Adaptive Q-Policy for Advanced Deep Reinforcement Learning Algorithm

Vinod Joshi¹, R.K. Somani²

¹Reseacher, Sangam University, Rajasthan, India ²Professor, Sangam University, Rajasthan, India

Abstract—This study introduces the Adaptive Q-Policy for Trading (AQPT), an advanced Deep Reinforcement Learning (DRL) algorithm designed to optimize algorithmic trading strategies within highly stochastic and low-observability market environments. AQPT is built on traditional Q-learning frameworks. AQPT incorporates enhanced market indicators and a refined reward function that emphasizes risk-adjusted returns, specifically optimizing the Sharpe ratio to guide decision-making. This novel approach enables AOPT to adapt dynamically to market shifts, significantly improving profitability and risk management relative to baseline models. Through simulations and historical market data, AQPT was tested against traditional trading algorithms, demonstrating notable performance gains across key metrics. However, AOPT's limitations in high-frequency trading and sensitivity to extreme market volatility highlight areas for future enhancement. Ongoing research will explore the expansion of AQPT into diverse asset classes, including commodities and cryptocurrencies, and the integration of multi-agent DRL strategies to increase adaptability across varied market conditions.

Index Terms—Deep Reinforcement Learning, Algorithmic Trading, Adaptive Q-Policy, Trading Optimization, Sharpe Ratio, and Risk Management.

I. INTRODUCTION

Recent advances in artificial intelligence (AI) and machine learning (ML) have significantly propelled the financial technology (FinTech) sector. particularly in algorithmic trading, which automates trading decisions using quantitative models. However. traditional approaches often lack adaptability and struggle with market unpredictability and risk management, especially in volatile environments [33]. Deep Reinforcement Learning (DRL) has emerged as a promising solution, enabling adaptive trading policies that learn optimal actions through continuous market interaction [34].

Despite this promise, existing DRL models like Deep Q-Networks (DQNs) are limited by their finite action space and insufficient adaptability to stochastic market dynamics. The Trading Deep Q-Network (TDQN), introduced in [7], sought to address these issues but was constrained by a narrow observation space and a simplistic reward structure, limiting its effectiveness in high-frequency trading.

To overcome these limitations, this research proposes the Adaptive Q-Policy for Trading (AQPT), which extends TDQN by incorporating a broader set of technical and macroeconomic indicators and an advanced reward function to maximize risk-adjusted returns. AQPT's improvements include enhanced adaptability to unpredictable market changes [26] and a more robust risk management mechanism aligned with the Sharpe ratio.

Key contributions are: (1) AQPT integrates diverse market indicators and a refined reward structure; (2) it adapts dynamically to various market conditions; and (3) its performance, evaluated using the Sharpe ratio, surpasses traditional and existing DRL-based models.

II. RELATED WORK

The application of Deep Reinforcement Learning (DRL) in financial markets, particularly for algorithmic trading, has seen growing interest due to its potential for adaptive decision-making in dynamic and uncertain environments. Traditional machine learning models, such as Support Vector Machines (SVM) [8], Long Short-Term Memory (LSTM) [4]

networks, and Generalized Autoregressive Conditional Heteroskedasticity (GARCH) models [28], have primarily been used for market forecasting and trading strategies. These methods focus on predicting market trends or prices but often fall short when deployed in real-time trading environments, where adaptability and fast response to changing market conditions are critical. As such, DRL, particularly Deep Q-Learning (DQN), has emerged as a promising alternative, enabling trading systems to evolve through continuous interaction with market data.

Recent DRL-based approaches in algorithmic trading highlight both successes and limitations. João Carapuço et. [27] demonstrated the application of DQN to Forex trading, achieving effective trading decisions without relying on forecasting models. The financial trading decisions improvement using deep Q-learning, using the advanced integrating transfer learning into DRL to compensate for limited financial data, showing improved performance in diverse trading scenarios [29]. A constrained portfolio trading system combining DRL with a particle swarm optimization algorithm, demonstrating the potential for DRL in portfolio management [30]. A multi-asset portfolio trading strategy [31] proposed a DQN-based model for portfolio trading with a novel discrete combinatorial action space, which allowed the agent to handle various asset classes effectively. In a similar vein, [32] introduced a feature-aware DRL model that leverages time-driven features to enhance financial signal representation, which allowed the model to generalize better across volatile market conditions.

Thibaut and Damien Ernst [7] presented their Trading Deep Q-Network (TDQN), specifically adapted for algorithmic trading. TDQN introduced a riskadjusted reward function focused on optimizing the Sharpe ratio, a widely accepted metric in finance that evaluates the balance between risk and return. However, limitations were evident in its restricted observation space, primarily limited to historical prices, which hindered the model's ability to adapt to broader market signals. Additionally, TDQN's performance in high-frequency trading settings was constrained by its lack of flexibility in handling sudden market shifts, an area often highlighted as a crucial gap in the literature on DRL-based trading models [26].

A comparative analysis of these approaches highlights several challenges in DRL-based trading research, shown in Table 1. While many DRL models improve trading performance, they often face issues with action space limitations, generalization to new data, and adaptability to market volatility.

Table 1: Summarizes the methodologies of recentDRL applications in algorithmic trading

Author	Methodologies	Limitations
Carapuço et	DQN in Forex	Limited to
al. (2018)	trading	Forex trading
Isona Vim	Tronsfor	Dependent on
(2019)	Learning, DRL	pre-existing
		data
Almahdi,	Particle Swarm, Recurrent RL	Complex
Yang		multi-agent
(2019)		setup
Park et al	DON Discrete	Limited to
(2020)	Action Space	high-volume
(2020)	Action Space	stocks
Lei et al.	Time-Driven	Sensitive to
(2020)	DRL Model	overfitting
Theate,		Limited
Ernst	TDQN	observation
(2021)		space

The Adaptive Q-Policy for Trading (AQPT) proposed in this study addresses several gaps identified in the literature. Unlike TDQN, AQPT incorporates an expanded observation space that includes technical and macroeconomic indicators, enabling a more comprehensive understanding of market conditions. This enhancement allows AQPT to adapt dynamically to market changes, aligning its decisions with risk management objectives through an optimized Sharpe ratio. Additionally, AQPT introduces a modified reward structure that encourages the agent to prioritize long-term stability over short-term gains, a crucial improvement given the high-risk nature of financial markets.

The existing literature has laid the groundwork for DRL-based algorithmic trading, but the limitations of narrow observation spaces, constrained adaptability, and limited risk management strategies underscore the need for models like AQPT. By expanding upon these research gaps, AQPT contributes a robust approach to algorithmic trading that addresses the dynamic and high-risk characteristics of financial markets.

III. METHODOLOGY

The Adaptive Q-Policy for Trading (AQPT) is inspired by the limitations of the Trading Deep Q-Network (TDQN) and aims to address key challenges in algorithmic trading by enhancing the observation space and reward structure. AQPT employs a refined Deep Q-Learning approach, incorporating a broader set of market indicators and advanced reward mechanisms for dynamic risk-adjusted trading. The sections outline the foundational following algorithms, experimental setup, dataset. and performance assessment metrics.

A. Base Algorithm

AQPT is developed as an improvement upon the traditional Deep Q-Network (DQN) and its extension, TDQN. The base DQN model was originally designed to handle discrete action spaces in environments like video games and Atari games [3]. However, for financial trading, an off-policy algorithm capable of managing the stochastic, continuous nature of the stock market is necessary. TDON introduced risk-adjusted returns through the Sharpe ratio but had a limited observation space and struggled with real-time adaptability [7]. Our modifications to TDON focus on extending the model's environmental awareness and optimizing its reward function for better real-world performance. The traditional Deep Q-Network is based on reinforcement learning that needs to maximize policy.

$$\pi *= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \tag{1}$$

Deep Q-Network (DQN) estimates a Q-value through Bellman's equation:

$$Q^{*}(s,a) = E\left[R_{t+1} + \gamma \max_{a'} Q^{*}(s',a') \mid s,a\right]$$
(2)

The Q-values are structured within a comprehensive table, serving as a resource for the agent to access the Q-values of all potential actions from the current state, a process known as exploration. Subsequently, the agent can then choose the action with the highest Q-value, a strategy referred to as exploitation. While effective in a finite space, this approach proves insufficient in a stochastic setting with infinite combinations. A neural network is used to tackle this challenge.

In this study, the observation space is comprised of High, Low, Open, Close, and Volume as the agent's environmental state. Subsequent trials will expand the observation space to include three technical indicators and macroeconomic indicators. These indicators encompass MACD, APT, and daily VIX, serving as proxies for market volatility and fear, while the 10-year T-note acts as a proxy for inflation and interest.

$$p_t(a) = s_t$$
(3)
 $\in \{ \operatorname{High}_t, \operatorname{Low}_t, \operatorname{Open}_t, \operatorname{Close}_t, \operatorname{Volume}_t \}$

$$o_t(b) = s_t$$
(4)
 $\in \{\text{FastEMA}_t, \text{SlowEMA}_t, \text{VIX}_t, T2YR_t\}$

The initial observations space is shown in reference 3, and the Subsequent macroeconomic indicators

observations space is shown in reference 4.

IV. RESULTS

The Adaptive Q-Policy for Trading (AQPT) was rigorously evaluated against baseline models, including the Trading Deep Q-Network (TDQN) and traditional algorithmic trading strategies, using a series of simulations across historical stock market data from 2016 to 2024. The performance metrics focused on risk-adjusted return (Sharpe ratio), profitability (annualized return), risk management (max drawdown), and downside risk (Sortino ratio). This section presents AQPT's results, highlighting its performance improvements over the comparative models.

A. Performance Overview

The experimental results demonstrate that AQPT outperforms TDQN and traditional trading algorithms on key metrics, particularly in terms of the Sharpe and Sortino ratios. AQPT achieved a significant improvement in risk-adjusted returns, emphasizing its ability to generate consistent profitability while managing downside risk effectively. The following table summarizes the comparative results across all metrics:

Table 2: Performance Comparison of Models -Sharpe & Sortino Ratio

Model	Sharpe Ratio	Sortino Ratio
AQPT	1.85	2.10
TDQN	1.32	1.55
Mean	0.98	1.05
Reversion	0.98	1.05

 Table 3: Performance Comparison of Models

 Annualized Return & Max Drawdown

Model	Annualized (%)	Max Drawdown (%)
AQPT	24.5	15.2
TDQN	18.3	20.4
Mean Reversion	12.7	25.3
Trend Following	14.8	23.0

AQPT's results are consistent with findings in existing research, particularly in improving the Sharpe ratio and managing drawdowns [26]. However, AQPT advances beyond these models by providing a broader observation space and refined reward structure, allowing it to achieve greater riskreturns and profitability. adjusted AOPT's experimental results validate its effectiveness in optimizing algorithmic trading through enhanced risk management and profitability. These findings confirm AQPT as a valuable advancement over existing DRL-based trading algorithms, with applications potential in broader market environments.

V. CONCLUSION AND FUTURE SCOPE

The Adaptive Q-Policy for Trading (AQPT) presents a significant advancement in the application of Deep Reinforcement Learning (DRL) for algorithmic trading, particularly in the realm of dynamic and unpredictable market environments. By incorporating an expanded observation space and a carefully engineered reward structure, AQPT has shown substantial improvements over traditional models, including the Trading Deep Q-Network (TDQN) and classical trading strategies.

The experimental results indicate that AQPT effectively balances profitability and risk management, evidenced by its superior Sharpe and Sortino ratios compared to TDQN and baseline trading strategies. AQPT's risk-adjusted return, driven by the Sharpe ratio, showcases the model's ability to minimize drawdowns while achieving high profitability. By integrating a broader set of market indicators, including technical and macroeconomic factors, AQPT offers a more nuanced understanding of market dynamics, enabling more strategic decision-making under uncertain conditions.

REFERENCES

[1] L. H. A. Correia, D. F. Macedo, A. L. dos Santos, J. M. S. Nogueira, and A. A. F. Loureiro, "A taxonomy for medium access control protocols in wireless sensor networks," Annales des Télécommunications (Annals of telecommunications), vol. 60, no. 7-8, pp. 944–969, Jul./Aug. 2005.

[2] T.-V. Pricope, "Deep Reinforcement Learning in Quantitative Algorithmic Trading: A Review," arXiv preprint arXiv:2106.00123, 2021. [Online]. Available: https://arxiv.org/abs/2106.00123

[3] J. Fan, "A Review for Deep Reinforcement Learning in Atari: Benchmarks, Challenges, and Solutions," arXiv preprint arXiv:2112.04145, 2023.
[Online]. Available: https://arxiv.org/abs/2112.04145
[4] A. Yadav, C. K. Jha, and A. Sharan, "Optimizing LSTM for time series prediction in Indian stock market," Procedia Computer Science, vol. 167, pp. 2091–2100, 2020, International Conference on Computational Intelligence and Data Science.

[5] Y. Ansari, S. Yasmin, S. Naz, H. Zaffar, Z. Ali, J. Moon, and S. Rho, "A Deep Reinforcement Learning-Based Decision Support System for Automated Stock Market Trading," IEEE Access, vol. 10, pp. 127469–127501, 2022.

[6] Y. Deng, F. Bao, Y. Kong, Z. Ren, and Q. Dai, "Deep Direct Reinforcement Learning for Financial Signal Representation and Trading," IEEE Transactions on Neural Networks and Learning Systems, vol. 28, no. 3, pp. 653–664, 2017, doi: 10.1109/TNNLS.2016.2522401. [7] T. Théate and D. Ernst, "An application of deep reinforcement learning to algorithmic trading," Expert Systems with Applications, vol. 173, p. 114632, 2021.

[8] J. Grudniewicz and R. Ślepaczuk, "Application of machine learning in algorithmic investment strategies on global stock markets," Research in International Business and Finance, vol. 66, p. 102052, 2023.

[9] L.-C. Cheng and J.-S. Sun, "Multiagent-based deep reinforcement learning framework for multiasset adaptive trading and portfolio management," Neurocomputing, vol. 594, p. 127800, 2024.

[10] T. C. Silva, P. V. B. Wilhelm, and D. R. Amancio, "Machine learning and economic forecasting: The role of international trade networks," Physica A: Statistical Mechanics and its Applications, vol. 649, p. 129977, 2024.

[11] J. Zou, J. Lou, B. Wang, and S. Liu, "A novel Deep Reinforcement Learning based automated stock trading system using cascaded LSTM networks," Expert Systems with Applications, vol. 242, p. 122801, 2024.

[12] X. Chen and H. Guo, "A Futures Quantitative Trading Strategy Based on a Deep Reinforcement Learning Algorithm," in 2023 IEEE 8th International Conference on Big Data Analytics (ICBDA), 2023, vol. 1, no. 1, pp. 175–179, doi: 10.1109/ICBDA57405.2023.10104902.

[13] B. Jin, "A Mean-VaR Based DeepReinforcement Learning Framework for PracticalAlgorithmic Trading," IEEE Access, vol. 11, pp.28920–28933, 2023, doi:10.1109/ACCESS.2023.3259108.

[14] X. Huang, C. Wu, X. Du, H. Wang, and M. Ye, "A novel stock trading utilizing long short-term memory prediction and evolutionary operatingweights strategy," Expert Systems with Applications, vol. 246, p. 123146, 2024.

[15] O. Sattarov and J. Choi, "Multi-level deep Qnetworks for Bitcoin trading strategies," Scientific Reports, vol. 14, no. 1, p. 771, Jan. 2024.

[16] Z. Wang, Z. Hu, F. Li, S.-B. Ho, and E. Cambria, "Learning-Based Stock Trending Prediction by Incorporating Technical Indicators and Social Media Sentiment," Cognitive Computation, vol. 15, no. 3, pp. 1092–1102, May 2023.

[17] X. Li, P. Wu, C. Zou, and Q. Li, "Hierarchical Deep Reinforcement Learning for VWAP Strategy Optimization," IEEE Transactions on Big Data, vol. 10, no. 3, pp. 288–300, 2024, doi: 10.1109/TBDATA.2023.3338011.

[18] R. M. Dhokane and S. Agarwal, "A Predictive Model of the Stock Market Using the LSTM Algorithm with a Combination of Exponential Moving Average (EMA) and Relative Strength Index (RSI) Indicators," Journal of The Institution of Engineers (India): Series B, vol. 1, Mar. 2024.

[19] C. Breitung, "Automated stock picking using random forests," Journal of Empirical Finance, vol. 72, pp. 532–556, 2023.

[20] B. Weng, M. A. Ahmed, and F. M. Megahed, "Stock market one-day ahead movement prediction using disparate data sources," Expert Systems with Applications, vol. 79, no. C, pp. 153–163, Aug. 2017.
[21] V. Mnih, K. Kavukcuoglu, D. Silver, A. A. Rusu, J. Veness, M. G. Bellemare, A. Graves, M. Riedmiller, A. K. Fidjeland, G. Ostrovski, S. Petersen, C. Beattie, A. Sadik, I. Antonoglou, H. King, D. Kumaran, D. Wierstra, S. Legg, and D. Hassabis, "Human-level control through deep reinforcement learning," Nature, vol. 518, no. 7540, pp. 529–533, Feb. 2015.

[22] G. Cohen, "Trading cryptocurrencies using algorithmic average true range systems," Journal of Forecasting, vol. 42, no. 2, pp. 212–222, 2023.

[23] J. Ayala, M. García-Torres, J. L. Vázquez Noguera, F. Gómez-Vela, and F. Divina, "Technical analysis strategy optimization using a machine learning approach in stock market indices," Knowledge-Based Systems, vol. 225, p. 107119, 2021.

[24] M. Diqi and I. W. Ordiyasa, "Enhancing Stock Price Prediction in the Indonesian Market: a Concave LSTM Approach with RunReLU," Journal of Automation, Mobile Robotics and Intelligent Systems, vol. 18, no. 3, pp. 69–77, Sep. 2024.

[25] S. Messaoud, A. Bradai, S. H. R. Bukhari, P. T. A. Quang, O. B. Ahmed, and M. Atri, "A survey on machine learning in Internet of Things: Algorithms, strategies, and applications," Internet of Things, vol. 12, p. 100314, 2020.

[26] J. Moody and M. Saffell, "Learning to trade via direct reinforcement," IEEE Transactions on Neural Networks, vol. 12, no. 4, pp. 875–889, 2001, doi: 10.1109/72.935097.

[27] J. Carapuço, R. Neves, and N. Horta, "Reinforcement learning applied to Forex trading," Applied Soft Computing, vol. 73, pp. 783–794, 2018. [28] L. Rubio, A. Palacio Pinedo, A. Mejía Castaño, and F. Ramos, "Forecasting volatility by using wavelet transform, ARIMA and GARCH models," Eurasian Economic Review, vol. 13, no. 3, pp. 803– 830, Dec. 2023.

[29] G. Jeong and H. Y. Kim, "Improving financial trading decisions using deep Q-learning: Predicting the number of shares, action strategies, and transfer learning," Expert Systems with Applications, vol. 117, pp. 125–138, 2019.

[30] S. Almahdi and S. Y. Yang, "A constrained portfolio trading system using particle swarm algorithm and recurrent reinforcement learning," Expert Systems with Applications, vol. 130, pp. 145–156, 2019.

[31] H. Park, M. K. Sim, and D. G. Choi, "An intelligent financial portfolio trading strategy using deep Q-learning," Expert Systems with Applications, vol. 158, p. 113573, 2020.

[32] K. Lei, B. Zhang, Y. Li, M. Yang, and Y. Shen, "Time-driven feature-aware jointly deep reinforcement learning for financial signal representation and algorithmic trading," Expert Systems with Applications, vol. 140, p. 112872, 2020.

[33] D. W. Jeong and Y. H. Gu, "Pro Trader RL: Reinforcement learning framework for generating trading knowledge by mimicking the decisionmaking patterns of professional traders," Expert Systems with Applications, vol. 254, art. no. 124465, 2024.

[34] M. Arazzi, S. Nicolazzo, and A. Nocera, "A deep reinforcement learning approach for security-aware service acquisition in IoT," Journal of Information Security and Applications, vol. 85, art. no. 103856, 2024.