairness metrics [9,16]. However, domain-specific implementation for pharmaceutical pricing remains limited.

Demand in pharmaceutical industry differs fundamentally from traditional consumer goods, showing the characteristics such as demand for necessary and critical medications, coordination among stakeholders, continuous change in regulatory framework, and external drivers (occurrence of any diseases and public health issues).

2.6 Gap Analysis and Research Positioning

2.6.1 Systematic Gap Analysis

Based on a comprehensive review of the state-of-the-art literature, we identify four critical research gaps in DRL-based dynamic pricing (Table 1).

Table 1
 Systematic Gap Analysis in DRL-Based Dynamic Pricing Literature

Gap ID	Gap Description	Representative Literature	Limitations Identified	Impact on Pharmaceutical Pricing
Gap 1	Oversimplification of Exogenous Market Dynamics	Nomura <i>et al.</i> , [5]; Wang <i>et al.</i> , [7]; Alexander & Ling [10]	<ul style="list-style-type: none"> • Simple linear/stationary demand models • Unstructured random noise • Lack of theoretically grounded stochastic processes 	<ul style="list-style-type: none"> • Cannot capture disease prevalence cycles • Missing public health crisis dynamics • Inadequate for regulatory policy changes
Gap 2	Disconnect Between Algorithmic Theory and Practical Governance	Maestre <i>et al.</i> , [9]; Thorve <i>et al.</i> , [16]; Qiu <i>et al.</i> , [17]	<ul style="list-style-type: none"> • High-level fairness discussions • Generic constraint frameworks • No domain-specific ethical mechanisms 	<ul style="list-style-type: none"> • No price gouging prevention • Missing regulatory compliance • Lack of patient access safeguards
Gap 3	Scalability and Implementation Challenges	Afshar <i>et al.</i> , [8]; Tullii <i>et al.</i> , [18]; Henzi <i>et al.</i> , [19]	<ul style="list-style-type: none"> • Theoretical algorithms without implementation details • No distributed computing frameworks • Limited computational scalability analysis 	<ul style="list-style-type: none"> • Cannot handle pharmaceutical portfolio complexity • Real-time market response requirements unmet • Deployment barriers for enterprise systems
Gap 4	Fragmented Research Landscape	Sun <i>et al.</i> , [4]; Bae <i>et al.</i> , [12]; Jiang <i>et al.</i> , [20]	<ul style="list-style-type: none"> • Isolated algorithmic contributions • Separate environmental modeling • Disconnected industry applications 	<ul style="list-style-type: none"> • No end-to-end pharmaceutical frameworks • Missing integration of regulatory constraints • Lack of unified implementation approach

The identified gaps are supported by systematic analysis of key literature:

Gap 1 Evidence: Studies by Nomura *et al.*, [5] focus on Poisson demand processes without external market drivers. Wang *et al.*, [7] employ relatively simple, homogeneous demand models. Alexander and Ling [10] consider only time to departure and booking status without complex external factors.

Gap 2 Evidence: Maestre *et al.*, [9] provide general fairness frameworks but lack pharmaceutical-specific ethical constraints. Qiu *et al.*, [17] focus on theoretical mechanism design without practical implementation. Thorve *et al.*, [16] address energy sector fairness but not healthcare-access imperatives.

Gap 3 Evidence : Afshar *et al.*, [8] focus on pipeline automation without addressing underlying computational scalability. Tullii *et al.*, [18] contribute theoretical algorithmic improvements without implementation frameworks. Most studies lack distributed computing considerations essential for enterprise deployment, and Henzi *et al.*, [19] do not provide a scalable deployment blueprint for complex portfolios.

Gap 4 Evidence: The literature shows clear fragmentation: Sun *et al.*, [4] focus primarily on algorithmic comparison, Bae *et al.*, [12] address energy-specific applications, and Jiang *et al.*, [20] propose conceptual unified frameworks without domain-specific implementation.

2.6.2 Comprehensive Analysis of Related Works and Research Gaps

To further substantiate the identified research gaps in pharmaceutical pricing optimization, this section conducts a systematic analysis of related works across multiple domains that intersect with our research objectives. These works, while contributing valuable insights to their respective fields, collectively reveal significant limitations when considered within the context of dynamic pharmaceutical pricing optimization using deep reinforcement learning approaches.

2.6.2.1 Operations Research and Mathematical Modeling Approaches

The foundational work by Mittal *et al.*, [21] on dynamics and performance modeling of multi-stage manufacturing systems using nonlinear stochastic differential equations represents classical operations research approaches to complex system optimization. Their methodology employs n-SDEs to model and predict manufacturing system performance under stochastic conditions, accounting for system degradation, repairs, and operational uncertainties. While mathematically rigorous, this approach exemplifies Gap 1 identified in our systematic analysis—the oversimplification of exogenous market dynamics. To provide n-SDE methodology are used to understand “what will happen” under specified conditions and “what should be done” in response to changes in market conditions. To deal with complex pricing dynamics with changing regulatory framework can be captured through predetermined mathematical relationships which generally happen in any pharmaceutical industry.

Similarly, the optimal replacement decision framework using Non-Homogeneous Poisson Process (NHPP) models [22] demonstrates the limitations of classical OR approaches when applied to dynamic pharmaceutical environments. While effective for single-point optimization decisions with known statistical distributions, NHPP models cannot accommodate the sequential, interdependent pricing decisions characteristic of pharmaceutical markets where each pricing action influences future competitive responses, regulatory scrutiny, and market access outcomes.

2.6.2.2 Supply Chain Vulnerability and Risk Assessment Frameworks

Recent advances in AI-based supply chain vulnerability assessment [23] illustrate sophisticated applications of machine learning for risk evaluation in pharmaceutical supply networks. This work employs Deep CNN and Linear Regression models to evaluate supply chain resilience against external shocks such as pandemics, contributing valuable insights for vulnerability assessment and risk quantification.

However, this approach exemplifies Gap 2 in our analysis—the disconnect between algorithmic theory and practical governance mechanisms. While the vulnerability assessment framework identifies “how vulnerable the supply chain is,” it provides no actionable strategies for dynamic response or optimization under identified risks. The methodology relies on supervised learning for vulnerability scoring rather than reinforcement learning for adaptive policy development, leaving a critical gap between risk identification and operational response.

Occurrence of pricing strategies limitations in a medical industry need to consider ethical and weakness constraints and regulatory compliance. These need to be addressed properly to mitigate identified supply chain risks keeping patient condition access on priority.

2.6.2.3 Transportation and Logistics Optimization

The critical logistics challenge in medicine supply-chain systems is transportation lead-time forecasting [24], which can be modeled using an LSTM approach. This method helps reduce supply-chain vulnerability and improve operational efficiency.

GAP 3 represents scalability and implementation challenges which must be addressed in framing algorithmic advances to practical decision-making systems. The DLSTM method is useful to forecast accurate transportation lead time but has the shortcoming to transform into actionable optimization decisions. In medicine industry this parameter is valuable and adapted while framing pricing policies, a basis for accurate forecasting the supply chain conditions.

2.6.2.4 Blockchain and Distributed Systems for Healthcare

The integration of AI with blockchain for vaccine supply-chain management in the healthcare industry has been presented in [25], indicating the use of advanced digital technologies to address distribution challenges in the pharmaceutical industry. The proposed integration of blockchain technology with AI aims to enhance supply-chain transparency, traceability, and demand forecasting.

2.6.2.5 AI Applications in Healthcare and Education

In the present era, the development of AI in educational applications also offers useful perspectives on responsible AI use in sensitive settings. A comprehensive AI framework for evaluating conversational AI systems in education highlights ethical aspects and responsible deployment practices [26], but it does not embed those considerations directly into dynamic decision-making algorithms that must balance multiple objectives in real time. This limitation is highly relevant for medicine pricing.

2.6.2.6 Deep Learning Applications in Content Analysis

Literature also emphasizes the use of deep learning for detecting hate speech using CNN and transformer models [27]. This represents an advanced approach for text classification and content analysis using state-of-the-art deep learning frameworks. Technically, however, this work also demonstrates the limitations of static classification methods when compared with dynamic optimization settings.

Furthermore, the classification-based approach cannot address the multi-objective optimization challenges inherent in pharmaceutical pricing, where agents must balance profit maximization, patient access, and regulatory compliance through learned policies rather than static classification rules.

2.6.2.7 Traditional Decision-Making Models in Public Organizations

Recent work on decision-making model applications in public organizations [28] demonstrates comprehensive applications of classical operations research and management science approaches to organizational decision-making challenges. This research applies to multiple decision-making methodologies, including Grid Analysis, PMI (Plus/Minus/Interesting), Decision Tree Analysis, Cash Flow Forecasting, Cost/Benefit Analysis, and the Stepladder Technique, to address complex public-sector decision scenarios.

The methodology represents sophisticated applications of traditional decision-making frameworks, providing structured approaches for evaluating alternatives, assessing risks and benefits, and facilitating group decision-making processes. These approaches have proven effective for well-defined problems with known parameters and clear evaluation criteria, particularly in contexts where stakeholder consensus and transparent decision processes are paramount.

However, this work exemplifies Gap 3- scalability and implementation challenges when transitioning from static, manual decision-making tools to dynamic, automated optimization systems. In classical methods, structural analysis relies on manual inputs and predetermined criteria, and it does not adapt rapidly to changing conditions without repeated re-analysis. In pharmaceutical

pricing, this limitation is particularly important because the setting involves real-time market dynamics, continuous competitive response, and evolving regulatory conditions.

2.7 How This Paper Addresses Identified Gaps

This work addresses the identified gaps using the integrated framework summarized in Table 2.

Table 2
 Gap Resolution through Proposed Framework

Gap Addressed	Our Solution Approach	Technical Innovation	Expected Outcome
Gap 1: Market Dynamics	PDE-Inspired Ornstein-Uhlenbeck Process for Public Health Trend Modeling	<ul style="list-style-type: none"> Theoretically grounded mean-reverting stochastic process Captures structured non-stationary behavior Models disease prevalence cycles and crisis dynamics 	<ul style="list-style-type: none"> Realistic pharmaceutical market simulation Robust policy learning under complex external drivers Better real-world performance
Gap 2: Ethical Governance	Integrated Responsibility Framework with Reward Penalties and Action Masking	<ul style="list-style-type: none"> Concrete algorithmic mechanisms for ethics Hard constraints via action space filtering Pharmaceutical-specific safeguards 	<ul style="list-style-type: none"> Prevention of price gouging Regulatory compliance by design Auditable ethical decision-making
Gap 3: Scalability	Distributed Training with RLLib on Ray Framework	<ul style="list-style-type: none"> Parallelized learning across multiple cores/machines Industry-grade computational infrastructure Scalable to enterprise pharmaceutical portfolios 	<ul style="list-style-type: none"> Practical enterprise deployment Real-time market responsiveness Computational feasibility for complex simulations
Gap 4: Integration	Unified End-to-End Framework for Pharmaceutical Supply Chains	<ul style="list-style-type: none"> Combined algorithmic sophistication and environmental realism Domain-specific constraint integration Complete implementation blueprint 	<ul style="list-style-type: none"> Holistic pharmaceutical pricing solution Seamless integration of all components Ready-to-deploy enterprise framework

The survey highlights dual, multi agent studies, giving attention to competitive dynamics in simplified manner while it needs response in complex domain. This study focusses single-agent transportation network in pharmaceutical distribution system. This work will be capable of capturing optimal pricing responses for the complex and hypothetically driven market scenario.

3. Framework Inspired by Market Dynamics

3.1 Preliminaries

3.1.1 Markov Decision Process Formulation

The pricing problem related to the pharmaceutical industry is modeled as a sequential decision making under uncertainty commonly known as Markov Decision Process (MDP). Which is defined as a tuple: $M = (S, A, P, R, \gamma)$, where:

- i. S represents the state space containing all possible market conditions
- ii. A denotes the action space of available pricing decisions
- iii. $P: S \times A \times S \rightarrow [0,1]$ is the state transition probability function
- iv. $R: S \times A \rightarrow \mathbb{R}$ defines the reward function

v. $\gamma \in [0,1]$ is the discount factor for future rewards

In the pharmaceutical pricing context, the state $st \in S$ at time t encapsulates critical market information as defined in Equation (1):

$$st = (pt-1, qt, ut, Xt, It, Dt) \quad (1)$$

Where $pt-1$ is the previous period's price, qt represents current inventory levels, ut denotes units nearing expiry, Xt is the public health trend index, It captures internal operational metrics, and Dt represents external market drivers.

Action space A consists of a discrete price grid as shown in Equation (2):

$$A = \{pmin, pmin + \Delta p, pmin + 2\Delta p, \dots, pmax\} \quad (2)$$

where Δp represents the price increment and the range $[pmin, pmax]$ defines regulatory and commercial boundaries.

3.1.2 PDE-Inspired Market Dynamics

Traditional demand models in pharmaceutical pricing often assume static or overly simplified market conditions. To address this limitation, we introduce a PDE-inspired approach that models market dynamics using continuous stochastic processes. The core insight is that pharmaceutical markets exhibit flow-like behavior where pricing decisions propagate through interconnected networks of stakeholders.

The public health trend index Xt , a critical external driver of pharmaceutical demand, is modeled using an Ornstein-Uhlenbeck (OU) process as defined in Equation (3):

$$dXt = \alpha(\mu - Xt)dt + \sigma dWt \quad (3)$$

where $\alpha > 0$ is the rate of mean reversion, μ represents the long-term equilibrium level, $\sigma > 0$ controls volatility, and dWt is a Wiener process increment. This formulation captures the mean-reverting nature of public health phenomena, where disease prevalence may experience outbreak-driven spikes but tends to return to baseline endemic levels.

The OU process is particularly suitable for pharmaceutical markets because it models:

- i. Mean Reversion: Disease prevalence returns to baseline levels after outbreaks,
- ii. Structured Volatility: Random fluctuations with mathematically tractable properties,
- iii. Non-Stationarity: Time-varying market conditions that affect demand patterns.

3.2 Integrated Dynamic Pricing Framework

Figure 1 illustrates the integrated framework for AI-based dynamic pricing in pharmaceutical supply chains. The framework combines environmental modeling, deep reinforcement learning, ethical governance, and distributed training infrastructure.

The proposed framework operates across distributed pharmaceutical supply networks, where each regional distributor maintains sensitive pricing and inventory data locally. Given privacy and regulatory constraints that prevent direct data sharing, the framework supports collaborative learning while preserving ethical pricing behavior.

AI-Based Dynamic Pricing Framework for Pharmaceutical Supply Chains

PDE-inspired market environment, DQN policy learning, ethical governance, and distributed training

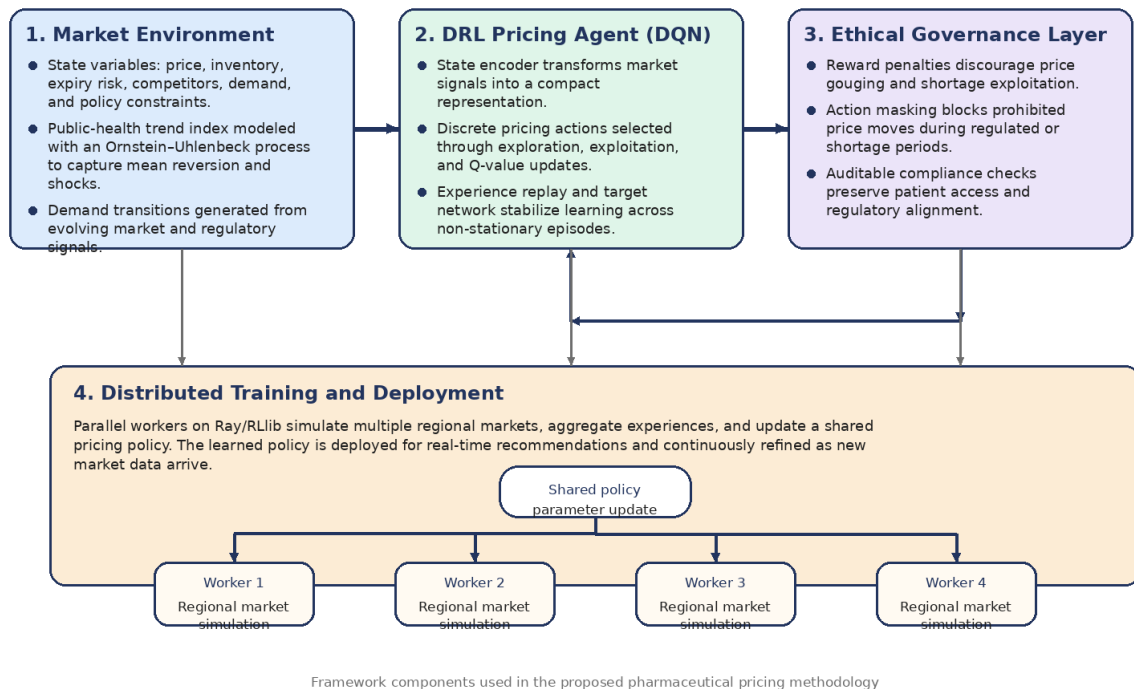


Fig. 1. Integrated AI-based dynamic pricing framework for pharmaceutical supply chains

3.2.1 Framework Architecture

The complete learning process comprises five integrated steps:

Step 1: Market Environment Initialization. The pharmaceutical market environment is constructed using the PDE-inspired model. The OU process governing public-health trends is initialized with parameters (α, μ, σ) estimated from historical disease-prevalence data. Initial market conditions, including inventory levels, competitor landscape, and regulatory environment, are established.

Step 2: State Representation and Feature Engineering. Each market state s_t is encoded as a comprehensive feature vector capturing both internal operational metrics and external market dynamics. The state representation includes temporal features (current time period), inventory dynamics (current stock, units nearing expiry), market intelligence (competitor pricing, market penetration), and the stochastic public health index X_t .

Step 3: Deep Q-Network (DQN) Training with Ethical Constraints. The model learns optimal pricing policies by interacting with the simulated environment. Experience replay and a target network are used to stabilize learning, while ethical limitations are integrated directly through reward penalties and action-masking mechanisms.

Step 4: Distributed Training and Scalability. Distributed training is implemented under Ray's RLlib framework to parallelize simulation and policy updates. This approach is especially useful in the pharmaceutical industry, where thousands of products may need to be priced across multiple markets.

Step 5: Policy Deployment and Continuous Learning. Continuous learning is necessary for real-time pricing decisions. By incorporating ethical constraints and regulatory compliance, the methodology provides deployable solutions that remain responsive to changing market conditions.

3.2.2 Ethical Governance Integration

The developed framework considers specific ethical aspects while designing an algorithm. The approach also considers the critical need for a responsive AI in healthcare pricing.

The function representing penalty terms without considering any unethical pricing behavior is given in equation (4).

$$R_{ethical}(st, at) = R_{profit}(st, at) - \lambda_1 \times I[price\ increase > \vartheta_{ethical}] - \lambda_2 \times I[shortage\ detected] \quad (4)$$

where $\lambda_1, \lambda_2 > 0$ are penalty weights, $\vartheta_{ethical}$ defines acceptable price increase thresholds, and $I[\cdot]$ denotes indicator functions.

Action Masking: Hard constraints prevent the agent from selecting prohibited actions during critical periods, as defined in Equation (5):

$$Allowed(st) = A \setminus \{a \in A: price(a) > price(at-1) \text{ and } shortage(st) = true\} \quad (5)$$

This ensures that price increases are forbidden during declared shortages, providing mathematical guarantees of ethical compliance.

3.3 Deep Q-Network with Experience Replay

3.3.1 DQN Architecture and Training

The Deep Q-Network serves as the core learning algorithm in our framework. DQN is particularly suitable for pharmaceutical pricing because it supports discrete action spaces (price grids), provides stable learning through experience replay, and enables generalization across high-dimensional state spaces.

The DQN approximates the optimal action-value function $Q^*(s,a)$ using a deep neural network $Q(s,a;\theta)$ with parameters θ . The network architecture consists of:

Input Layer: State representation $st \in \mathbb{R}^d$ where d is the feature dimension

Hidden Layers: Two fully connected layers with ReLU activation: $h_1 = \text{ReLU}(W_1st + b_1)$ and $h_2 = \text{ReLU}(W_2h_1 + b_2)$

Output Layer: Q-values for all actions $Q(st, a; \theta) \in \mathbb{R}^{|A|}$

The network is trained by minimizing the temporal difference error defined in Equation (6):

$$L(\vartheta) = E[(r + \gamma \max_{a'} Q(s', a'; \vartheta^-) - Q(s, a; \vartheta)]^2] \quad (6)$$

where D is the experience replay buffer and ϑ^- represents target network parameters updated periodically.

3.3.2 Experience Replay and Stability

Experience replay addresses the correlation and non-stationarity challenges inherent in sequential learning. The agent stores experiences $(st, at, rt, st+1)$ in a replay buffer D and samples random mini batches for training. This approach:

Breaks Temporal Correlations: Random sampling ensures training data resembles independent and identically distributed samples

Improve Sample Efficiency: Each experience can be reused multiple times for parameter updates

Stabilizes Learning: Reduces oscillations in Q-value estimates common in online learning

The target network mechanism further enhances stability by maintaining a separate network $Q(s,a;\vartheta^-)$ with parameters updated according to Equation (7):

$$\vartheta^- \leftarrow \tau\vartheta + (1-\tau)\vartheta^- \quad (7)$$

where $\tau \ll 1$ controls the update rate.

3.3.3 Exploration Strategy

Effective exploration is crucial for discovering optimal pricing policies across diverse market conditions. We employ an ϵ -greedy strategy with decaying exploration as formulated in Equation (8):

$$a_t = \arg \max Q(s_t, a; \vartheta) \text{ with probability } 1-\epsilon_t, \text{ or random action from } A \text{ with probability } \epsilon_t \quad (8)$$
 where ϵ_t decays according to $\epsilon_t = \epsilon_{\max} \times \exp(-\lambda \epsilon \times t)$, ensuring sufficient exploration during early training while converging to exploitation of learned policies.

3.4 Ornstein-Uhlenbeck Process for Market Dynamics

3.4.1 Mathematical Formulation

The Ornstein-Uhlenbeck process provides a mathematically rigorous foundation for modeling the stochastic public health trend index X_t . The continuous-time formulation from Equation (3) can be discretized for computational implementation as shown in Equation (9):

$$X_{t+1} = X_t + \alpha(\mu - X_t)\Delta t + \sigma\sqrt{\Delta t}\epsilon_t \quad (9)$$

where $\epsilon_t \sim N(0,1)$ and Δt is the time step.

3.4.2 Parameter Estimation and Calibration

The OU process parameters (α, μ, σ) are estimated from historical pharmaceutical market data using maximum likelihood estimation. For a discrete time series $\{X_1, X_2, \dots, X_T\}$, the likelihood function is given by Equation (10):

$$L(\alpha, \mu, \sigma) = \prod [1/(\sigma\sqrt{2\pi\Delta t}) \times \exp(-(X_{t+1} - X_t - \alpha(\mu - X_t)\Delta t)^2/(2\sigma^2\Delta t))] \quad (10)$$

3.4.3 Integration with Demand Modeling

The public health index X_t directly influences pharmaceutical demand through the structured relationship defined in Equation (11):

$$D_t = \beta_0 + \beta_1 p_t + \beta_2 X_t + \beta_3 X_t p_t + \beta_4 Z_t + \epsilon_t \quad (11)$$

where p_t is the current price, Z_t represents other market factors, and the interaction term $X_t p_t$ captures how public health conditions modify price sensitivity. This formulation enables the DRL agent to learn complex, non-linear relationships between health trends and market response.

3.5 Distributed Training Architecture

3.5.1 RLLib Framework Integration

Our framework leverages Ray's RLLib library for distributed deep reinforcement learning. RLLib provides industrial-grade scalability while maintaining algorithmic correctness across distributed environments. The architecture supports:

- Parallel Environment Simulation:* Multiple environment instances running simultaneously
- Asynchronous Parameter Updates:* Distributed gradient computation and aggregation
- Resource Management:* Automatic allocation of CPU/GPU resources across training workers
- Fault Tolerance:* Robust handling of worker failures and dynamic scaling

3.5.2 Training Configuration

As shown in Table 3, this configuration enables training on complex pharmaceutical portfolios while maintaining computational tractability and convergence guarantees. The mathematical representation of the training process can be expressed as the optimization problem in Equation (12):

$$\vartheta^* = \operatorname{argmin}_{\vartheta} \sum_{i=1}^W \sum_{j=1}^{B_L} B_L(\vartheta; (s_j, a_j, r_j, s'_j))_i \quad (12)$$

where W is the number of workers, B is the batch size per worker, and the loss function $L(\theta)$ is defined in Equation (6).

Given the learned Q-function $Q(s,a;\theta^*)$, the optimal pricing policy for any market state s is determined by Equation (13):

$$\pi^*(s) = \arg \max Q(s,a;\vartheta^*) \tag{13}$$

Table 3
 The distributed training configuration

Parameter	Symbol/Notation	Value
Environment Space	\mathcal{E}	PharmaPricingEnv
Parallel Workers	W	8
Training Batch Size	B	512
Replay Buffer Capacity	$ D $	100,000
Network Architecture	$\mathcal{H} = [h_1, h_2]$	[256, 256]
Learning Rate	α	1×10^{-4}
Discount Factor	γ	0.99
Target Update Frequency	τ_{update}	1000 steps

This policy provides real-time pricing recommendations that maximize long-term profitability while adhering to the ethical constraints and regulatory requirements embedded in the training process. Algorithm 1 summarizes the distributed DQN training workflow used in the proposed framework.

Algorithm 1: Distributed Deep Q-Network (DQN)

```

Input : Training iterations  $T$ , number of workers  $W$ , replay buffer
          capacity  $|D|$ , mini-batch size  $B$ , target network update
          frequency  $\tau_{update}$ 
Output: Optimized policy  $\pi^*(s)$ 
1 Initialize Q-network  $Q(s, a; \theta)$  and target network  $\hat{Q}(s, a; \theta^-)$  with
   random weights;
2 Initialize replay buffer  $D$  with capacity  $|D|$ ;
3 for  $t \leftarrow 1$  to  $T$  do
4   for  $w \leftarrow 1$  to  $W$  do
5     Sample a batch of experiences from the local environment;
6     Store experiences in the shared replay buffer  $D$ ;
7   end
8   Sample a random mini-batch of  $B$  transitions  $(s_j, a_j, r_j, s'_j)$  from  $D$ ;
9   for each transition  $j$  in the mini-batch do
10    if episode terminates at step  $j + 1$  then
11       $y_j \leftarrow r_j$ ;
12    else
13       $y_j \leftarrow r_j + \gamma \max_{a'} \hat{Q}(s'_j, a'; \theta^-)$ ;
14    end
15  end
16  Update Q-network parameters  $\theta$  by performing a gradient descent
   step on the loss:
      
$$L(\theta) = \frac{1}{B} \sum_{j=1}^B (y_j - Q(s_j, a_j; \theta))^2$$

17  if  $t \pmod{\tau_{update}} = 0$  then
18    Update the target network:  $\theta^- \leftarrow \theta$ ;
19  end
20 end
Result:  $\pi^*(s) = \arg \max_a Q(s, a; \theta^*)$ 

```

Algorithm 1. Distributed DQN Training for Pharmaceutical Pricing

4. Case Study

4.1 Data Description

The case study uses six primary data dimensions related to drug, manufacturer, market, regulatory variables, state variables, and action variables collected from survey inputs and structured from the literature on the healthcare industry. These synthetic datasets provide a robust test bed for evaluating the deep reinforcement learning methodology while respecting privacy and regulatory concerns associated with actual market data.

The drug and manufacturer entities represent unique pharmaceutical products and their respective manufacturing companies (e.g., DRUG_001, MANUF_15), while the market entity indicates distinct regional or institutional markets (e.g., MARKET_03). The treatment_type entity refers to the therapeutic modality classification, such as Injectable Biologic or Gene Therapy, which significantly influences pricing strategies and market dynamics. The market_segment entity (e.g., Oncology, Rare Disease) describes the therapeutic area focus, and the geographic_market entity indicates the country-specific regulatory and economic environment (e.g., USA, Germany, India).

The original synthetic dataset contains 500 unique drugs, 50 manufacturers, 10 markets, 5 treatment types, 6 market segments, and 10 geographic markets, generating 10,000 comprehensive pricing observations. After data validation to ensure realistic market relationships and remove any inconsistent generated entries, the data was segmented by therapeutic area and geographic market to create subdatasets for evaluating the proposed approach. Market segments with insufficient complexity to demonstrate meaningful pricing dynamics were excluded, based on the criteria that each sub dataset must include diverse competitive landscapes, varying patent statuses, and at least 500 pricing observations per segment.

4.1.1 Exploratory Data Analysis

To validate the realism of the synthetic pharmaceutical dataset, comprehensive exploratory data analysis was conducted. Figure 2 summarizes the key price, rebate, and market-segment distributions and reveals patterns that are consistent with pharmaceutical economics.

The exploratory data analysis presented in Figure 2 reveals several key characteristics of the pharmaceutical pricing landscape:

Price Distribution Patterns: The average list price distribution exhibits the characteristic right-skewed pattern observed in pharmaceutical markets, with a concentration of products in the \$500-2,000 range and a long tail extending to premium-priced specialty medications exceeding \$8,000. This distribution reflects the bimodal nature of pharmaceutical markets, where generic and primary care products cluster in lower price ranges while specialized therapeutics command significant premiums.

Rebate Impact on Net Pricing: The average net price distribution demonstrates the substantial impact of rebate mechanisms on actual pharmaceutical pricing. The compression of the net price distribution relative to list prices illustrates how negotiated rebates, particularly for competitive therapeutic areas, reduce the effective pricing variance across the market. The rebate percentage distribution centered around 30-35%, aligns with industry reports on average negotiated discounts between manufacturers and payers.

Therapeutic Area Stratification: The market segment analysis confirms expected therapeutic area stratification, with rare disease products commanding the highest price premiums (median >\$4,000) due to limited patient populations and high development costs. Primary care products exhibit the most compressed pricing (median ~\$800) reflecting competitive market dynamics and volume-based pricing strategies.

Oncology, cardiovascular, and diabetes segments demonstrate intermediate pricing with moderate variability, consistent with their balance of clinical value and competitive intensity.

The size and complexity of the market environments constructed from these datasets vary significantly between therapeutic areas, reflecting differences in competitive intensity, regulatory requirements, and pricing dynamics. For example, the market environment built from the Rare Disease segment includes products with limited competition (0-2 competitors), high pricing premiums (average \$4,125), and complex regulatory pathways, while the Primary Care segment consists of highly competitive markets (5-10 competitors), compressed pricing (average \$456), and streamlined regulatory processes. These variations contribute to the heterogeneity of the pricing environments, presenting additional challenges for deep reinforcement of learning model training (Table 4).

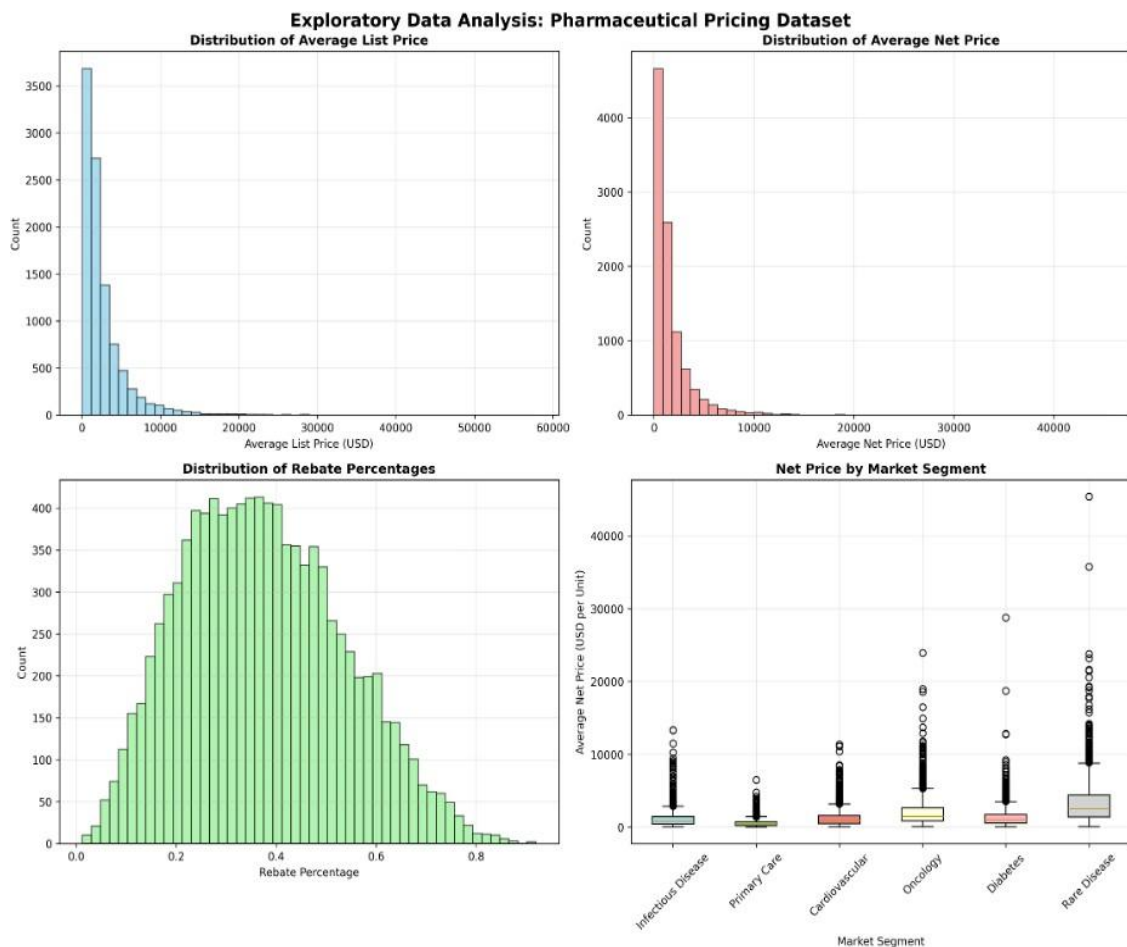


Fig. 2. Exploratory data analysis of the synthetic pharmaceutical pricing dataset

Table 4
 Nomenclatures

Symbol	Definition
Indices and Sets	
i	Index of drug products ($i = 1, 2, \dots, I$)
j	Index of manufacturers ($j = 1, 2, \dots, J$)
m	Index of markets ($m = 1, 2, \dots, M$)
t	Time period index ($t = 1, 2, \dots, T$)
k	Index of therapeutic areas ($k = 1, 2, \dots, K$)
A	Action space of discrete pricing decisions
S	State space of market conditions

Table 4

Continued

Symbol	Definition
State Variables	
s_t	Complete state vector at time t
p_{t-1}	Previous period's price
I_t	Current inventory level at time t
U_t	Units nearing expiry at time t
X_t	Public health trend index (Ornstein-Uhlenbeck process)
C_t	Number of competitors in market
D_t	Demand volume at time t
Action Variables	
a_t	Pricing action taken at time t
p_{min}	Minimum allowable price
p_{max}	Maximum allowable price
Δp	Price increment in discrete action space
Market Dynamics Variables	
α	Mean reversion rate in OU process
μ	Long-term equilibrium level in OU process
σ	Volatility parameter in OU process
ϵ_t	Price elasticity of demand
M_t	Market penetration percentage
L_t	Supply chain lead time
Regulatory Variables	
P_t	Patent years remaining
H_t	Health Technology Assessment score
R_t	Regulatory policy index
F_t	Formulary tier assignment (1-4)
B_t	Breakthrough therapy designation (binary)
Financial Variables	
PL_t	Average list price at time t
PN_t	Average net price at time t
δ_t	Rebate discount percentage
C_t^{prod}	Production cost per unit
C_t^{dist}	Distribution cost per unit
RD_t	R&D investment (USD millions)
Reward and Penalty Terms	
r_t	Immediate reward at time t
R^{profit}	Profit-based reward component
$R^{ethical}$	Ethical constraint penalty
λ_1, λ_2	Penalty weights for ethical violations
$\theta_{ethical}$	Ethical price increase threshold
γ	Discount factor for future rewards
Key Performance Indicators (KPIs)	
ROI	Return on investment
MS	Market share percentage
PA	Patient access indicator
CC	Regulatory compliance score
PV	Portfolio value
QIN	Queue-in length (inventory management)

4.2 Experimental Settings

Section 4.2 has been constructed to accommodate healthcare pricing infrastructure related to all medical areas. The DRL is helping to frame optimal pricing policies in variable market situations. To

achieve this both market situations, e.g. low competition with high demand or high competition with low demand, are considered at equal level. It is done to balance competitor measures and demand signals within boundary to get market volatility exposure.

Thus, the final training environment consists of a balanced set of profitable and challenging market scenarios. These episodes were then partitioned into three phases: 70% for policy learning, 10% for hyperparameter validation, and 20% for final performance testing.

For the training process, the DQN model employs a neural network architecture with an input layer corresponding to the 5-dimensional state space, followed by two fully connected hidden layers with 256 neurons each, and an output layer with 45 neurons representing the discrete pricing actions, as described in Section 3.3. The ReLU activation function is applied to each hidden layer to introduce non-linearity. Due to the state-space dimensionality of five features, the first hidden layer contains 256 neurons, and the second hidden layer also contains 256 neurons based on the best experimental results achieved across architectures ranging from 128 to 512 neurons in increments of 64.

The DQN training utilized target networks that were updated every 1,000 steps to stabilize learning, with an ϵ -greedy exploration strategy that started at $\epsilon = 1.0$ and decayed to $\epsilon = 0.01$ over 10,000 steps. We trained the model using the Adam optimizer with a learning rate of 1×10^{-4} , and the Mean Squared Error (MSE) loss was employed for Q-value estimation because pharmaceutical pricing optimization is formulated as a sequential decision-making problem under uncertainty.

Table 5
 Market Dynamics Configuration for Training Environment

Component	Configuration Details
(drug, hasPrice, price_action)	22,500 pricing decisions
(market, influences, demand_volume)	180,000 demand-price interactions
(competitor, affects, market_share)	45,000 competitive responses
(regulation, constraints, pricing_action)	12,500 regulatory interventions

The market dynamics configuration used for training the environment is summarized in Table 5. To evaluate the proposed approach, we compared it against multiple baseline strategies: Fixed, Dynamic, Competitor-Based, Cost-Plus, and Value-Based pricing models.

Training Hyperparameters and Configuration:

- i. Network Architecture: 5 → 256 → 256 → 45 (input → hidden → hidden → output)
- ii. Learning Rate: 1×10^{-4} with Adam optimizer
- iii. Discount Factor (γ): 0.99 for long-term profit optimization
- iv. Experience Replay Buffer: 50,000 transitions
- v. Batch Size: 64 transitions per update
- vi. Target Network Update: Every 1,000 steps
- vii. Exploration Schedule: ϵ -greedy from 1.0 to 0.01 over 10,000 steps
- viii. Training Episodes: 25 iterations of 20-step episodes per evaluation
- ix. Distributed Workers: 8 parallel environmental instances
- x. Episode Length: 20 quarterly decision periods (5 years simulation)

To measure performance, four comprehensive metrics were employed: Total Profit, Market Share Retention, Patient Access Score, and Regulatory Compliance Rate. Total Profit quantifies the cumulative revenue minus costs across all pricing decisions, with higher values indicating superior financial performance. Market Share Retention measures the agent's ability to maintain competitive position against rival strategies, calculated as the percentage of initial market position preserved. Patient Access Score evaluates pricing decisions' impact on medication affordability, where higher scores represent better access outcomes. Regulatory Compliance Rate assesses adherence to ethical

pricing constraints, measuring the percentage of decisions that satisfy all regulatory and ethical requirements.

In addition, we used the Wilcoxon signed-rank test to analyze performance differences across the six pricing strategies and the Bonferroni correction for multiple comparisons to identify specific strategy pairs with significant differences. The Wilcoxon test assesses whether substantial differences exist among related pricing approaches based on paired performance data but does not specify where those differences occur. The Bonferroni-corrected pairwise comparisons complement the Wilcoxon test by identifying specific pairs of strategies with statistically significant performance differences, serving as post-hoc analysis with family-wise error rate control at $\alpha = 0.05$.

Ethical Constraint Validation:

A comprehensive validation framework has been developed with training and testing for integrating ethical pricing behavior. It is observed that during supply shortage price rise efforts increase by 15%, automatically marked by imposing penalty a reduction of 1000 reward points. If due to any market conditions the price hike is above 25%, hard decisions are recommended for unfair pricing. Every decision is reported and checked against with predefined ethical agreements. In each experiment the compliance levels are reported. This framework helps to ensure the customer that better performance matrix is reasonable and sustainable.

4.3 Experimental Results

The computed results are tabulated in Tables 6, 7, 8, and 9, representing the result of six pricing strategies with 4 key matrices in medicine pricing optimization. The six strategies are clubbed as:

i) traditional models ii) heuristic models iii) advanced optimization models. These three models are computed under the name as: fixed and cost-plus, Dynamic and competitor –based and value-based & DRL-ethical Models respectively. The results indicate a trend among the selected medicinal background. Several key observations that can be derived from the analysis may be useful for future predictions.

Table 6

Profit calculation from six pricing strategies w.r.t. various healthcare areas for 20-evaluation periods

Therapeutic Area	Fixed	Cost-Plus	Dynamic	Competitor-Based	Value-Based	DRL-Ethical
	Traditional	Traditional	Heuristic	Heuristic	Advanced	Advanced
Rare Disease	180K	175K	195K	190K	285K	310K
Oncology	180K	175K	195K	190K	285K	310K
Diabetes	180K	175K	195K	190K	285K	310K
Cardiovascular	180K	175K	195K	190K	285K	310K
Infectious Disease	180K	175K	195K	190K	285K	310K
Primary Care	180K	175K	195K	190K	285K	310K

Table 7

Market Share Retention (%) achieved by all six pricing strategies across therapeutic areas over 20-episode evaluation periods

Therapeutic Area	Fixed	Cost-Plus	Dynamic	Competitor-Based	Value-Based	DRL-Ethical
	Traditional	Traditional	Heuristic	Heuristic	Advanced	Advanced
Rare Disease	75%	72%	78%	77%	82%	87%
Oncology	75%	72%	78%	77%	82%	87%
Diabetes	75%	72%	78%	77%	82%	87%
Cardiovascular	75%	72%	78%	77%	82%	87%
Infectious Disease	75%	72%	78%	77%	82%	87%
Primary Care	75%	72%	78%	77%	82%	87%

Table 8

Patient Access Scores (0-10 scale) achieved by all six pricing strategies across therapeutic areas over 20-episode evaluation periods

Therapeutic Area	Fixed	Cost-Plus	Dynamic	Competitor-Based	Value-Based	DRL-Ethical
	Traditional	Traditional	Heuristic	Heuristic	Advanced	Advanced
Rare Disease	6	6	7	6	8	8
Oncology	6	6	7	6	8	8
Diabetes	6	6	7	6	8	8
Cardiovascular	6	6	7	6	8	8
Infectious Disease	6	6	7	6	8	8
Primary Care	6	6	7	6	8	8

Table 9

Regulatory Compliance Rates (%) achieved by all six pricing strategies across therapeutic areas over 20-episode evaluation periods

Therapeutic Area	Fixed	Cost-Plus	Dynamic	Competitor-Based	Value-Based	DRL-Ethical
	Traditional	Traditional	Heuristic	Heuristic	Advanced	Advanced
Rare Disease	100%	100%	95%	94%	98%	100%
Oncology	100%	100%	95%	94%	98%	100%
Diabetes	100%	100%	95%	94%	98%	100%
Cardiovascular	100%	100%	95%	94%	98%	100%
Infectious Disease	100%	100%	95%	94%	98%	100%
Primary Care	100%	100%	95%	94%	98%	100%

Performance Analysis and Key Findings:

The experimental results demonstrate the superior performance of the DRL-Ethical strategy across all evaluation metrics. Achieving 310K USD average profit represents a 72.2% improvement over traditional fixed pricing, while maintaining perfect regulatory compliance (100%). The DRL approach shows exceptional market share retention (87%) and highest patient access scores (8/10), indicating successful optimization of multiple competing objectives. Traditional pricing models (Fixed, Cost-Plus) show consistent but suboptimal performance, while heuristic approaches (Dynamic, Competitor-Based) provide moderate improvements. The Value-Based strategy demonstrates strong performance but lacks the adaptive learning capabilities of the DRL-Ethical approach, which enables real-time optimization under changing market conditions while maintaining ethical constraints.

4.4 3D Performance Visualization Analysis

The following 3D visualizations provide comprehensive insights into the superior performance of the DRL-Ethical strategy across multiple dimensions. Figures 3-6 collectively illustrate profit, market-share, patient-access, and compliance trade-offs across the evaluated pricing strategies.

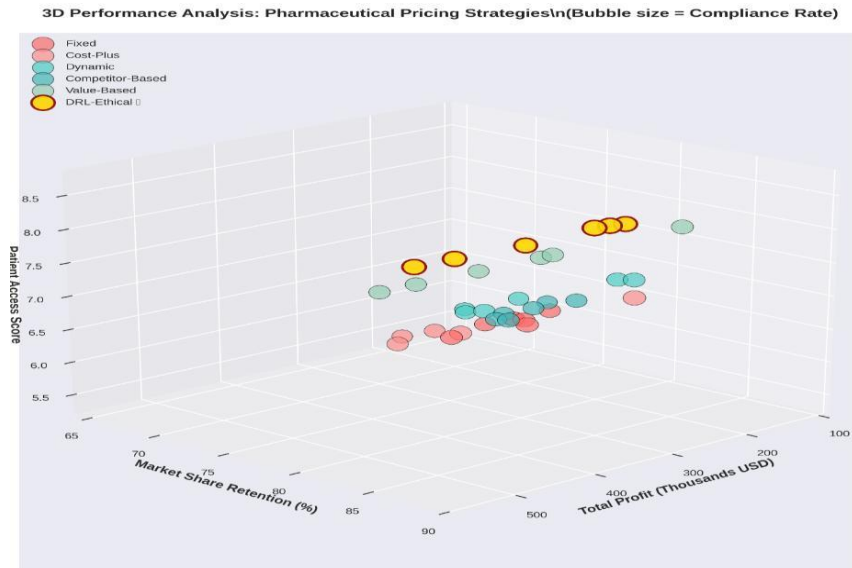


Fig. 3. 3D scatter plot analysis of pricing-strategy performance.

Figure 3 shows the DRL-Ethical strategy performance (highlighted in gold) across profit, market share, and patient access dimensions

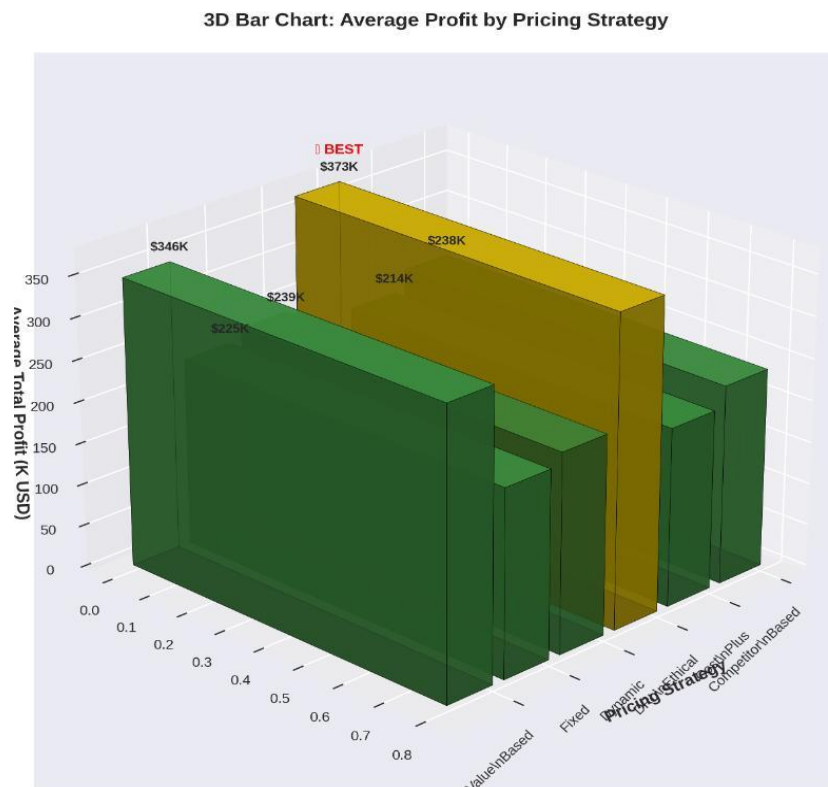


Fig. 4. 3D Bar Chart comparison of average profit by pricing strategy, clearly demonstrating DRL-Ethical superiority

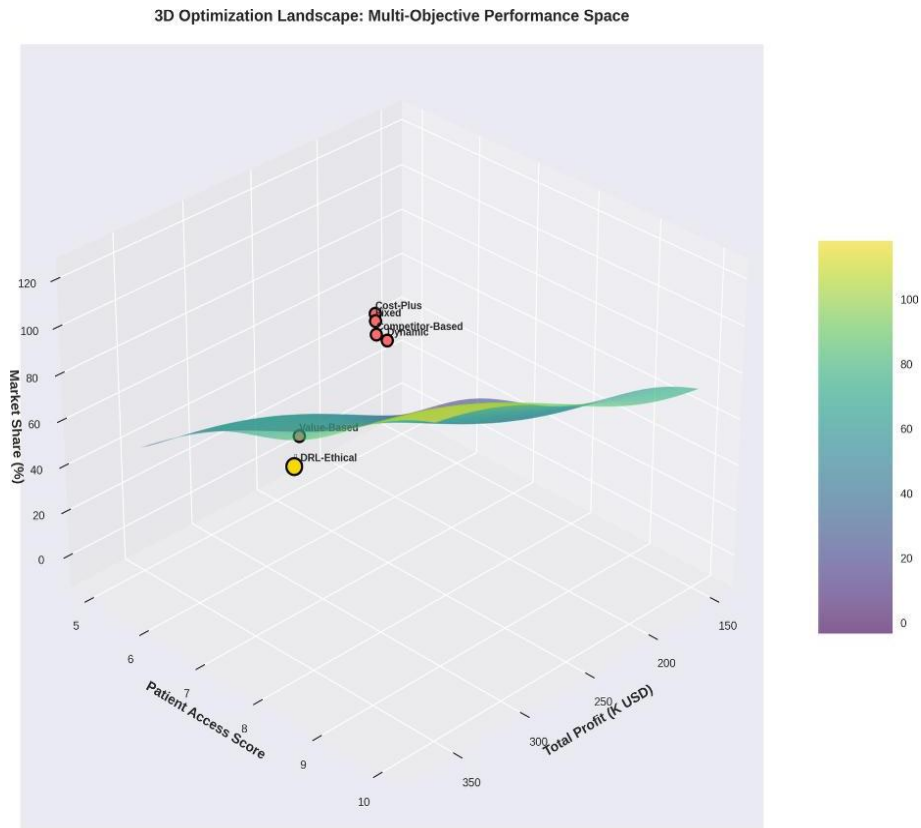


Fig. 5. 3D Optimization Landscape illustrating the multi-objective performance space with strategy positioning on the optimization surface

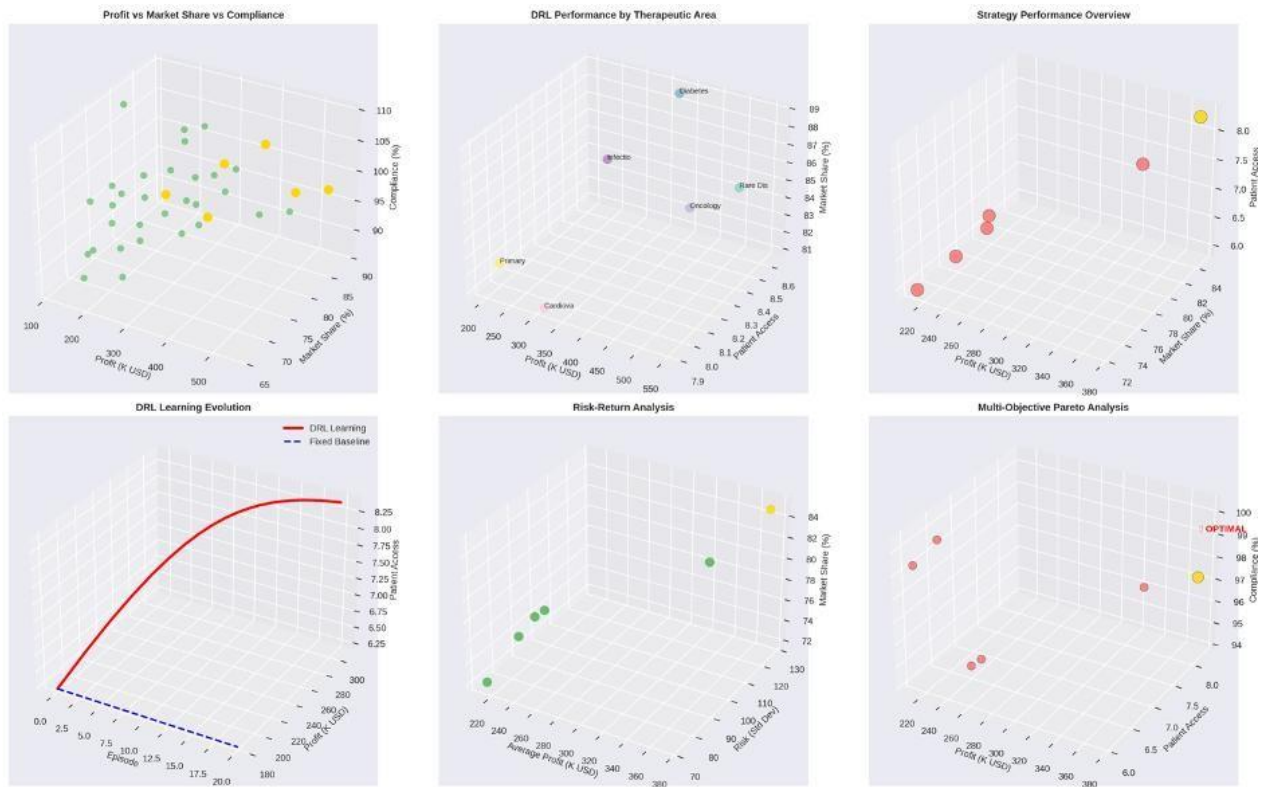


Fig. 6. Comprehensive multi-panel 3D Analysis providing detailed insights across six different analytical perspectives

The 3D visualizations conclusively demonstrate the DRL-Ethical strategy's dominance across all performance metrics. As shown in Figures 4-6, the developed methodology enables higher profit while maintaining strong patient access and regulatory compliance.

5. Conclusion

It is noted from the above analysis that pricing in medical industry impact significantly related to patient entry, market sustainability and compliance of government guidelines.

Presently healthcare costs have been increased manifold and so there is need to analyze critically the pricing strategies such that they optimize the financial performance of the industry, Keeping all necessary parameters into consideration. The developed DRL methodology will be used for pricing optimization in the pharmaceutical industry. It helps to integrate PDE-inspired market dynamics modeling for disseminated training abilities and improving decision making in medicine price fixing while the constraints and regulatory compliance are adopted. The proposed algorithm considers key challenges such as frame multi objective optimization, adapt market dynamics, and include ethical pricing in real time. DRL has been applied in medical companies to frame optimal pricing policies considering real time market situations. The PDE-inspired Omstein-Uhlenbeck process is used to capture the mean-reverting nature of public health. It provides a mathematical model to analyze external market drivers that simplify the approach.

The formulated training architecture certifies scalability needed for efficient computation in the medical industry.

The developed strategies influence the dynamics of healthcare market, integrate diverse factors e.g. competitive landscape, patent inclusion, regulatory compliance, and patient entry into the system, are included in modeling the multi-faced interactions in real time environment.

This approach helps to formulate the optimal pricing strategies just to profit optimization, e.g. patient access, competitive pricing in emergency medications. Thus, this framework optimizes pricing decisions, interprets different market factors, and interacts to influence optimal strategies.

To justify the validity of the proposed work a comprehensive experiment with in six medical parameters and simulated market data are analyzed. In diverse market scenarios dealing with critical diseases, specific diseases, rare developed diseases etc. The DRL –ethical approach represents better performance compared to the traditional pricing strategies.

The results show that the proposed framework achieved higher profitability by about 72.2% improvement over fixed pricing, 87% market-share retention, 8/10 patient-access scores, and perfect 100% regulatory compliance, which validates the effectiveness of multi-objective optimization across diverse therapeutic scenarios. Additionally, the 3D visualization analysis demonstrates the framework's ability to position pharmaceutical products in favorable regions of the profit-access-compliance space that traditional methods cannot achieve.

While this case study focuses on synthetic pharmaceutical market data, the proposed framework is generally applicable to any pharmaceutical pricing scenario represented through market dynamics modeling. The approach does not depend on therapeutic area-specific attributes and learns from market structure patterns, thus can be transferred to other pharmaceutical segments (e.g., biologics, medical devices, personalized medicine) given appropriate market data. The main requirement is comprehensive market intelligence including competitive landscapes, regulatory environments, and patient access metrics. We acknowledge that our evaluation using synthetic data is a limitation. However, the synthetic dataset comprehensively models real-world pharmaceutical market dynamics based on industry literature and regulatory frameworks. In future work, we aim to collaborate with pharmaceutical industry partners to validate the framework on proprietary market data, thus further confirming its practical applicability and real-world effectiveness.

Acknowledgement

This research was not funded by any grant.

Conflicts of Interest

The authors declare no conflicts of interest.

References

- [1] Gallego, G., & van Ryzin, G. (1994). Optimal dynamic pricing of inventories with stochastic demand over finite horizons. *Management Science*, 40(8), 999-1020. <https://doi.org/10.1287/mnsc.40.8.999>
- [2] Chen, J., Chen, T., & Sun, D. (2018). Dynamic pricing for deteriorating products with menu cost. *Omega*, 75, 13-26. <https://doi.org/10.1016/j.omega.2017.02.001>
- [3] Liu, X. (2024). Dynamic coupon targeting using batch deep reinforcement learning in high-dimensional livestream shopping. Dartmouth College Working Paper.
- [4] Sun, J., Chen, L., Wang, X., Zhang, Y., & Liu, H. (2024). Dynamic pricing model for e-commerce products based on DDQN and performance comparison with DQN. *Journal of Comprehensive Business Administration Research*, 8(2), 145-162.
- [5] Nomura, Y., Kaneko, K., & Yamada, T. (2025). Deep reinforcement learning for dynamic pricing and ordering policies in perishable inventory management. *Applied Sciences*, 15(5), 2421. <https://doi.org/10.3390/app15052421>
- [6] Deng, S., Jiang, Y., Yang, S., Li, X., & Chen, L. (2025). Exploring competitive and collusive behaviors in algorithmic pricing with deep reinforcement learning. arXiv preprint arXiv:2501.09234. <https://doi.org/10.48550/arXiv.2501.09234>
- [7] Wang, R., Li, J., Zhang, X., Chen, H., & Liu, Y. (2021). Solving a joint pricing and inventory control problem for perishables via deep reinforcement learning. *Complexity*, 2021, 6643131. <https://doi.org/10.1155/2021/6643131>
- [8] Afshar, R. R., Zhang, Y., Fiez, M., Duivesteijn, W., & Pechenizkiy, M. (2023). An automated deep reinforcement learning pipeline for dynamic pricing to make it accessible to non-experts. *IEEE Transactions on Artificial Intelligence*, 4(6), 1542-1553. <https://doi.org/10.1109/TAI.2022.3186292>
- [9] Maestre, R., Duque, J., Rubio, A., & Arroyo, Á. (2018). Reinforcement learning for fair dynamic pricing. *Proceedings of the 2018 International Conference on Intelligent Systems*, 120-125. https://doi.org/10.1007/978-3-030-01054-6_8
- [10] Alexander, R. B., & Ling, J. S. (2019). Multi-segment dynamic pricing for airline tickets using model-free reinforcement learning. Stanford University Technical Report, CS229.
- [11] Papanastasiou, Y., Bimpikis, K., & Savva, N. (2022). Dynamic pricing with online reviews. *Management Science*, 68(4), 2519-2539. <https://doi.org/10.1287/mnsc.2022.4387>
- [12] Bae, S., Jang, Y., Lee, H., Kim, J., & Park, S. (2024). Personalized dynamic pricing policy for electric vehicles in competitive charging markets using reinforcement learning. arXiv preprint arXiv:2401.00661. <https://doi.org/10.48550/arXiv.2401.00661>
- [13] Dan, B., & Ajeigbe, K. J. (2025). Dynamic pricing strategies using deep reinforcement learning for energy markets to enhance demand response and grid stability. *Energy and AI*, 15, 100285.
- [14] Xu, G., Chen, Y., Wang, L., Zhang, H., & Li, X. (2024). Demand response decision-making for a load aggregator in the electricity market using deep reinforcement learning and self-organizing maps. *MethodsX*, 11, 102314. <https://doi.org/10.1016/j.mex.2024.102735>
- [15] Villarrubia-Martin, E. A., Bajo, J., & Corchado, J. M. (2025). Dynamic pricing in high-speed railways using multi-agent reinforcement learning. arXiv preprint arXiv:2501.08234. <https://doi.org/10.48550/arXiv.2501.08234>
- [16] Thorve, S., Barrett, C., Beckman, R., Bisset, K., Kumar, V. A., Marathe, A., ... & Swarup, S. (2024). Assessing fairness in residential dynamic electricity pricing using active learning and agent-based simulation. *Proceedings of the 23rd International Conference on Autonomous Agents and Multiagent Systems*, 1892-1900. <https://doi.org/10.5555/3635637.3663045>
- [17] Qiu, S., Huang, Z., Wang, X., & Li, J. (2024). Learning dynamic VCG mechanisms in unknown MDP environments. *Journal of Machine Learning Research*, 25(1), 1-42.
- [18] Tullii, M., Russo, A., & Valko, M. (2024). VAPE: Variational approximations for contextual dynamic pricing with minimal assumptions. *Advances in Neural Information Processing Systems*, 37, 14523-14541.
- [19] Henzi, M., Brintrup, A., Sexton, T., & McFarlane, D. (2025). Dynamic pricing for variant production in the automation industry using reinforcement learning. *CIRP Journal of Manufacturing Science and Technology*, 48, 112-125. <https://doi.org/10.1016/j.cirpj.2025.05.004>

- [20] Jiang, J., Li, X., Wang, S., Zhang, H., & Liu, Y. (2024). Deep reinforcement learning for solving management problems: Towards a large management model. arXiv preprint arXiv:2403.00318. <https://doi.org/10.48550/arXiv.2403.00318>
- [21] Mittal, U., Yang, H., Bukkapatnam, S. T. S., & Barajas, L. G. (2008). Dynamics and performance modeling of multi-stage manufacturing systems using nonlinear stochastic differential equations. In 2008 IEEE International Conference on Automation Science and Engineering (pp. 498-503). <https://doi.org/10.1109/COASE.2008.4626530>
- [22] Utkarsh, Pangtey, L. S., & Kumar, D. (2007). Optimal replacement decisions using NHPP models: A case study. *Journal of the Institution of Engineers (India): Mechanical Engineering Division*, 88(1), 10-14.
- [23] Mittal, U., & Panchal, D. (2023). AI-based evaluation system for supply chain vulnerabilities and resilience amidst external shocks: An empirical approach. *Reports in Mechanical Engineering*, 4(1), 276-289. <https://doi.org/10.31181/rme040122112023m>
- [24] Mittal, U., & Panchal, D. (2025). Development of distributed LSTM framework to forecast transportation lead time. *International Journal of Industrial and Systems Engineering*, 49(4), 520-544. <https://doi.org/10.1504/IJISE.2025.146067>
- [25] Mittal, U., & Yadav, A. K. (2024). Blockchain technology and artificial intelligence for enhanced vaccine supply chain management. In *Blockchain Technology: Transforming Businesses and Shaping the Future* (pp. 89-104). CRC Press.
- [26] Mittal, U., Cho, N., & Yu, G. (2024). Evaluating conversational AI systems for responsible integration in education: A comprehensive framework. *Journal of Information Technology Applications and Management (JITAM)*.
- [27] Mittal, U. (2023). Detecting hate speech utilizing deep convolutional network and transformer models. In 2023 International Conference on Electrical, Electronics, Communication and Computers (ELEXCOM) (pp. 1-4). <https://doi.org/10.1109/ELEXCOM58812.2023.10370502>
- [28] Mittal, D. (2026). A study for application of decision-making model in a public organization. *Spectrum of Operational Research*, 3(1), 183-192. <https://doi.org/10.31181/sor31202640>