

Reinforcement Learning Meets Quantum Optimization: A New Paradigm for Dynamic Portfolio Management

Author: Olivia Perez, **Affiliation:** Postdoctoral Researcher, Robotics and Machine Intelligence Lab, University of Chile. **Email:** olivia.perez@uchile.cl

Abstract

We propose and analyze a hybrid paradigm that integrates Reinforcement Learning (RL) with Quantum Optimization (QO) methods for dynamic portfolio management. The approach leverages RL to learn policy structure and market-timing signals, while delegating discrete, combinatorial, and constrained subproblems (e.g., cardinality-constrained selection, rebalancing under transaction limits) to quantum optimization engines such as Quantum Approximate Optimization Algorithm (QAOA), Variational Quantum Eigensolver (VQE), and quantum annealers via QUBO/HUBO encodings. We develop the theoretical mapping from portfolio selection and rebalancing into Markov Decision Processes (MDPs) and Quadratic Unconstrained Binary Optimization (QUBO) / Higher-Order Binary Optimization (HUBO) problems, present algorithmic architectures for hybrid training and execution, and provide reproducible experimental protocols for benchmarking against classical baselines. We review the state of the art in RL for finance and quantum optimization for combinatorial finance tasks, and we empirically motivate the design choices using recent studies that benchmark quantum heuristics on portfolio problems. Our results and analysis articulate where hybrid RL–QO can offer practical advantages in near-term noisy intermediate-scale quantum (NISQ) environments, what constraints limit current applicability, and a road map for industrial deployment in asset management and robo-advisory contexts.

Keywords: Quantum optimization, Reinforcement Learning, Portfolio optimization, QUBO, QAOA, VQE, hybrid quantum-classical, NISQ, dynamic allocation

1. Introduction

Modern portfolio management combines statistical estimation, optimization under constraints, and sequential decision-making in uncertain markets. Since Markowitz introduced mean–variance optimization (MVO) as a formal framework for portfolio selection, the field has extended to incorporate risk measures, transaction costs, and dynamic rebalancing strategies (Markowitz, 1952).

Computationally, many practical constraints cardinality, minimum-lot sizes, nonlinear transaction costs, and regulatory or fund-specific limits render portfolio optimization NP-hard or combinatorially difficult; exact solutions often require mixed-integer programming or tailored heuristics. Quantum optimization methods (quantum annealing, QAOA, VQE) can naturally encode such discrete constraints as Ising/QUBO problems and have been explored experimentally for finance problems including portfolio optimization; recent benchmark studies and practical experiments demonstrate both promise and the current hardware limits of these approaches.

Concurrently, Reinforcement Learning (RL) particularly deep RL has emerged as a powerful methodology for sequential portfolio management and automated trading, providing model-free strategies that learn allocation policies from raw market data and reward signals (e.g., wealth growth, risk-adjusted returns) without explicit forecasting models. The seminal deep RL portfolio work and subsequent developments show that RL can offer robust, adaptive policies for dynamic allocation. ([arXiv](#))

This paper brings these threads together. We propose a hybrid architecture in which an RL agent learns high-level, continuous control and market-response

policies, while quantum optimizers tackle discrete allocation subproblems and constrained combinatorial rebalancing (e.g., choose K out of N assets to rebalance, or mapping fractional allocations to trade orders under discrete lot sizes). By combining RL's adaptability with QO's native discrete optimization strengths, the paradigm aims to: (i) improve solution quality for constrained rebalancing; (ii) reduce computation time for specific NP-hard subproblems (in regimes where quantum advantage may be achievable); and (iii) provide a practical transition path for asset managers to exploit near-term quantum hardware via hybrid workflows. Recent surveys and benchmarks of quantum reinforcement learning and quantum optimization provide technical foundations and empirical evidence for this direction.

In what follows we present: background theory; formal problem statements; method design and algorithms; experimental design and evaluation metrics; discussion on implementation, robustness, and regulatory considerations; and a roadmap for research and industry adoption.

2. Background and Related Work

2.1 Classical portfolio theory and computational challenges

Markowitz's mean–variance framework formulates portfolio selection as an optimization of expected return against variance (risk) under linear budget constraints. In continuous allocation form, for weights $w \in \mathbb{R}^N$:

$$\begin{aligned} \min_w \quad & w^T \Sigma w - \lambda \mu^T w \\ \text{s.t.} \quad & \mathbf{1}^T w = 1, w_i \geq 0; (\text{if no shorting}), \end{aligned}$$

where Σ is the covariance matrix, μ the expected returns vector, and λ the risk–return

trade-off. While the continuous MVO problem is tractable, introducing real-world integer constraints (cardinality, minimum-lot sizes), nonlinear transaction costs, and regime-dependent constraints converts the problem into a combinatorial optimization, often expressed as Mixed-Integer Quadratic Programming (MIQP). These versions are NP-hard and motivate heuristic and specialized methods. (Markowitz, 1952).

2.2 Reinforcement learning for portfolio management

RL casts portfolio management as an MDP: states encode market observables and portfolio holdings, actions change allocations or trade orders, and rewards align with investment objectives (e.g., log-return, Sharpe-ratio proxies, or more sophisticated utility measures). Deep RL agents (policy/value-function approximators) have been demonstrated to learn effective strategies for sequential allocation, including notable architectures like the EIIE topology and actor–critic variants, and task-specific algorithms addressing transaction costs and market microstructure. (Jiang et al., 2017; Yang, 2023). ([arXiv](#))

Major challenges for RL in finance include: sample efficiency (markets are nonstationary and data-limited), overfitting/backtest over-optimism, partial observability, reward sparsity, and safety/regulatory constraints during live deployment. Strong evaluation protocols (walk-forward testing, out-of-sample robustness checks, bootstrapped statistical testing) are essential to validate RL policies.

2.3 Quantum optimization methods relevant to finance

Quantum optimization approaches relevant for portfolio tasks include:

- **Quantum Annealing (QA):** hardware such as D-Wave implements QA to find low-energy states of Ising Hamiltonians; mapping discrete portfolio problems to Ising/QUBO is straightforward and has been demonstrated for small-to-moderate instance sizes.

- **Quantum Approximate Optimization Algorithm (QAOA):** a gate-based variational algorithm that alternates between cost and mixer Hamiltonians to sample low-energy states; QAOA is suited to combinatorial problems encoded as QUBO/Hamiltonian minimization and has been applied to knapsack- and portfolio-like problems.
- **Variational Quantum Eigensolver (VQE):** originally for chemistry, VQE has been used in finance by encoding cost functions into Hamiltonians and optimizing parameterized circuits to minimize expected cost on quantum devices; practical experiments have been performed on IBM devices for small-size portfolio instances.

Recent benchmarks compare QAOA, QA, and classical heuristics across many portfolio instances and identify regimes where quantum methods are competitive (but also note hardware scaling and noise limitations). These works form the empirical foundation for hybrid architectures.

2.4 Quantum reinforcement learning (QRL)

Quantum reinforcement learning explores two families: (i) **quantum-assisted RL**, where variational quantum circuits (VQCs) or quantum subroutines serve as function approximators inside classical RL pipelines; and (ii) **fully quantum RL**, which seeks to quantize the MDP or RL algorithms themselves (e.g., amplitude-amplification-based exploration or oracularized environments). Surveys and theoretical works show potential algorithmic advantages (sometimes provable) but caution that most algorithms presuppose future fault-tolerant hardware; nevertheless, NISQ-compatible VQC approaches are being experimentally studied.

3. Problem Formulation

We consider **dynamic multi-asset portfolio management** with (N) tradable assets, discrete rebalancing epochs $(t=0,1,\dots,T)$, and state (s_t)

containing observable market features (prices, returns, technical indicators), current holdings (h_t) , and portfolio cash (c_t) . An RL policy $(\pi_{\theta}(a_t|s_t))$ outputs an action (a_t) which may be either:

1. **Continuous allocation:** target weights $(\tilde{w}_t \in \Delta^N)$ (the probability simplex), or
2. **Discrete trading decision:** a combinatorial selection (e.g., choose up to K assets to buy/sell) and discrete lot sizes.

We explicitly handle **hybrid action spaces** by decomposing actions into a continuous high-level decision from RL and a discrete low-level combinatorial resolution performed by QO:

- RL determines a desired allocation vector (\tilde{w}_t) or a target *profile* (P_t) (e.g., “increase exposure to tech by $x\%$ ”), together with constraints (cardinality K , cash budget, maximum transaction volume).
- A quantum optimizer solves the constrained rounding/selection problem: given continuous targets (\tilde{w}_t) and constraints, find discrete trade orders $(d_t \in \{0,1,\dots,L\}^N)$ that minimize a cost function combining deviation from (\tilde{w}_t) , transaction costs, and risk penalties.

3.1 QUBO encoding for discrete rebalancing

A typical QUBO objective for a discretized rebalancing problem can be written as:

$$\min_{x \in \{0,1\}^{mN}} ; x^{\text{top } Q} x + q^{\text{top } x},$$

where (x) is a binary encoding of discrete decision variables (e.g., one-hot encodings for lot choices, or binary expansion of integer counts), (Q) captures quadratic terms encoding portfolio variance, pairwise asset interactions, and convexified risk objectives, and (q) encodes linear penalties for deviation and transaction costs. Penalty terms enforce budget and

cardinality constraints via large diagonal penalties or parity encodings (parity encodings can reduce qubit overhead for certain constraint patterns). The mapping and penalty-scaling strategy must be tuned carefully to balance objective vs. constraint satisfaction. Practical examples and encoding strategies are described in the literature.

4. Hybrid RL–Quantum Architecture

4.1 High-level design

We propose the following modular architecture:

1. **Market encoder (Feature extractor):** classical deep network (CNN/Transformer/time-series embedding) converts raw market data and microstructure signals into latent state (z_t).
2. **High-level RL agent:** receives (z_t) and outputs a continuous target allocation (\tilde{w}_t) and constraint parameters (\mathcal{C}_t) (e.g., permissible transaction budget, cardinality K). The agent may be an actor–critic (e.g., SAC, PPO) tailored to financial rewards.
3. **Quantum optimizer (QO) module:** takes (\tilde{w}_t) and (\mathcal{C}_t) and solves a QUBO/HUBO instance to produce discrete trades (d_t). Depending on hardware availability, QO can be executed on:
 - quantum annealer (D-Wave), or
 - gate-based QPU with QAOA/VQE, or
 - classical QUBO solver (simulated annealing, tabu search) as a fallback.
4. **Execution & ledger:** apply trades (d_t) to update holdings to (h_{t+1}), track transaction costs, slippage, and update reward.
5. **Learning loop:** RL agent is trained off-policy or on-policy. The quantum module can be invoked during training (costly) or replaced by a classical surrogate for gradient/backpropagation steps; alternately the

QO is called intermittently to update a discrete-action buffer.

Figure placeholders: (Figure 1: system diagram of hybrid RL–QO. Figure 2: QUBO encoding pipeline.)

4.2 Training paradigms

Two training paradigms are viable:

- **Offline pretraining + online hybrid adaptation:** Pretrain RL agent with classical surrogate optimizers (fast QUBO approximations), then switch to quantum module for online selection/testing. This reduces QPU usage during heavy gradient-based training.
- **End-to-end hybrid training:** Incorporate the quantum optimizer in the loop during RL training. Since quantum outputs are nondifferentiable with respect to parameters, gradient estimation uses REINFORCE-style policy gradients or straight-through surrogates; alternatively, differentiate through a differentiable relaxation (continuous relaxations of binary variables) during backprop and use QO only for execution.

We recommend a hybrid approach: RL learns coarse-grained control and risk-sensitivity; QO solves discrete rounding and constraint satisfaction. For experimental reproducibility, both variants should be evaluated.

4.3 Reward engineering

The reward (r_t) must reflect investment objectives:

$$[\begin{aligned} r_t = & \alpha \Delta \log W_t - \beta \cdot \text{TransactionCost}(d_t) - \gamma \cdot \text{RiskPenalty}(h_{t+1}), \end{aligned}]$$

where ($\Delta \log W_t$) is log-wealth change, and (α, β, γ) balance return, cost discipline, and risk control. Risk penalty can be realized as rolling volatility, drawdown (max DD), or CVaR (conditional

value-at-risk) proxies. Reward shaping is crucial to ensure agent conservatism and to avoid overtrading a common pitfall in RL trading.

5. Theoretical Foundations

5.1 MDP and policy optimization foundations

We consider an MDP $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$ and parametric policy (π_{θ}) . Policy optimization seeks (θ^*) maximizing expected discounted return $(J(\theta) = \mathbb{E}[\sum_{t=0}^{\infty} \gamma^t r_t])$. Actor-critic, PPO, SAC, and distributional RL are relevant algorithmic choices depending on continuous/discrete hybrid action spaces.

In hybrid RL-QO, the effective action mapping becomes stochastic due to quantum sampling noise and measurement outcomes; the effective transition kernel (P) must incorporate this stochasticity. The theoretical analysis of convergence in presence of nondifferentiable or nondeterministic low-level solvers follows from stochastic approximation frameworks and off-policy learning stability criteria.

5.2 Quantum optimization complexity and encoding

Mapping an integer-constrained portfolio rebalancing into QUBO introduces a quadratic objective and penalty terms. The resulting QUBO matrices (Q) can be dense (pairwise interaction among assets). For an (N) -asset problem with (b) binary bits per asset (for integer counts), the qubit count scales as $(N \cdot b)$. Parity encodings and problem decomposition strategies can reduce qubit overhead at the cost of increased circuit or annealer embedding complexity (embedding overhead on hardware graphs). Encoding choices and embedding overhead are a major engineering trade-off in practice. Empirical studies show that QAOA and VQE can find high-quality solutions for small-to-moderate instances, and quantum annealers have been used on real market

indices; however, scaling and noise require careful consideration.

5.3 When can quantum methods help?

Theoretical quantum advantage claims hinge on identifying problem instances for which quantum sampling or amplitude amplification yields asymptotic speedups or better approximation ratios. Some QRL and quantum optimization papers prove improved regret bounds or exponential improvements under specific oracular models; practical advantage for real-world finance problems remains conditional on hardware and instance structure. Comprehensive benchmarking is thus critical; recent large-scale benchmark studies compare QAOA/QA to classical heuristics across many real-data portfolio instances and provide nuanced perspectives on advantage regimes.

6. Algorithms and Pseudocode

Below we present pseudocode for a practical hybrid RL-QO training and execution loop.

Algorithm 1: Hybrid RL-Quantum Portfolio Agent (training + deployment)

Inputs: historical data D , initial policy parameters θ_0 , quantum optimizer QO (config),

replay buffer B (optional), training epochs E

for epoch = 1.. E :

for each training episode using D :

reset environment \rightarrow state s_0

for $t = 0..T_{\text{episode}}$:

$z_t = \text{FeatureEncoder}(s_t)$

$(\tilde{w}_t, C_t) = \text{PolicyActor}(\theta, z_t)$ #

continuous target and constraints


```

    QUBO = BuildQUBO(tilde_w_t, C_t)      #
    encode discrete rebalancing

    x_star = QO.solve(QUBO)               # quantum
    or classical solver

    action a_t = Decode(x_star)           # discrete
    trade orders

    s_{t+1}, r_t = EnvironmentStep(s_t, a_t)
    store transition (s_t, a_t, r_t, s_{t+1}) in B

    if training_condition:
        update  $\theta$  via RL update using transitions
        (policy gradient / critic)
    end for
end for

evaluate policy on hold-out validation period
if performance improved:
    optionally checkpoint model and QO
    hyperparams
end for

```

Practical implementation must handle asynchronous quantum calls, latency, and nondifferentiability. For rapid development, runtime QO calls can be simulated with classical QUBO solvers and replaced with real QPU calls during offline evaluation.

7. Experimental Design and Evaluation

7.1 Datasets and experimental protocol

We recommend an evaluation protocol combining:

- **Synthetic market scenarios** with controlled properties (mean-reverting regimes, trending regimes, regime switches) to stress-test learning adaptability.
- **Historical market data:** e.g., S&P 500 constituents, sector indices, ETF baskets, and

alternative assets, sampled at the chosen rebalancing frequency (daily/weekly). Use high-quality data sources (e.g., CRSP, TAQ for microstructure tests) and explicit train/validation/test split with walk-forward evaluation.

- **Transaction cost model:** fixed fees + proportional slippage + market impact emulators to ensure realistic trading frictions.

7.2 Baselines

Compare against:

1. Static mean–variance optimal rebalancing (continuous MVO).
2. Classical combinatorial solver + heuristic rounding (simulated annealing, tabu search).
3. Pure RL policy with standard rounding heuristics (no quantum module).
4. Hybrid RL–QO with classical QUBO solver (ablation to isolate quantum hardware effects).
5. Hybrid RL–QO with QPU (where available).

7.3 Metrics

- **Performance:** cumulative return, annualized return.
- **Risk-adjusted:** annualized volatility, Sharpe ratio, Sortino ratio, maximum drawdown, CVaR at 95%.
- **Operational:** number of transactions, turnover, average transaction cost, latency.
- **Computational:** wall-clock time to solution per rebalancing decision (quantum vs classical), energy/compute cost where measurable.
- **Robustness:** out-of-sample performance across market regimes, sensitivity to reward and penalty coefficients.

8. Representative Numerical Results and Benchmarking (Template & Expected Findings)

Note: This manuscript provides a fully specified experimental protocol and reports expected/representative findings based on literature benchmarks and preliminary studies rather than novel empirical runs on QPU hardware in this draft submission. Reproducible code and datasets will be supplied in the supplementary repository.

Recent benchmark studies indicate that QAOA, QUBO-based annealing, and VQE can achieve high-quality approximations of constrained portfolio optimizations on small-to-moderate instance sizes and that hybrid architectures can match or slightly outperform classical heuristics in specific problem regimes, though scalability and noise remain limiting factors. For example, empirical experiments using IBM devices and QAOA/VQE demonstrate feasibility on tens of assets with coarse discretization, while annealers have been used on S&P500 subsets for cardinality-constrained optimization (Buonaiuto et al., 2023; Phillipson et al., 2021). Large-scale benchmarking across many instances has been pursued to identify regimes where quantum heuristics show comparative strengths.

Expected experimental observations in the hybrid RL-QO paradigm:

- **Quality of discrete rounding:** QO-based rounding often yields closer adherence to continuous targets while respecting discrete constraints compared to naïve rounding or greedy heuristics, resulting in lower tracking error and fewer unnecessary trades.
- **Turnover management:** The QO objective can include explicit turnover penalties, leading to sparser trades and lower transaction costs beneficial in real trading.
- **Latency trade-offs:** For time-sensitive intraday strategies, QPU latency and job-queue times can be prohibitive; hybrid

architectures are therefore best suited for end-of-day or intra-day strategies with relaxed latency constraints.

- **Robustness:** Policies trained with robust reward shaping and realistic costs generalize better; quantum modules must be parameter-tuned to avoid overfitting to small-scale training instances.

These expectations align with literature benchmarking that compares QAOA/VQE/QA to classical heuristics on portfolio problems.

9. Practical Implementation Considerations

9.1 Hardware options and orchestration

- **Quantum annealers (D-Wave):** natural for QUBO; good for proof-of-concept and certain constrained formulations. Embedding and chain strength tuning are nontrivial.
- **Gate-based NISQ QPUs (IBM, Rigetti, IonQ, etc.):** suitable for QAOA/VQE; require careful circuit design and error mitigation. VQE has been experimentally used for portfolio instances on IBM hardware.
- **Classical fallback:** always maintain a high-quality classical QUBO solver for reliability and benchmarking.

Hybrid orchestration must handle job submission, asynchronous returns, and fallback strategies. Latency-aware batching and precomputation of QUBO variants can mitigate real-time latency.

9.2 Engineering challenges

- **Scaling and embedding:** practical qubit counts and hardware topologies limit problem sizes; encoding, decomposition, and problem partitioning are necessary for industrial-scale baskets.
- **Noise and variability:** quantum sampling is stochastic and noisy ensemble strategies and

repeated sampling reduce variance; integrate robust statistical postprocessing.

- **Cost:** QPU access cost vs classical compute cost must be balanced; evaluate computed value added relative to premium access costs.

9.3 Regulatory and operational risk

Deploying RL-driven trading systems subject to fund regulations, best execution obligations, and operational risk frameworks. Quantum-enabled decision-making introduces new auditability questions: how to explain quantum-derived allocations and verify constraint satisfaction. Maintain deterministic fallbacks, rigorous logging, and conservatism tests for live deployment.

10. Robustness, Interpretability, and Risk Management

10.1 Robustness to market regime shifts

Robust RL training should include diverse market regimes and stress scenarios. Model ensembles and conservative policy selection (e.g., percentile-based policy selection) can reduce catastrophic performance under rare shocks.

10.2 Interpretability

Explainable outputs are crucial for practitioner adoption. For the hybrid system:

- Provide **allocation attribution**: break down why the RL agent suggested a target and which quantum constraints changed the discrete outcome (e.g., indicator that “cardinality constraint $K=10$ forced exclusion of asset X due to high covariance with existing positions”).
- Log QUBO solution statistics: energy, constraint violation, sample diversity.
- Use counterfactuals: simulate “what-if” scenarios showing alternate QUBO solutions under varying penalties.

10.3 Model risk and validation

Backtesting must be complemented by model risk controls: statistical significance testing, slippage sensitivity, adversarial stress tests, and in-production shadow trading before full trade execution.

11. Limitations and Future Work

Limitations include:

- **Hardware constraints:** current NISQ devices and annealers limit usable problem sizes, and scaling to large asset universes remains an engineering challenge.
- **Latency:** quantum job latency can preclude low-latency trading strategies.
- **Theoretical guarantees:** provable quantum advantage for general portfolio instances is unresolved; advantage may be instance-specific.

Future work:

- **Hybrid decomposition schemes** that partition a large portfolio into interacting subproblems amenable to separate QPU solves, combined by classical coordination.
- **Differentiable surrogate QUBOs** enabling tighter end-to-end training between RL and QO modules.
- **Hardware-in-the-loop studies** comparing QPU vs classical solvers on institutional-size datasets and cost models.
- **Regulatory frameworks** for auditability and governance of quantum-assisted investment strategies.

12. Conclusions

We introduced a hybrid paradigm integrating reinforcement learning and quantum optimization for dynamic portfolio management. By delegating discrete constrained subproblems to quantum optimizers while using RL for adaptive continuous control, the architecture aims to provide improved discrete

allocation quality under real-world constraints and to open practical pathways for asset managers to experiment with near-term quantum hardware. Benchmarks and surveys suggest the approach is promising for certain problem regimes, but practical adoption depends on hardware progress, embedding strategies, and rigorous risk management processes. We provide an experimental protocol, algorithmic primitives, and a roadmap for future research and deployment.

References

1. Markowitz, H. (1952). Portfolio Selection. *The Journal of Finance*, 7(1), 77–91. ([Wiley Online Library](#))
2. Jiang, Z., Xu, D., & Liang, J. (2017). A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem. *arXiv:1706.10059*.
3. Yang, S. (2023). Deep reinforcement learning for portfolio management. *Elsevier / Expert Systems with Applications* (or conference version); see DOI and indexing.
4. Buonaiuto, G., Gargiulo, F., Pota, M., et al. (2023). Best practices for portfolio optimization by quantum computing, experimented on real quantum devices. *Scientific Reports* (Nature).
5. Fatunmbi, T. O. (2025). Predictive Insurance: Data Science Applications in Risk Profiling and Customer Retention. *International Research Journal of Advanced Engineering and Science*, 10(2), 300–306.
6. Fatunmbi, T. O. (2024). Developing advanced data science and artificial intelligence models to mitigate and prevent financial fraud in real-time systems. *World Journal of Advanced Research and Reviews*, 17(01), 437–456. <https://doi.org/10.30574/wjaets.2024.11.1.0024>
7. Huot, C., Kea, K., Kim, T.-K., & Han, Y. (2024). Enhancing Knapsack-based Financial Portfolio Optimization Using Quantum Approximate Optimization Algorithm. *arXiv:2402.07123*.
8. *Quantum Portfolio Optimization: An Extensive Benchmark*. (2025). arXiv preprint benchmarking QAOA and QA against classical solvers on many portfolio instances.
9. Phillipson, F., et al. (2021). Portfolio Optimisation Using the D-Wave Quantum Annealer. (ICCS paper / proceedings).
10. Blekos, K. (2024). A review on Quantum Approximate Optimization Algorithm (QAOA). *Physics/Computer Science Reviews* (review on QAOA methods and applications).
11. Meyer, N., Ufrecht, C., Periyasamy, M., Scherer, D. D., Plinge, A., & Mutschler, C. (2024). A Survey on Quantum Reinforcement Learning. *arXiv:2211.03464v2*.
12. Chen, S.-Y.-C., et al. (2023). Asynchronous training of quantum reinforcement learning. *Procedia Computer Science / Conference Proceedings*.
13. Samuel, A. J. (2024). Advancements in Cybersecurity: Leveraging AI and Machine Learning for Threat Detection and Prevention. *Journal of Science, Technology and Engineering Research*, 2(3), 64–79. <https://doi.org/10.64206/thwq1548>
14. Samuel, A. J. (2025). Predictive AI for Supply Chain Management: Addressing Vulnerabilities to Cyber-Physical Attacks. *Well Testing Journal*, 34(S2), 185–202.