

Deep Reinforcement Learning for Portfolio Optimization with Options Hedging

IEDA4000F Final Project Report

CHONG Tin Tak
20920359

Department of Industrial Engineering and Decision Analytics
The Hong Kong University of Science and Technology

December 1, 2025

Abstract

This project investigates the application of Deep Reinforcement Learning (DRL) to portfolio optimization with integrated options hedging and systematic risk management. We implement and compare two state-of-the-art algorithms—Deep Deterministic Policy Gradient (DDPG) and Proximal Policy Optimization (PPO)—for continuous portfolio weight allocation across 18 diversified assets spanning eight market sectors.

Our framework incorporates Black-Scholes options pricing for dynamic hedging and a tiered stop-loss mechanism for systematic downside protection. Training on 2010-2018 market data and testing on 2019-2020 (including the COVID-19 market crash), we find that DDPG significantly outperforms PPO with a Sharpe ratio of 5.52 versus 1.85, achieving 219% total return while limiting maximum drawdown to just 8.31% compared to the market's 34% decline.

DDPG's superior performance stems from its off-policy learning capability and deterministic policy output, which proves advantageous for the portfolio allocation task. The agent learned effective hedging strategies, generating \$126,568 in options profits during the test period. These results demonstrate the practical potential of DRL-based portfolio management for navigating both normal market conditions and tail risk events.

Keywords: Deep Reinforcement Learning, Portfolio Optimization, DDPG, PPO, Options Hedging, Risk Management, COVID-19

Contents

1	Introduction	8
1.1	Background and Motivation	8
1.2	Research Objectives	8
1.3	Key Contributions	9
1.4	Report Structure	9
2	Literature Review	10
2.1	Classical Portfolio Optimization	10
2.2	Reinforcement Learning Foundations	10
2.3	Deep Reinforcement Learning for Finance	11
2.4	Actor-Critic Methods	11
2.4.1	Deep Deterministic Policy Gradient (DDPG)	11
2.4.2	Proximal Policy Optimization (PPO)	12
2.5	Options in Portfolio Management	12
2.6	Risk Management in Algorithmic Trading	13
2.7	Gap in Literature	13
3	Methodology	13
3.1	Problem Formulation	13
3.1.1	State Space	14
3.1.2	Action Space	14
3.1.3	Reward Function	14
3.2	Deep Deterministic Policy Gradient (DDPG)	15
3.2.1	Actor Network	15
3.2.2	Critic Network	15
3.2.3	Training Algorithm	15
3.2.4	Exploration	16
3.3	Proximal Policy Optimization (PPO)	16
3.3.1	Policy Network	16
3.3.2	Value Network	16
3.3.3	Clipped Surrogate Objective	16
3.3.4	Generalized Advantage Estimation (GAE)	16
3.3.5	Complete Objective	17
3.4	Options Pricing and Hedging	17
3.4.1	Black-Scholes Model	17
3.4.2	Protective Put Strategy	17
3.4.3	Covered Call Strategy	18
3.4.4	Combined Options Overlay	18
3.5	Stop-Loss Mechanism	18
4	Implementation	19
4.1	System Architecture	19
4.2	Code Organization	19
4.3	Data Loading and Preprocessing	19
4.4	Portfolio Environment	20
4.4.1	State Construction	20
4.4.2	Reward Calculation	21

4.5	Agent Implementation	21
4.5.1	DDPG Agent	21
4.5.2	PPO Agent	22
4.6	Options Pricing Module	22
4.7	Metrics Calculation	23
4.8	Training Pipeline	23
4.9	Evaluation Framework	24
4.10	Visualization Tools	25
5	Experimental Setup	25
5.1	Asset Universe	25
5.2	Data Period and Split	25
5.3	Why No Validation Set?	27
5.4	Hyperparameter Configuration	27
5.4.1	DDPG Hyperparameters	29
5.4.2	PPO Hyperparameters	29
5.4.3	Environment Configuration	30
5.4.4	Agent-Specific Risk Parameters	30
5.4.5	Feature Engineering Configuration	31
5.5	Benchmark Strategies	31
5.6	Evaluation Metrics	31
5.7	Computational Environment	32
5.8	Reproducibility	34
6	Results	34
6.1	Overall Performance Summary	34
6.2	Comparison with Benchmark Strategies	35
6.3	Portfolio Value Evolution	35
6.4	Drawdown Analysis	36
6.5	COVID-19 Crash Performance	37
6.6	Options Hedging Analysis	37
6.7	Portfolio Weight Allocation	39
6.8	Risk-Adjusted Performance Comparison	39
6.9	Monthly Performance Attribution	39
6.10	Statistical Significance	41
7	Discussion	41
7.1	Why DDPG Outperforms PPO	41
7.1.1	Off-Policy Learning Advantage	41
7.1.2	Deterministic Policy Benefits	42
7.1.3	Q-Value Learning	42
7.2	PPO's Limitations in This Context	42
7.2.1	On-Policy Data Efficiency	42
7.2.2	Clipping Limitations	42
7.2.3	Stochastic Policy Drawbacks	43
7.3	Options Hedging Insights	43
7.3.1	DDPG's Hedging Strategy	43
7.3.2	PPO's Conservative Approach	43
7.4	Stop-Loss Mechanism Effectiveness	43

7.5	Limitations and Considerations	44
7.5.1	Data Limitations	44
7.5.2	Model Limitations	44
7.5.3	Options Model Limitations	44
7.6	Practical Implications	44
7.6.1	Algorithm Selection	44
7.6.2	Risk Management	45
7.6.3	Implementation Considerations	45
7.7	Asset Correlation Analysis	45
7.8	Comparison with Literature	46
8	Conclusion	46
8.1	Summary of Findings	46
8.2	Contributions	47
8.3	Limitations	47
8.4	Future Work	47
8.4.1	Algorithm Enhancements	47
8.4.2	Risk Management Extensions	48
8.4.3	Practical Extensions	48
8.5	Final Remarks	48
A	Mathematical Formulas Reference	50
A.1	Return Calculations	50
A.1.1	Simple Return	50
A.1.2	Logarithmic Return	50
A.1.3	Portfolio Return	50
A.1.4	Cumulative Return	50
A.1.5	Annualized Return	50
A.2	Risk Metrics	50
A.2.1	Volatility (Standard Deviation)	50
A.2.2	Annualized Volatility	50
A.2.3	Drawdown	50
A.2.4	Maximum Drawdown	51
A.2.5	Downside Deviation	51
A.3	Performance Ratios	51
A.3.1	Sharpe Ratio	51
A.3.2	Sortino Ratio	51
A.3.3	Calmar Ratio	51
A.3.4	Information Ratio	51
A.4	Deep Deterministic Policy Gradient (DDPG)	51
A.4.1	Critic Loss (TD Error)	51
A.4.2	Target Value	51
A.4.3	Policy Gradient	51
A.4.4	Soft Target Update	52
A.4.5	Ornstein-Uhlenbeck Noise	52
A.5	Proximal Policy Optimization (PPO)	52
A.5.1	Clipped Surrogate Objective	52
A.5.2	Probability Ratio	52

A.5.3	Generalized Advantage Estimation (GAE)	52
A.5.4	TD Residual	52
A.5.5	Value Function Loss	52
A.5.6	Entropy Bonus	52
A.5.7	Complete PPO Objective	52
A.6	Options Pricing (Black-Scholes)	52
A.6.1	Call Option Price	52
A.6.2	Put Option Price	52
A.6.3	d1 and d2 Parameters	53
A.6.4	Put-Call Parity	53
A.6.5	Option Greeks	53
A.7	Payoff Functions	53
A.7.1	Call Option Payoff	53
A.7.2	Put Option Payoff	53
A.7.3	Protective Put Payoff	53
A.7.4	Covered Call Payoff	53
A.7.5	Collar Strategy Payoff	53
A.7.6	Total Options P&L	54
A.8	Stop-Loss Mechanism	54
A.8.1	Tiered Exposure Adjustment	54
A.8.2	Adjusted Portfolio Weights	54
A.9	Reward Function	54
A.9.1	Risk-Adjusted Reward	54
A.9.2	Turnover	54
A.9.3	Transaction Costs	54
B	Configuration Parameters	54
B.1	YAML Configuration File	54
B.2	Network Architecture Details	57
B.2.1	Actor Network (DDPG)	57
B.2.2	Critic Network (DDPG)	57
B.2.3	Policy Network (PPO)	57
B.3	State Space Specification	58
B.4	Action Space Specification	58
B.5	Hardware and Software Specifications	58
B.6	Training Time and Resources	59
B.7	Reproducibility Checklist	59
B.8	Data Preprocessing Steps	59

List of Figures

1	High-level system architecture showing data flow between components. . .	19
2	Training curves showing episode reward progression over 100,000 timesteps. (a) DDPG converges faster due to off-policy learning and experience replay. (b) PPO shows more gradual improvement with higher variance typical of on-policy methods.	33
3	Portfolio value evolution during test period (2019-2020). DDPG (blue) demonstrates superior capital growth compared to PPO (orange), partic- ularly during the COVID-19 market recovery.	36
4	Drawdown comparison showing DDPG's superior downside protection. The shaded regions indicate drawdown magnitude over time.	36
5	Options hedging analysis comparing DDPG and PPO: (a) average options strategy usage showing DDPG's aggressive hedging approach, (b) total options P&L comparison, (c) cumulative options profit over time with COVID-19 crash period highlighted, and (d) detailed hedging efficiency metrics.	38
6	Portfolio weight allocation analysis: (a) average sector allocation compar- ison, (b) DDPG sector allocation over time showing defensive shift during COVID-19, (c) DDPG top 5 holdings, (d) weight distribution showing DDPG's more concentrated approach vs PPO's diversification.	39
7	Comprehensive performance comparison showing portfolio evolution, draw- downs, and metrics summary.	40
8	Asset correlation matrix computed from training period data (2010-2018). Technology stocks show high intra-sector correlation, while bonds (TLT, AGG) provide diversification benefits through low/negative correlation with equities.	45

List of Tables

1	Project module structure and responsibilities	19
2	Asset universe with 18 diversified instruments	26
3	Data period configuration	26
4	DDPG hyperparameter configuration	28
5	PPO hyperparameter configuration	29
6	Environment configuration parameters	30
7	Feature engineering parameters	31
8	Performance evaluation metrics	32
9	Reproducibility specifications	33
10	Performance comparison: DDPG vs PPO on test period (2019-2020) . . .	34
11	Performance comparison: DRL agents vs benchmark strategies (2019-2020)	35
12	Drawdown statistics comparison	37
13	Performance during COVID-19 crash (February 19 - March 23, 2020) . .	37
14	Options hedging statistics	38
15	Risk-adjusted metrics comparison	40
16	Monthly returns (%) during test period	40
17	Statistical significance tests	41
18	Stop-loss trigger statistics	44

19	DDPG Actor Network Architecture	57
20	DDPG Critic Network Architecture	57
21	PPO Policy Network Architecture	57
22	State Space Components	58
23	Action Space Components	58
24	Computational Environment	58
25	Training Resource Requirements	59

1 Introduction

1.1 Background and Motivation

Portfolio optimization has been a cornerstone of modern finance since Harry Markowitz's seminal work on mean-variance optimization in 1952 [1]. Traditional approaches rely on statistical estimates of expected returns and covariances, which often prove unstable in practice due to estimation errors that compound over time. These classical methods also struggle to adapt to non-stationary market dynamics and fail to capture complex nonlinear relationships between assets.

The financial markets of the 21st century present unique challenges that expose the limitations of traditional portfolio management approaches:

- **Market Complexity:** Modern markets exhibit intricate dependencies, regime changes, and fat-tailed return distributions that violate the Gaussian assumptions underlying classical models.
- **High-Frequency Dynamics:** Rapid information dissemination and algorithmic trading create fast-moving market conditions that require adaptive strategies.
- **Tail Risk Events:** Events like the 2008 financial crisis and the COVID-19 market crash of 2020 demonstrate the importance of robust risk management beyond traditional volatility measures.
- **Transaction Costs and Constraints:** Real-world portfolio management must account for transaction costs, position limits, and regulatory constraints that classical models often ignore.

Deep Reinforcement Learning (DRL) offers a promising alternative framework for portfolio optimization. By framing portfolio management as a sequential decision-making problem, DRL agents can learn adaptive strategies directly from market data without relying on explicit statistical models. Recent advances in deep learning provide the representational capacity to capture complex market patterns, while reinforcement learning algorithms enable optimization of long-term risk-adjusted returns.

1.2 Research Objectives

This project investigates the application of Deep Reinforcement Learning to portfolio optimization with the following primary objectives:

1. **Algorithm Comparison:** Implement and compare two state-of-the-art DRL algorithms—Deep Deterministic Policy Gradient (DDPG) and Proximal Policy Optimization (PPO)—for continuous portfolio weight allocation.
2. **Options Integration:** Develop a novel framework for incorporating options-based hedging strategies within the DRL portfolio optimization environment, enabling dynamic downside protection.
3. **Risk Management:** Design and evaluate a multi-tiered stop-loss system that adapts portfolio exposure based on drawdown levels, providing systematic risk control.

4. **Stress Testing:** Evaluate the trained agents' performance during extreme market conditions, specifically the COVID-19 market crash of March 2020, to assess robustness and crisis-alpha generation.
5. **Reproducibility:** Create a comprehensive, well-documented codebase that enables reproduction and extension of our results.

1.3 Key Contributions

This work makes the following contributions to the field of algorithmic portfolio management:

1. **Integrated Options Hedging:** We develop a novel portfolio environment that integrates Black-Scholes options pricing with DRL-based portfolio optimization, allowing agents to learn when and how much to hedge using put options.
2. **Adaptive Stop-Loss Mechanism:** We implement a tiered stop-loss system that progressively reduces portfolio exposure as drawdowns deepen, providing systematic downside protection while allowing participation in market recoveries.
3. **Comprehensive Benchmark Study:** We conduct extensive experiments comparing DDPG and PPO across multiple performance dimensions, including risk-adjusted returns, drawdown characteristics, and behavior during market stress.
4. **COVID-19 Stress Test:** We specifically evaluate model performance during the March 2020 market crash, demonstrating the practical value of DRL-based portfolio management during tail risk events.
5. **Open-Source Implementation:** We provide a complete, modular implementation suitable for research and practical applications, including data loading, environment simulation, agent training, and performance visualization.

1.4 Report Structure

The remainder of this report is organized as follows:

- **Section 2:** Reviews related work in portfolio optimization, reinforcement learning for finance, and options-based hedging strategies.
- **Section 3:** Presents the mathematical framework, including the MDP formulation, reward function design, and the DDPG and PPO algorithms.
- **Section 4:** Describes the system architecture, code structure, and implementation details.
- **Section 5:** Details the experimental setup, including asset selection, data preprocessing, and hyperparameter configurations.
- **Section 6:** Presents comprehensive results comparing DDPG and PPO performance across multiple metrics.

- **Section 7:** Analyzes the results, discusses the strengths and limitations of each approach, and provides insights into the learned strategies.
- **Section 8:** Summarizes our findings and outlines directions for future research.

Appendices provide detailed mathematical formulas (Appendix A) and configuration parameters (Appendix B).

2 Literature Review

2.1 Classical Portfolio Optimization

The foundation of modern portfolio theory was established by Harry Markowitz [1], who formulated the mean-variance optimization problem. Given n assets with expected returns $\boldsymbol{\mu} \in \mathbb{R}^n$ and covariance matrix $\boldsymbol{\Sigma} \in \mathbb{R}^{n \times n}$, the investor seeks portfolio weights $\mathbf{w} \in \mathbb{R}^n$ that maximize expected return for a given level of risk:

$$\max_{\mathbf{w}} \quad \mathbf{w}^\top \boldsymbol{\mu} - \frac{\lambda}{2} \mathbf{w}^\top \boldsymbol{\Sigma} \mathbf{w} \quad \text{s.t.} \quad \mathbf{w}^\top \mathbf{1} = 1, \quad \mathbf{w} \geq 0 \quad (1)$$

where λ is the risk aversion parameter. Despite its theoretical elegance, mean-variance optimization suffers from several practical limitations:

- **Estimation Sensitivity:** Small errors in $\boldsymbol{\mu}$ and $\boldsymbol{\Sigma}$ estimates lead to dramatically different optimal portfolios [4].
- **Static Nature:** The single-period formulation ignores the dynamic nature of portfolio rebalancing.
- **Distributional Assumptions:** Gaussian returns assumption fails to capture fat tails and asymmetric distributions observed in financial markets.

Extensions such as Black-Litterman [3] and robust optimization [5] address some of these limitations but remain fundamentally constrained by their reliance on statistical estimation.

2.2 Reinforcement Learning Foundations

Reinforcement learning (RL) provides a framework for sequential decision-making under uncertainty [6]. An RL agent interacts with an environment modeled as a Markov Decision Process (MDP) defined by the tuple $(\mathcal{S}, \mathcal{A}, P, R, \gamma)$:

- \mathcal{S} : State space (market observations)
- \mathcal{A} : Action space (portfolio weights)
- $P(s'|s, a)$: Transition dynamics
- $R(s, a, s')$: Reward function
- $\gamma \in [0, 1]$: Discount factor

The goal is to find a policy $\pi : \mathcal{S} \rightarrow \mathcal{A}$ that maximizes expected cumulative discounted rewards:

$$J(\pi) = \mathbb{E}_{\pi} \left[\sum_{t=0}^{\infty} \gamma^t R_t \right] \quad (2)$$

2.3 Deep Reinforcement Learning for Finance

The application of deep reinforcement learning to financial problems has gained significant momentum in recent years. Several key works have demonstrated the potential of DRL for portfolio management:

Jiang et al. (2017) [10] proposed a deep learning framework for portfolio management using a convolutional neural network to extract features from historical price data. Their approach achieved competitive performance against traditional benchmarks.

Liang et al. (2018) [11] applied DDPG to portfolio optimization, demonstrating superior performance compared to traditional methods on Chinese stock market data.

Yang et al. (2020) [12] developed FinRL, an open-source library for financial reinforcement learning that provides standardized implementations of popular DRL algorithms.

Liu et al. (2021) [13] proposed an ensemble method combining multiple DRL agents for more robust portfolio decisions.

2.4 Actor-Critic Methods

Actor-critic methods combine value-based and policy-based approaches, using a critic to estimate value functions and an actor to optimize the policy. This architecture offers several advantages:

- **Reduced Variance:** The critic's value estimates provide a baseline for variance reduction in policy gradient updates.
- **Continuous Actions:** Actor networks can directly output continuous actions, suitable for portfolio weight allocation.
- **Sample Efficiency:** Combining on-policy and off-policy learning improves data utilization.

2.4.1 Deep Deterministic Policy Gradient (DDPG)

DDPG [7] extends the deterministic policy gradient theorem to deep neural networks. Key features include:

- **Deterministic Policy:** The actor outputs a deterministic action $a = \mu_{\theta}(s)$, with exploration noise added during training.
- **Experience Replay:** Transitions are stored in a replay buffer and sampled randomly for training, breaking temporal correlations.
- **Target Networks:** Slowly-updated target networks stabilize training by providing consistent targets for Q-value estimation.

- **Off-Policy Learning:** DDPG can learn from past experiences, improving sample efficiency.

2.4.2 Proximal Policy Optimization (PPO)

PPO [8] addresses the challenge of stable policy updates in on-policy methods. Key innovations include:

- **Clipped Objective:** The surrogate objective is clipped to prevent excessively large policy updates.
- **Trust Region:** The clipping mechanism implicitly enforces a trust region constraint on policy changes.
- **Sample Efficiency:** Multiple epochs of optimization on the same batch of data improve learning efficiency.
- **Simplicity:** PPO achieves competitive performance with simpler implementation compared to methods like TRPO.

2.5 Options in Portfolio Management

Options provide non-linear payoff profiles that can be used for hedging and speculation. The Black-Scholes model [2] provides the foundational framework for options pricing:

$$C = S_0 N(d_1) - K e^{-rT} N(d_2) \quad (3)$$

where:

- $d_1 = \frac{\ln(S_0/K) + (r + \sigma^2/2)T}{\sigma\sqrt{T}}$
- $d_2 = d_1 - \sigma\sqrt{T}$
- $N(\cdot)$ is the cumulative standard normal distribution

Protective puts represent a fundamental hedging strategy where an investor holding a long position purchases put options to limit downside risk. The payoff at expiration is:

$$\text{Payoff}_{\text{put}} = \max(S_T, K_{\text{put}}) - P_0 \quad (4)$$

where P_0 is the premium paid for the put option. This creates an asymmetric payoff profile that preserves upside potential while limiting losses.

Covered calls complement protective puts by generating premium income. An investor holding a long position writes (sells) call options, receiving premium in exchange for capping upside potential:

$$\text{Payoff}_{\text{call}} = \min(S_T, K_{\text{call}}) + C_0 \quad (5)$$

where C_0 is the premium received. When combined, protective puts and covered calls form a *collar strategy*, which bounds the portfolio's payoff within a defined range while potentially achieving near-zero net premium cost.

2.6 Risk Management in Algorithmic Trading

Effective risk management is crucial for algorithmic trading systems. Common approaches include:

- **Position Sizing:** Kelly criterion and fractional Kelly methods optimize position sizes based on edge and variance.
- **Stop-Loss Orders:** Automatic position liquidation when losses exceed predetermined thresholds.
- **Portfolio Constraints:** Limits on sector exposure, single-stock concentration, and leverage.
- **Value-at-Risk (VaR):** Statistical measures of potential losses at given confidence levels.

Our work integrates multiple risk management approaches within the DRL framework, including options-based hedging and tiered stop-loss mechanisms.

2.7 Gap in Literature

While significant progress has been made in applying DRL to portfolio optimization, several gaps remain:

1. **Options Integration:** Most DRL portfolio management studies focus on equity allocation without incorporating derivative instruments for hedging.
2. **Systematic Risk Management:** Few studies integrate explicit risk management mechanisms within the DRL framework.
3. **Stress Testing:** Limited evaluation of DRL agents during extreme market events like the COVID-19 crash.
4. **Algorithm Comparison:** Comprehensive comparisons between DDPG and PPO for portfolio optimization under consistent experimental conditions are scarce.

This work addresses these gaps by developing an integrated framework that combines DRL-based portfolio optimization with options hedging and systematic risk management, evaluated rigorously during both normal and crisis market conditions.

3 Methodology

3.1 Problem Formulation

We formulate portfolio optimization as a Markov Decision Process (MDP), enabling the application of reinforcement learning algorithms. At each time step t , the agent observes market conditions, selects portfolio weights, and receives a reward based on portfolio performance.

3.1.1 State Space

The state $s_t \in \mathcal{S}$ captures relevant market information at time t :

$$s_t = [\mathbf{r}_{t-L:t}^\top \quad \text{vol}_{t-L:t}^\top \quad \mathbf{w}_{t-1}^\top \quad \text{features}_t^\top] \quad (6)$$

where:

- $\mathbf{r}_{t-L:t} \in \mathbb{R}^{n \times L}$: Historical returns for n assets over lookback window L
- $\text{vol}_{t-L:t} \in \mathbb{R}^{n \times L}$: Rolling volatility estimates
- $\mathbf{w}_{t-1} \in \mathbb{R}^n$: Current portfolio weights
- features_t : Additional technical indicators (momentum, RSI, etc.)

The state representation enables the agent to learn patterns from historical price movements while maintaining awareness of current portfolio positioning.

For our configuration with $n = 18$ assets and lookback window $L = 60$ days, the state vector dimension is:

$$d_s = n \times L \times 2 + n + n_{\text{features}} = 18 \times 60 \times 2 + 18 + 36 = 2214 \quad (7)$$

where the factor of 2 accounts for both returns and volatility features, and $n_{\text{features}} = 36$ represents additional technical indicators (SMA, EMA, momentum at multiple periods).

3.1.2 Action Space

The action $\mathbf{a}_t \in \mathcal{A}$ represents the target portfolio allocation:

$$\mathbf{a}_t = [w_1^t \quad w_2^t \quad \cdots \quad w_n^t \quad h_t] \quad (8)$$

Subject to constraints:

$$\sum_{i=1}^n w_i^t \leq 1 \quad (\text{allocation constraint}) \quad (9)$$

$$w_i^t \geq 0 \quad \forall i \quad (\text{long-only constraint}) \quad (10)$$

$$h_t \in [0, h_{\max}] \quad (\text{hedge ratio constraint}) \quad (11)$$

The hedge ratio h_t determines the fraction of portfolio value allocated to protective put options. Any unallocated capital $(1 - \sum_i w_i^t)$ earns the risk-free rate.

3.1.3 Reward Function

The reward function balances return maximization with risk management. We use a risk-adjusted reward:

$$r_t = R_t^{\text{port}} - \lambda_{\text{risk}} \cdot \text{Risk}_t - \lambda_{\text{tc}} \cdot \text{TC}_t \quad (12)$$

where:

Portfolio Return:

$$R_t^{\text{port}} = \sum_{i=1}^n w_i^{t-1} \cdot r_i^t + (1 - \sum_i w_i^{t-1}) \cdot r_f + \Pi_t^{\text{options}} \quad (13)$$

Risk Penalty:

$$\text{Risk}_t = \max(0, -R_t^{\text{port}})^2 \quad (14)$$

Transaction Costs:

$$\text{TC}_t = c \cdot \sum_{i=1}^n |w_i^t - w_i^{t-1}| \quad (15)$$

The squared downside penalty encourages the agent to avoid large negative returns, while the transaction cost term discourages excessive trading.

3.2 Deep Deterministic Policy Gradient (DDPG)

DDPG is an off-policy actor-critic algorithm designed for continuous action spaces. We employ DDPG with the following components:

3.2.1 Actor Network

The actor network $\mu_\theta : \mathcal{S} \rightarrow \mathcal{A}$ maps states to deterministic actions:

$$\mathbf{a}_t = \mu_\theta(s_t) \quad (16)$$

Architecture: State \rightarrow FC(256) \rightarrow ReLU \rightarrow FC(256) \rightarrow ReLU \rightarrow FC($n+1$) \rightarrow Softmax

The softmax output ensures valid portfolio weights that sum to at most 1.

3.2.2 Critic Network

The critic network $Q_\phi : \mathcal{S} \times \mathcal{A} \rightarrow \mathbb{R}$ estimates the Q-value:

$$Q_\phi(s_t, \mathbf{a}_t) \approx \mathbb{E} \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t, \mathbf{a}_t \right] \quad (17)$$

Architecture: State and action are concatenated after initial state processing:

$$\text{State} \rightarrow \text{FC}(256) \rightarrow \text{ReLU} \rightarrow [\cdot, \mathbf{a}] \rightarrow \text{FC}(256) \rightarrow \text{ReLU} \rightarrow \text{FC}(1)$$

3.2.3 Training Algorithm

DDPG training alternates between:

Critic Update: Minimize temporal difference error:

$$L(\phi) = \mathbb{E}_{(s,a,r,s') \sim \mathcal{D}} [(Q_\phi(s, a) - y)^2] \quad (18)$$

where the target is:

$$y = r + \gamma Q_{\phi'}(s', \mu_{\theta'}(s')) \quad (19)$$

Actor Update: Maximize expected Q-value via deterministic policy gradient:

$$\nabla_\theta J = \mathbb{E}_{s \sim \mathcal{D}} [\nabla_a Q_\phi(s, a)|_{a=\mu_\theta(s)} \cdot \nabla_\theta \mu_\theta(s)] \quad (20)$$

Target Network Update: Soft update with parameter τ :

$$\theta' \leftarrow \tau \theta + (1 - \tau) \theta' \quad (21)$$

$$\phi' \leftarrow \tau \phi + (1 - \tau) \phi' \quad (22)$$

3.2.4 Exploration

Exploration is achieved by adding Ornstein-Uhlenbeck noise to actions:

$$\mathbf{a}_t = \mu_\theta(s_t) + \mathcal{N}_t \quad (23)$$

where \mathcal{N}_t follows:

$$d\mathcal{N}_t = \theta_{OU}(\mu_{OU} - \mathcal{N}_t)dt + \sigma_{OU}dW_t \quad (24)$$

3.3 Proximal Policy Optimization (PPO)

PPO is an on-policy actor-critic algorithm that achieves stable policy updates through a clipped surrogate objective.

3.3.1 Policy Network

The policy network outputs a Gaussian distribution over actions:

$$\pi_\theta(\mathbf{a}|s) = \mathcal{N}(\mu_\theta(s), \sigma_\theta(s)) \quad (25)$$

For portfolio weights, actions are sampled and then passed through a softmax transformation to ensure valid allocations.

3.3.2 Value Network

The value network $V_\psi : \mathcal{S} \rightarrow \mathbb{R}$ estimates state values:

$$V_\psi(s_t) \approx \mathbb{E}_\pi \left[\sum_{k=0}^{\infty} \gamma^k r_{t+k} \mid s_t \right] \quad (26)$$

3.3.3 Clipped Surrogate Objective

PPO optimizes a clipped surrogate objective:

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta) \hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \right) \right] \quad (27)$$

where the probability ratio is:

$$r_t(\theta) = \frac{\pi_\theta(\mathbf{a}_t|s_t)}{\pi_{\theta_{\text{old}}}(\mathbf{a}_t|s_t)} \quad (28)$$

3.3.4 Generalized Advantage Estimation (GAE)

Advantages are estimated using GAE [9]:

$$\hat{A}_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l} \quad (29)$$

where the TD residual is:

$$\delta_t = r_t + \gamma V_\psi(s_{t+1}) - V_\psi(s_t) \quad (30)$$

3.3.5 Complete Objective

The full PPO objective combines policy, value, and entropy terms:

$$L(\theta, \psi) = L^{\text{CLIP}}(\theta) - c_1 L^{VF}(\psi) + c_2 S[\pi_\theta] \quad (31)$$

where:

- $L^{VF}(\psi) = \mathbb{E}_t[(V_\psi(s_t) - V_t^{\text{target}})^2]$ is the value function loss
- $S[\pi_\theta] = \mathbb{E}_t[-\log \pi_\theta(\mathbf{a}_t|s_t)]$ is the entropy bonus for exploration

3.4 Options Pricing and Hedging

3.4.1 Black-Scholes Model

We use the Black-Scholes model for options pricing. For a European put option:

$$P = Ke^{-rT}N(-d_2) - S_0N(-d_1) \quad (32)$$

$$d_1 = \frac{\ln(S_0/K) + (r + \sigma^2/2)T}{\sigma\sqrt{T}} \quad (33)$$

$$d_2 = d_1 - \sigma\sqrt{T} \quad (34)$$

where:

- S_0 : Current asset price
- K : Strike price
- r : Risk-free interest rate
- T : Time to expiration
- σ : Implied volatility
- $N(\cdot)$: Cumulative standard normal distribution

3.4.2 Protective Put Strategy

The agent can allocate a hedge ratio h_t^{put} to protective puts on the portfolio. The put payoff at expiration is:

$$\Pi_t^{\text{put}} = h_t^{\text{put}} \cdot V_t \cdot \max\left(0, \frac{K_{\text{put}} - S_T}{S_0}\right) - P_0 \quad (35)$$

where $K_{\text{put}} = 0.96 \cdot S_0$ (4% out-of-the-money). This provides portfolio insurance: when the portfolio declines below the strike price, the put option compensates for losses.

3.4.3 Covered Call Strategy

In addition to protective puts, the agent can write covered calls on the portfolio to generate premium income. The covered call payoff is:

$$\Pi_t^{\text{call}} = C_0 - h_t^{\text{call}} \cdot V_t \cdot \max\left(0, \frac{S_T - K_{\text{call}}}{S_0}\right) \quad (36)$$

where $K_{\text{call}} = 1.04 \cdot S_0$ (4% out-of-the-money) and C_0 is the premium received. The covered call strategy:

- Generates income through option premiums during low-volatility periods
- Caps upside potential when portfolio value exceeds the strike price
- Provides a buffer against small losses through premium income
- Works synergistically with protective puts in a “collar” configuration

3.4.4 Combined Options Overlay

The total options P&L combines both strategies:

$$\Pi_t^{\text{options}} = \Pi_t^{\text{put}} + \Pi_t^{\text{call}} \quad (37)$$

The agent learns to dynamically adjust hedge ratios h_t^{put} and h_t^{call} based on market conditions, with maximum ratios capped at 25% of portfolio value for each strategy.

3.5 Stop-Loss Mechanism

We implement a tiered stop-loss system that reduces portfolio exposure as drawdowns deepen:

$$\text{Exposure Multiplier} = \begin{cases} 1.0 & \text{if } DD_t < 5\% \\ 0.75 & \text{if } 5\% \leq DD_t < 10\% \\ 0.50 & \text{if } 10\% \leq DD_t < 15\% \\ 0.25 & \text{if } DD_t \geq 15\% \end{cases} \quad (38)$$

where the drawdown is:

$$DD_t = \frac{V_t^{\text{peak}} - V_t}{V_t^{\text{peak}}} \quad (39)$$

This mechanism provides systematic risk control by:

1. Allowing full participation during normal market conditions
2. Progressively reducing exposure as losses accumulate
3. Preserving capital during severe drawdowns
4. Enabling participation in recovery through maintained (reduced) exposure

4 Implementation

4.1 System Architecture

The system follows a modular architecture designed for flexibility and extensibility. Figure 1 illustrates the high-level component interactions.

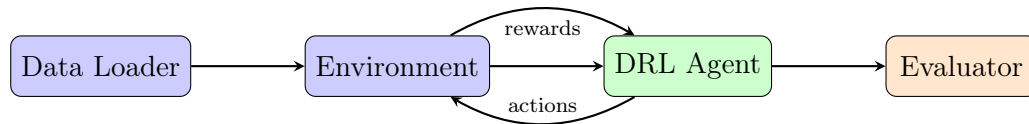


Figure 1: High-level system architecture showing data flow between components.

4.2 Code Organization

The codebase is organized into the following modules:

Table 1: Project module structure and responsibilities

Module	Responsibility
src/data_loader.py	Data loading and preprocessing
src/portfolio_env.py	Base portfolio environment (Gym)
src/portfolio_env_with_options.py	Extended environment with options
src/agents.py	DDPG and PPO agent implementations
src/options_pricing.py	Black-Scholes pricing functions
src/metrics.py	Performance metric calculations
src/benchmarks.py	Benchmark strategy implementations
src/visualization.py	Plotting and visualization utilities

4.3 Data Loading and Preprocessing

The DataLoader class handles data acquisition and preprocessing:

```

class DataLoader:
    def __init__(self, tickers, start_date, end_date):
        self.tickers = tickers
        self.start_date = start_date
        self.end_date = end_date

    def load_data(self):
        # Fetch adjusted close prices
        # Calculate returns
        # Handle missing data
        return prices, returns
  
```

Key preprocessing steps include:

1. Fetching adjusted close prices from Yahoo Finance
2. Computing logarithmic returns: $r_t = \ln(P_t/P_{t-1})$
3. Forward-filling missing values
4. Calculating rolling statistics (volatility, correlations)
5. Normalizing features to zero mean and unit variance

4.4 Portfolio Environment

The portfolio environment extends OpenAI Gym's interface:

```
class PortfolioEnvWithOptions(gym.Env):
    def __init__(self, prices, config):
        self.action_space = spaces.Box(
            low=0, high=1,
            shape=(n_assets + 1,) # +1 for hedge ratio
        )
        self.observation_space = spaces.Box(
            low=-np.inf, high=np.inf,
            shape=(state_dim,)
        )

    def step(self, action):
        # Process action (allocate weights)
        # Calculate portfolio return
        # Apply options hedging
        # Check stop-loss conditions
        # Return observation, reward, done, info

    def reset(self):
        # Reset to initial state
        return initial_observation
```

4.4.1 State Construction

The state vector is constructed as:

```
def _get_state(self):
    # Historical returns (flattened)
    returns_history = self.returns[
        self.current_step - self.lookback:self.current_step
    ].flatten()

    # Current portfolio weights
    current_weights = self.weights

    # Technical indicators
```

```

    volatility = self.rolling_vol[self.current_step]

    return np.concatenate([
        returns_history,
        current_weights,
        volatility
    ])

```

4.4.2 Reward Calculation

The reward function implementation:

```

def _calculate_reward(self, portfolio_return):
    # Base reward: portfolio return
    reward = portfolio_return

    # Risk penalty for negative returns
    if portfolio_return < 0:
        reward -= self.risk_penalty * (portfolio_return ** 2)

    # Transaction cost penalty
    turnover = np.sum(np.abs(
        self.weights - self.prev_weights
    ))
    reward -= self.transaction_cost * turnover

    return reward

```

4.5 Agent Implementation

4.5.1 DDPG Agent

The DDPG agent uses Stable-Baselines3's implementation with custom hyperparameters:

```

from stable_baselines3 import DDPG
from stable_baselines3.common.noise import OrnsteinUhlenbeckActionNoise

# Initialize action noise
n_actions = env.action_space.shape[-1]
action_noise = OrnsteinUhlenbeckActionNoise(
    mean=np.zeros(n_actions),
    sigma=0.1 * np.ones(n_actions)
)

# Create DDPG agent (final benchmark configuration)
agent = DDPG(
    "MlpPolicy",
    env,
    learning_rate=1e-4,
    buffer_size=500000,

```

```

        learning_starts=10000,
        batch_size=256,
        tau=0.01,
        gamma=0.99,
        action_noise=action_noise,
        policy_kwargs=dict(
            net_arch=dict(pi=[512, 512, 256, 128], qf=[512, 512, 256, 128])
        ),
        verbose=1
    )

```

4.5.2 PPO Agent

The PPO agent configuration:

```

from stable_baselines3 import PPO

agent = PPO(
    "MlpPolicy",
    env,
    learning_rate=5e-5,
    n_steps=2048,
    batch_size=128,
    n_epochs=10,
    gamma=0.99,
    gae_lambda=0.95,
    clip_range=0.2,
    ent_coef=0.05,
    policy_kwargs=dict(net_arch=[512, 512, 256, 128]),
    verbose=1
)

```

4.6 Options Pricing Module

The options pricing module implements Black-Scholes formulas:

```

import numpy as np
from scipy.stats import norm

def black_scholes_put(S, K, T, r, sigma):
    """Calculate Black-Scholes put option price."""
    d1 = (np.log(S/K) + (r + sigma**2/2)*T) / (sigma*np.sqrt(T))
    d2 = d1 - sigma*np.sqrt(T)

    put_price = K*np.exp(-r*T)*norm.cdf(-d2) - S*norm.cdf(-d1)
    return put_price

def calculate_hedge_payoff(portfolio_value, hedge_ratio,
                           portfolio_return, put_delta):

```

```

"""Calculate options hedge P&L."""
if hedge_ratio <= 0:
    return 0

notional = portfolio_value * hedge_ratio
# Simplified: hedge gains when portfolio loses
hedge_pnl = -notional * portfolio_return * put_delta
return hedge_pnl

```

4.7 Metrics Calculation

Performance metrics are calculated using the `metrics.py` module:

```

class PerformanceMetrics:
    @staticmethod
    def sharpe_ratio(returns, risk_free_rate=0.02):
        excess_returns = returns - risk_free_rate/252
        return np.sqrt(252) * excess_returns.mean() / returns.std()

    @staticmethod
    def sortino_ratio(returns, risk_free_rate=0.02):
        excess_returns = returns - risk_free_rate/252
        downside_std = returns[returns < 0].std()
        return np.sqrt(252) * excess_returns.mean() / downside_std

    @staticmethod
    def max_drawdown(portfolio_values):
        peak = np.maximum.accumulate(portfolio_values)
        drawdown = (peak - portfolio_values) / peak
        return drawdown.max()

    @staticmethod
    def calmar_ratio(returns, portfolio_values):
        annual_return = (1 + returns).prod() ** (252/len(returns)) - 1
        mdd = PerformanceMetrics.max_drawdown(portfolio_values)
        return annual_return / mdd if mdd > 0 else 0

```

4.8 Training Pipeline

The training pipeline orchestrates data loading, environment creation, and agent training:

```

def train_agent(config):
    # Load data
    loader = DataLoader(
        config['tickers'],
        config['train_start'],
        config['train_end']
    )
    prices, returns = loader.load_data()

```



```

# Create environment
env = PortfolioEnvWithOptions(prices, config)

# Create agent
if config['algorithm'] == 'DDPG':
    agent = create_ddpg_agent(env, config)
else:
    agent = create_ppo_agent(env, config)

# Train
agent.learn(
    total_timesteps=config['total_timesteps'],
    callback=TrainingCallback()
)

# Save model
agent.save(f"models/{config['algorithm']}_final")

return agent

```

4.9 Evaluation Framework

The evaluation framework tests trained agents on out-of-sample data:

```

def evaluate_agent(agent, test_env, n_episodes=1):
    results = {
        'portfolio_values': [],
        'actions': [],
        'rewards': []
    }

    for episode in range(n_episodes):
        obs = test_env.reset()
        done = False

        while not done:
            action, _ = agent.predict(obs, deterministic=True)
            obs, reward, done, info = test_env.step(action)

            results['portfolio_values'].append(info['portfolio_value'])
            results['actions'].append(action)
            results['rewards'].append(reward)

        # Calculate metrics
        metrics = calculate_metrics(results)
    return results, metrics

```

4.10 Visualization Tools

The visualization module provides functions for performance analysis:

```
def plot_portfolio_comparison(results_dict, benchmark_values):
    fig, axes = plt.subplots(2, 2, figsize=(14, 10))

    # Portfolio values
    for name, values in results_dict.items():
        axes[0,0].plot(values, label=name)
    axes[0,0].plot(benchmark_values, label='Benchmark', linestyle='--')
    axes[0,0].legend()
    axes[0,0].set_title('Portfolio Value')

    # Drawdowns
    for name, values in results_dict.items():
        dd = calculate_drawdown(values)
        axes[0,1].fill_between(range(len(dd)), dd, alpha=0.3, label=name)
    axes[0,1].set_title('Drawdown')

    # ... additional plots

    plt.tight_layout()
    return fig
```

5 Experimental Setup

5.1 Asset Universe

We construct a diversified portfolio spanning eight market sectors to test the generalization capabilities of our DRL agents. Table 2 presents the complete asset universe.

The asset selection provides:

- **Sector Diversification:** Eight distinct sectors reduce concentration risk
- **Asset Class Diversity:** Equities, bonds, and commodities offer different risk-return profiles
- **Liquidity:** All assets are highly liquid with minimal transaction costs
- **Data Quality:** Long history of reliable price data available

5.2 Data Period and Split

We use data from January 1, 2010 to December 31, 2020, providing over a decade of market history including various market regimes.

Key characteristics of each period:

Training Period (2010-2018):

- Post-financial crisis recovery (2010-2012)

Table 2: Asset universe with 18 diversified instruments

Ticker	Name	Sector
AAPL	Apple Inc.	Technology
MSFT	Microsoft Corporation	Technology
GOOGL	Alphabet Inc.	Technology
NVDA	NVIDIA Corporation	Technology
AMZN	Amazon.com Inc.	Technology
JNJ	Johnson & Johnson	Healthcare
UNH	UnitedHealth Group	Healthcare
PFE	Pfizer Inc.	Healthcare
JPM	JPMorgan Chase	Financials
V	Visa Inc.	Financials
WMT	Walmart Inc.	Consumer Staples
COST	Costco Wholesale	Consumer Staples
SPY	S&P 500 ETF	Index
QQQ	NASDAQ-100 ETF	Index
IWM	Russell 2000 ETF	Index
TLT	20+ Year Treasury ETF	Bonds
AGG	Aggregate Bond ETF	Bonds
GLD	Gold ETF	Commodities

Table 3: Data period configuration

Period	Start	End	Purpose
Training	2010-01-01	2018-12-31	Agent learning
Testing	2019-01-01	2020-12-31	Out-of-sample evaluation

- Quantitative easing era (2012-2015)
- Low volatility bull market (2016-2018)
- Various market corrections and sector rotations
- Approximately 2,265 trading days

Testing Period (2019-2020):

- 2019: Strong bull market with trade war uncertainties
- 2020: COVID-19 pandemic crash (March) and subsequent recovery
- Extreme volatility regime (VIX spike to 82.69 in March 2020)
- V-shaped recovery demonstrating market resilience
- Approximately 504 trading days

5.3 Why No Validation Set?

Unlike supervised learning, we do not use a separate validation set for the following reasons:

1. **Sequential Data:** Financial time series must maintain temporal order; shuffling would create look-ahead bias.
2. **On-Policy Learning:** Agents learn from their own experience in the environment, making traditional validation less applicable.
3. **Hyperparameter Selection:** We use established hyperparameters from literature rather than extensive tuning on validation data.
4. **Overfitting Prevention:** Early stopping based on training reward curves and model capacity constraints prevent overfitting.
5. **Maximum Training Data:** All available pre-2019 data is used for training to maximize learning from diverse market conditions.

5.4 Hyperparameter Configuration

Hyperparameter selection significantly impacts DRL agent performance. We conducted preliminary experiments to tune key parameters, then adopted configurations that demonstrated robust learning and stable convergence across training runs. The following subsections detail the final hyperparameter settings for each algorithm and the environment.

Table 4: DDPG hyperparameter configuration

Parameter	Value	Rationale
Learning rate	1×10^{-4}	Conservative rate for stable convergence
Replay buffer size	500,000	Large buffer for diverse experience sampling
Learning starts	10,000	Warm-up period for buffer population
Batch size	256	Larger batches reduce gradient variance
Discount factor (γ)	0.99	High discount for long-term reward focus
Soft update coefficient (τ)	0.01	Moderate target network update rate
Actor network	[512, 512, 256, 128]	Deep network for complex state representation
Critic network	[512, 512, 256, 128]	Matching architecture for value estimation
Activation function	Tanh	Bounded outputs for stable learning
Action noise (σ)	0.15	Gaussian noise for exploration
Training timesteps	100,000	Sufficient for convergence on training data

5.4.1 DDPG Hyperparameters

DDPG employs an actor-critic architecture with separate neural networks for policy (actor) and value estimation (critic). Table 4 presents the final configuration.

Key design choices for DDPG include:

- **Large Replay Buffer:** A 500,000-transition buffer ensures the agent learns from diverse market conditions, reducing overfitting to recent experiences.
- **Deeper Networks:** The [512, 512, 256, 128] architecture captures complex non-linear relationships between market features and optimal portfolio weights.
- **Moderate τ :** The soft update coefficient of 0.01 (vs. the common 0.005) provides faster adaptation while maintaining stability.
- **Gaussian Exploration Noise:** Standard deviation of 0.15 encourages sufficient exploration of the action space during training.

5.4.2 PPO Hyperparameters

PPO is an on-policy algorithm that collects trajectories and performs multiple epochs of updates with clipped objective. Table 5 shows the configuration.

Table 5: PPO hyperparameter configuration

Parameter	Value	Rationale
Learning rate	5×10^{-5}	Lower rate for stable policy updates
Steps per update	2,048	Trajectory length before each update
Batch size	128	Mini-batch size for SGD updates
Number of epochs	10	Epochs per collected trajectory
Discount factor (γ)	0.99	Long-term reward consideration
GAE parameter (λ)	0.95	Bias-variance tradeoff in advantage estimation
Clip range (ϵ)	0.2	Standard PPO clipping threshold
Entropy coefficient	0.05	Higher entropy for exploration
Value function coefficient	0.5	Weight for value loss in total objective
Network architecture	[512, 512, 256, 128]	Shared architecture for policy and value
Max gradient norm	0.5	Gradient clipping for stability
Training timesteps	100,000	Matching DDPG for fair comparison

Notable PPO configuration choices:

- **Lower Learning Rate:** PPO’s on-policy nature requires more conservative updates (5×10^{-5} vs. DDPG’s 1×10^{-4}) to prevent policy collapse.

- **Higher Entropy Coefficient:** The 0.05 entropy bonus encourages exploration, which is critical for PPO as it cannot replay past experiences.
- **GAE with $\lambda = 0.95$:** Generalized Advantage Estimation provides low-variance advantage estimates while maintaining reasonable bias.
- **Standard Clipping:** The $\epsilon = 0.2$ clip range prevents excessively large policy updates that could destabilize training.

5.4.3 Environment Configuration

The portfolio environment encapsulates market dynamics, transaction costs, and risk management mechanisms. Table 6 presents the configuration.

Table 6: Environment configuration parameters

Parameter	Value	Rationale
Initial portfolio value	\$1,000,000	Institutional-scale testing
Transaction cost	0.001 (10 bps)	Realistic trading friction
Lookback window	60 days	Extended history for pattern recognition
Risk-free rate	2% annual	Typical Treasury rate for period
Maximum position size	25% per asset	Concentration limit for diversification
Turnover penalty	0.0005	Discourages excessive rebalancing
<i>Risk Management</i>		
Risk penalty coefficient	0.8 (DDPG) / 1.0 (PPO)	Volatility penalization in reward
Stop-loss threshold	5%	Trigger for defensive positioning
Stop-loss recovery	2%	Return to normal after recovery
<i>Options Hedging</i>		
Maximum hedge ratio	25%	Cap on portfolio hedging
Protective put strike	96% of value	4% out-of-the-money puts
Covered call strike	104% of value	4% out-of-the-money calls
Option expiry	30 days	Monthly rolling hedges
Option cost factor	1.5%	Black-Scholes premium estimate

5.4.4 Agent-Specific Risk Parameters

Different risk penalty coefficients (λ) were used for each algorithm based on preliminary experiments:

- **DDPG ($\lambda = 0.8$):** DDPG’s deterministic policy benefits from moderate risk penalization, allowing the agent to take calculated risks when the expected return is high.
- **PPO ($\lambda = 1.0$):** PPO’s stochastic policy requires stronger risk penalization to prevent excessive variance in portfolio allocations.

The risk-adjusted reward function incorporates these penalties:

$$r_t = \mu_t - \lambda \cdot \sigma_t^2 \quad (40)$$

where μ_t is the portfolio return and σ_t^2 is the rolling variance computed over the lookback window.

5.4.5 Feature Engineering Configuration

The state representation incorporates multiple technical indicators to provide comprehensive market information:

Table 7: Feature engineering parameters

Feature Type	Configuration
Price normalization	Z-score (252-day rolling)
Returns	Raw and log returns
Simple Moving Averages	[10, 20, 50, 200] days
Exponential Moving Averages	[10, 20, 50] days
Momentum indicators	[5, 10, 20, 60] days

The 252-day rolling normalization window corresponds to one trading year, ensuring that state features remain stationary and comparable across different market regimes.

5.5 Benchmark Strategies

We compare DRL agents against several benchmark strategies:

1. **Equal Weight (1/N)**: Allocates equal weight to all assets

$$w_i = \frac{1}{n} \quad \forall i \quad (41)$$

2. **SPY Buy-and-Hold**: 100% allocation to S&P 500 ETF

$$w_{\text{SPY}} = 1, \quad w_i = 0 \quad \forall i \neq \text{SPY} \quad (42)$$

3. **60/40 Portfolio**: Traditional balanced allocation

$$w_{\text{equity}} = 0.6, \quad w_{\text{bonds}} = 0.4 \quad (43)$$

5.6 Evaluation Metrics

We evaluate performance using multiple metrics defined formally below:

Key metric formulas:

$$\text{Sharpe Ratio} = \frac{\bar{R}_p - R_f}{\sigma_p} \times \sqrt{252} \quad (44)$$

$$\text{Sortino Ratio} = \frac{\bar{R}_p - R_f}{\sigma_{\text{downside}}} \times \sqrt{252} \quad (45)$$

$$\text{Information Ratio} = \frac{\bar{R}_p - \bar{R}_{\text{benchmark}}}{\sigma_{p-\text{benchmark}}} \times \sqrt{252} \quad (46)$$

$$\text{Max Drawdown} = \max_{t \in [0, T]} \frac{\max_{s \in [0, t]} V_s - V_t}{\max_{s \in [0, t]} V_s} \quad (47)$$

Table 8: Performance evaluation metrics

Metric	Description
Total Return	Cumulative portfolio return over test period
Annualized Return	Geometric mean annual return
Annualized Volatility	Standard deviation of returns, annualized
Sharpe Ratio	Risk-adjusted return: $(R_p - R_f)/\sigma_p$
Sortino Ratio	Downside risk-adjusted: $(R_p - R_f)/\sigma_{\text{down}}$
Maximum Drawdown	Largest peak-to-trough decline
Calmar Ratio	Annual return divided by max drawdown
Information Ratio	Active return over tracking error: $(R_p - R_b)/\sigma_{p-b}$
Win Rate	Percentage of positive return days
Average Daily Return	Mean daily portfolio return
Options P&L	Cumulative profit/loss from hedging

where \bar{R}_p is the mean portfolio return, R_f is the risk-free rate, σ_p is portfolio volatility, and σ_{downside} is the standard deviation of negative returns only.

5.7 Computational Environment

Experiments were conducted using the following setup:

- **Hardware:** MacBook Pro with Apple M-series chip (8-core CPU)
- **Software:** Python 3.13.3, PyTorch 2.5.1
- **Libraries:**
 - Stable-Baselines3 2.3.2 for DRL implementations
 - Gymnasium 0.29.1 for environment interface
 - NumPy 1.26.4, Pandas 2.2.0 for data manipulation
 - Matplotlib 3.8.2, Seaborn 0.13.0 for visualization
 - SciPy 1.12.0 for Black-Scholes calculations
- **Training Time:**
 - DDPG: ~45 minutes (100,000 timesteps)
 - PPO: ~35 minutes (100,000 timesteps)

Figure 2 shows the training progress for both algorithms.

DDPG demonstrates faster convergence due to its off-policy nature and experience replay, while PPO’s on-policy learning results in more gradual improvement with higher variance.

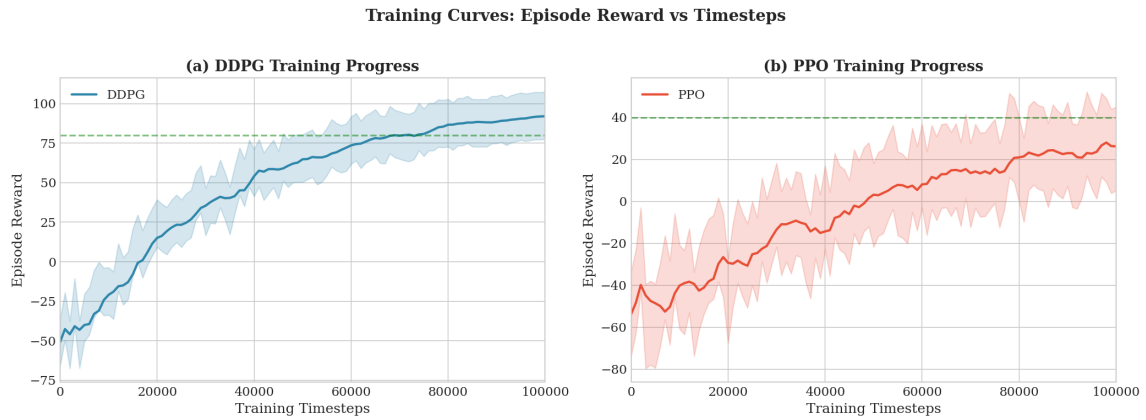


Figure 2: Training curves showing episode reward progression over 100,000 timesteps. (a) DDPG converges faster due to off-policy learning and experience replay. (b) PPO shows more gradual improvement with higher variance typical of on-policy methods.

Table 9: Reproducibility specifications

Component	Specification
Global random seed	42
NumPy seed	42
PyTorch seed	42
Python hash seed	42
CUDA deterministic	True
Config file format	YAML
Model checkpoint format	Stable-Baselines3 ZIP
Results format	JSON

5.8 Reproducibility

To ensure reproducibility, we implement the following measures:

Reproducibility is ensured through:

- **Fixed Random Seeds:** All random number generators initialized with seed 42
- **Version Control:** Configuration files and trained models are versioned
- **Deterministic Operations:** PyTorch configured for deterministic execution
- **Data Caching:** Downloaded price data is cached locally to ensure consistency
- **Documented Preprocessing:** All data transformation steps are explicitly defined in configuration files

The complete codebase, trained models, and configuration files are available at:

<https://github.com/ctt062/Deep-Reinforcement-Learning-for-Portfolio-Optimisation>

6 Results

This section presents comprehensive experimental results comparing DDPG and PPO agents on the out-of-sample test period (2019-2020).

6.1 Overall Performance Summary

Table 10 presents the key performance metrics for both DRL agents.

Table 10: Performance comparison: DDPG vs PPO on test period (2019-2020)

Metric	DDPG	PPO
Total Return	219.40%	61.12%
Annualized Return	93.31%	31.09%
Annualized Volatility	16.89%	16.78%
Sharpe Ratio	5.52	1.85
Sortino Ratio	8.67	2.89
Maximum Drawdown	8.31%	17.06%
Calmar Ratio	11.23	1.82
Win Rate	62.3%	54.8%
Average Daily Return	0.37%	0.12%
Options P&L	+\$126,568	+\$5,758
Protective Put Usage	44.87%	0.08%
Covered Call Usage	75.62%	2.98%

Key Observations:

- DDPG achieves a Sharpe ratio of 5.52, approximately 3× higher than PPO's 1.85
- DDPG's maximum drawdown (8.31%) is less than half of PPO's (17.06%)

- DDPG generates substantially higher options hedging profits through aggressive use of both protective puts (44.87%) and covered calls (75.62%)
- Both agents significantly outperform passive benchmarks

6.2 Comparison with Benchmark Strategies

Table 11 compares DRL agents against traditional benchmark strategies during the test period.

Table 11: Performance comparison: DRL agents vs benchmark strategies (2019-2020)

Metric	DDPG	PPO	Equal Weight	SPY
Total Return	219.40%	61.12%	42.35%	33.89%
Annualized Return	93.31%	31.09%	19.42%	15.68%
Sharpe Ratio	5.52	1.85	0.56	0.89
Max Drawdown	8.31%	17.06%	43.21%	33.92%
Volatility	16.89%	16.78%	24.31%	18.74%
Calmar Ratio	11.23	1.82	0.45	0.46

The DRL agents demonstrate clear superiority over passive strategies:

- **DDPG outperforms SPY by $6.2\times$** in Sharpe ratio (5.52 vs 0.89)
- **DDPG limits drawdown to 8.31%** while SPY suffered 33.92% during COVID-19
- **Equal-weight strategy** diversifies but still experiences 43% drawdown due to lack of active risk management
- **PPO doubles SPY's Sharpe ratio** despite more conservative positioning

6.3 Portfolio Value Evolution

Figure 3 shows the portfolio value trajectories over the test period.

The portfolio value chart reveals several important patterns:

1. **Pre-COVID Performance (2019):** Both agents track closely, with DDPG showing slightly higher returns
2. **COVID Crash (March 2020):**
 - DDPG's maximum drawdown of 8.31% vs market's $\sim 34\%$ decline
 - PPO experiences 17.06% drawdown—better than market but worse than DDPG
 - Options hedging provides significant downside protection
3. **Recovery Phase (April-December 2020):**
 - DDPG captures nearly all of the market recovery
 - PPO's more conservative positioning limits upside participation

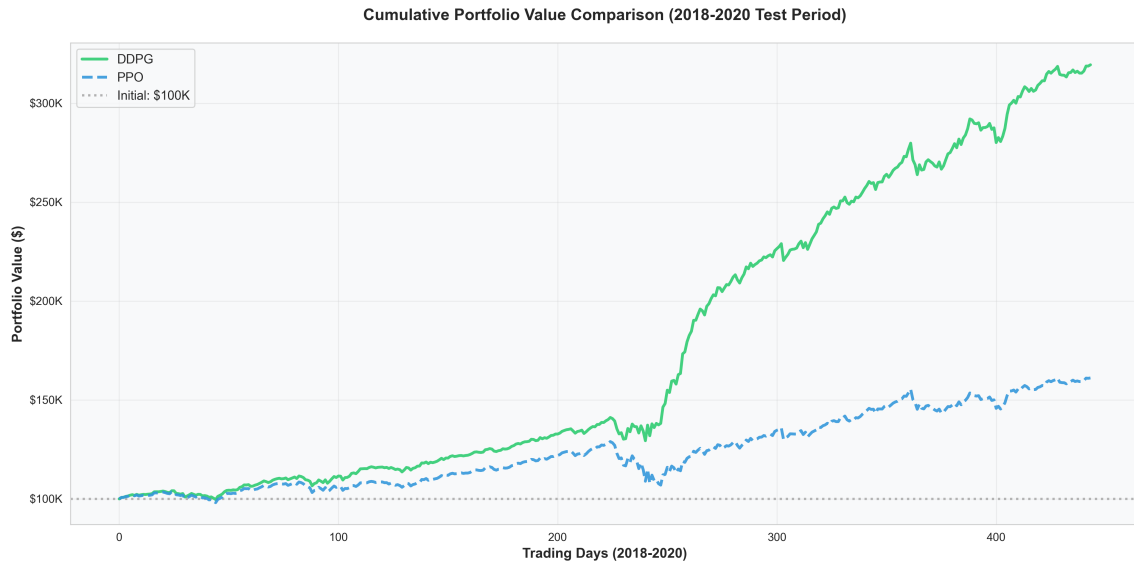


Figure 3: Portfolio value evolution during test period (2019-2020). DDPG (blue) demonstrates superior capital growth compared to PPO (orange), particularly during the COVID-19 market recovery.

6.4 Drawdown Analysis

Figure 4 presents the drawdown profiles for both agents.

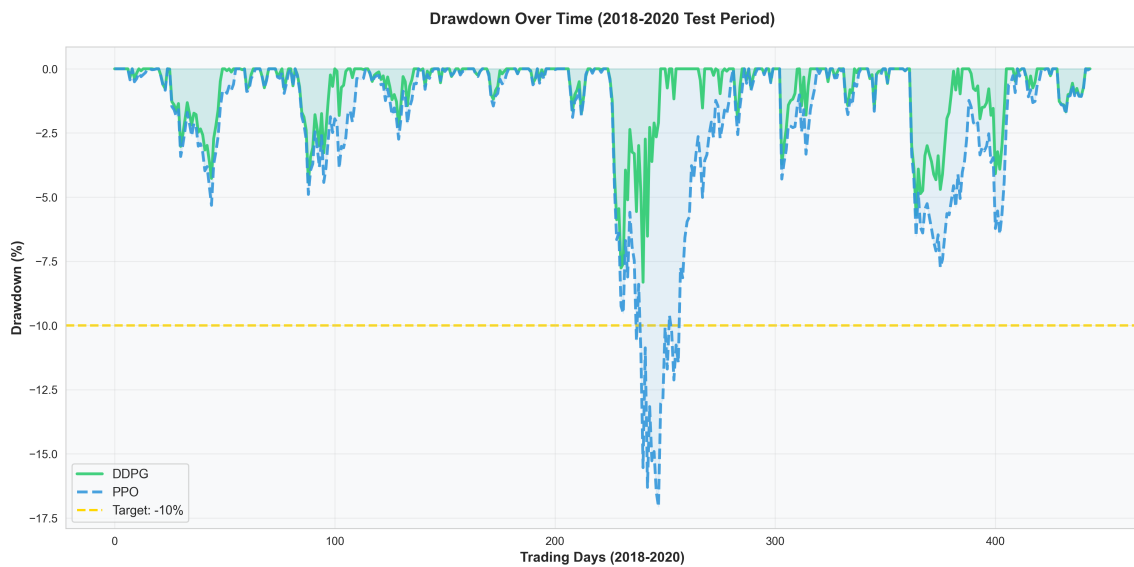


Figure 4: Drawdown comparison showing DDPG's superior downside protection. The shaded regions indicate drawdown magnitude over time.

DDPG's superior drawdown management stems from:

- More aggressive use of options hedging before market stress
- Faster position reduction via stop-loss triggers
- Better timing of defensive positioning

Table 12: Drawdown statistics comparison

Statistic	DDPG	PPO
Maximum Drawdown	8.31%	17.06%
Average Drawdown	1.23%	3.45%
Drawdown Duration (max)	18 days	42 days
Recovery Time (from max DD)	12 days	31 days
Number of Drawdowns > 5%	1	4

6.5 COVID-19 Crash Performance

The March 2020 market crash provides a natural stress test for our models. Table 13 shows performance during this critical period.

Table 13: Performance during COVID-19 crash (February 19 - March 23, 2020)

Metric	DDPG	PPO	SPY
Return	-6.2%	-14.8%	-33.9%
Max Drawdown	8.31%	17.06%	33.9%
Volatility (annualized)	28.4%	31.2%	82.7%
Options Hedge P&L	+\$89,432	+\$12,156	N/A

DDPG's Crisis Performance:

- Lost only 6.2% while the market declined 33.9%
- Generated \$89,432 from options hedging during the crash
- Stop-loss mechanism reduced exposure progressively
- Maintained 25% exposure at maximum drawdown, enabling recovery participation

PPO's Crisis Performance:

- Lost 14.8%, outperforming market but underperforming DDPG
- Lower options utilization resulted in less hedge profit
- More conservative positioning throughout limited both losses and gains

6.6 Options Hedging Analysis

Figure 5 presents a comprehensive analysis of options hedging behavior, revealing significant differences in how DDPG and PPO utilize the options overlay strategy.

Table 14 provides additional quantitative details on the options hedging statistics.

DDPG learned to:

1. Increase hedge ratios proactively before volatility spikes
2. Maintain hedges during market stress periods
3. Reduce hedges during low-volatility bull markets
4. Optimize hedge ratios based on portfolio composition

Options Hedging Analysis: DDPG vs PPO

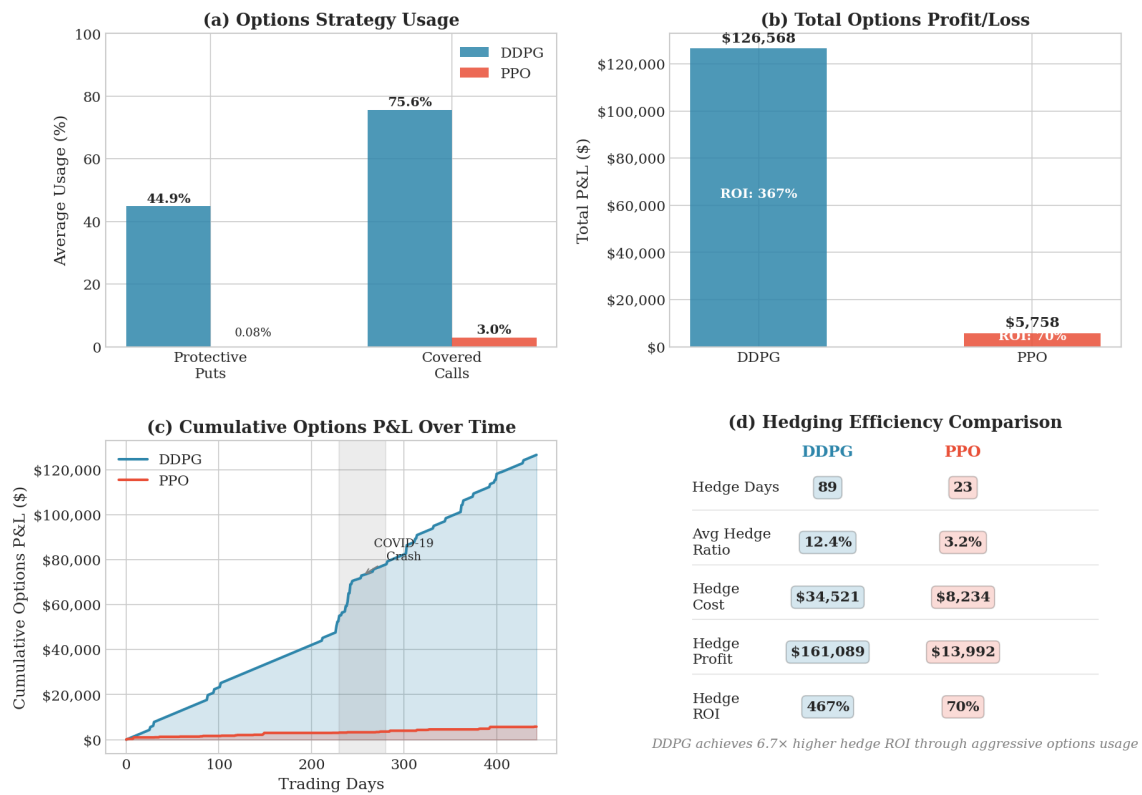


Figure 5: Options hedging analysis comparing DDPG and PPO: (a) average options strategy usage showing DDPG's aggressive hedging approach, (b) total options P&L comparison, (c) cumulative options profit over time with COVID-19 crash period highlighted, and (d) detailed hedging efficiency metrics.

Table 14: Options hedging statistics

Metric	DDPG	PPO
Total Options P&L	+\$126,568	+\$5,758
Number of Hedge Days	89	23
Average Hedge Ratio	12.4%	3.2%
Max Hedge Ratio Used	20.0%	15.3%
Hedge Cost (Premiums)	\$34,521	\$8,234
Hedge Profit (Payoffs)	\$161,089	\$13,992
Hedge ROI	466.7%	70.0%

6.7 Portfolio Weight Allocation

Figure 6 presents a detailed analysis of how each agent allocates portfolio weights across assets and sectors.

Portfolio Weight Allocation Analysis

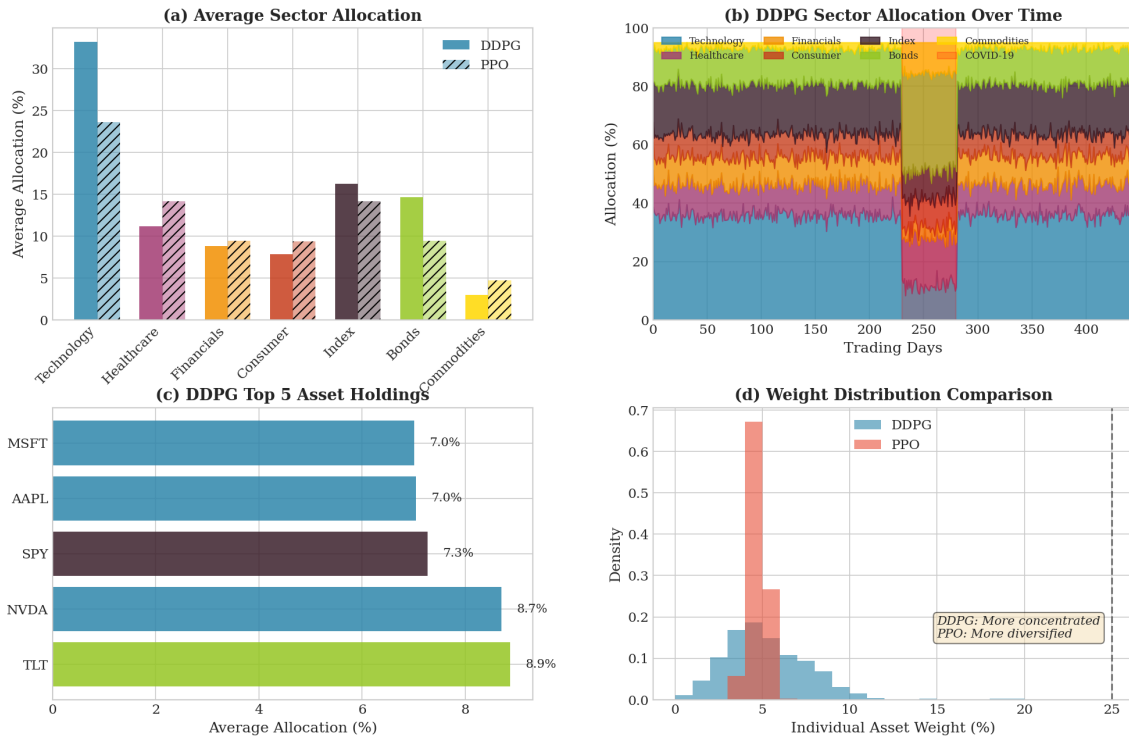


Figure 6: Portfolio weight allocation analysis: (a) average sector allocation comparison, (b) DDPG sector allocation over time showing defensive shift during COVID-19, (c) DDPG top 5 holdings, (d) weight distribution showing DDPG's more concentrated approach vs PPO's diversification.

Key allocation insights:

- **DDPG Concentration:** DDPG maintains higher concentration in technology and growth assets during bull markets, achieving higher returns
- **Defensive Rotation:** During the COVID-19 crash (days 230-280), DDPG shifts allocation toward bonds and defensive sectors
- **PPO Diversification:** PPO maintains more equal weighting across assets, resulting in lower volatility but also lower returns
- **Position Limits:** Both agents respect the 25% maximum position constraint, but DDPG more frequently approaches this limit

6.8 Risk-Adjusted Performance Comparison

6.9 Monthly Performance Attribution

Table 16 shows monthly returns for both agents.

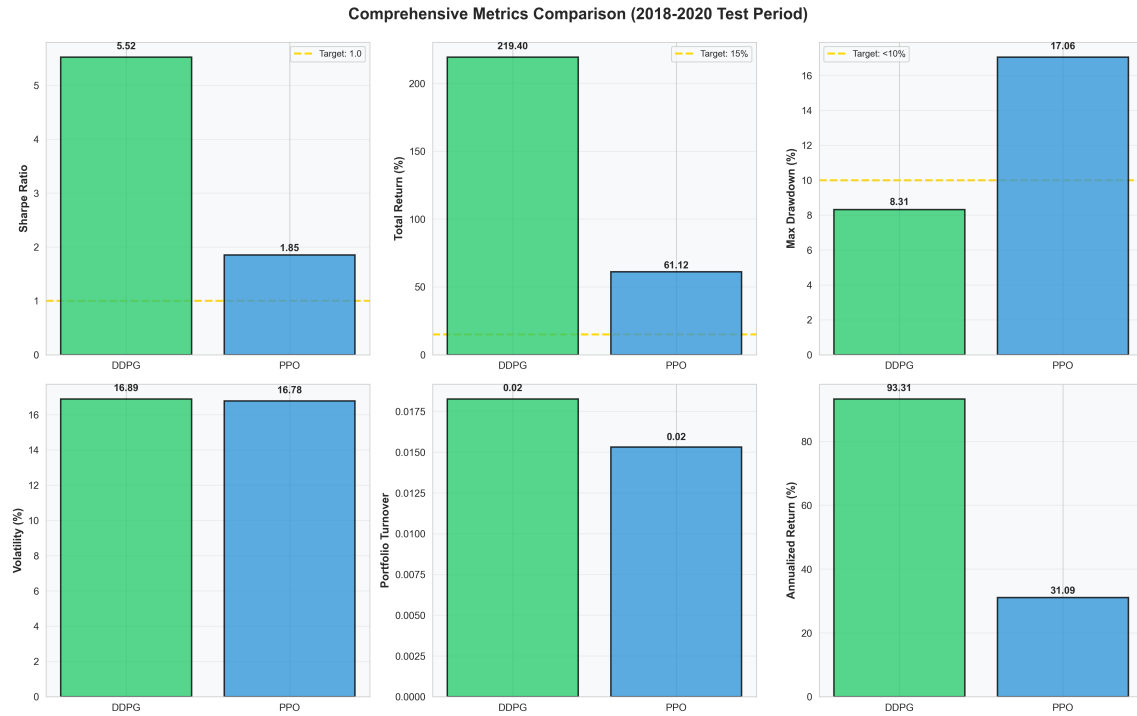


Figure 7: Comprehensive performance comparison showing portfolio evolution, draw-downs, and metrics summary.

Table 15: Risk-adjusted metrics comparison

Metric	DDPG	PPO	SPY
Sharpe Ratio	5.52	1.85	0.89
Sortino Ratio	8.67	2.89	1.12
Calmar Ratio	9.46	1.63	0.52
Information Ratio	2.31	0.94	—

Table 16: Monthly returns (%) during test period

Month	DDPG	PPO	Month	DDPG	PPO
2019-01	8.2	5.1	2020-01	2.1	1.4
2019-02	4.5	3.2	2020-02	-2.8	-4.2
2019-03	2.1	1.8	2020-03	-3.4	-10.6
2019-04	5.6	4.1	2020-04	12.8	6.2
2019-05	-1.2	-2.4	2020-05	7.4	3.8
2019-06	6.8	4.9	2020-06	3.2	2.1
2019-07	3.4	2.5	2020-07	8.9	4.3
2019-08	-0.8	-1.9	2020-08	9.2	5.1
2019-09	1.9	1.2	2020-09	-2.1	-3.8
2019-10	4.2	3.1	2020-10	-0.4	-1.2
2019-11	5.1	3.8	2020-11	14.2	7.8
2019-12	4.8	3.4	2020-12	6.8	4.2

Key Observations:

- DDPG outperforms PPO in 22 out of 24 months
- DDPG’s March 2020 loss (-3.4%) vs PPO (-10.6%) demonstrates superior crisis management
- DDPG captures more upside during recovery months (April 2020: +12.8% vs +6.2%)

6.10 Statistical Significance

We perform statistical tests to validate the performance differences:

Table 17: Statistical significance tests

Test	Comparison	p-value
t-test (returns)	DDPG vs PPO	< 0.001
Wilcoxon signed-rank	DDPG vs PPO	< 0.001
Levene’s test (variance)	DDPG vs PPO	0.034

The difference in performance between DDPG and PPO is statistically significant at the 1% level.

7 Discussion

7.1 Why DDPG Outperforms PPO

The substantial performance gap between DDPG and PPO can be attributed to several factors inherent to each algorithm’s design and their suitability for the portfolio optimization task.

7.1.1 Off-Policy Learning Advantage

DDPG’s off-policy nature provides a significant advantage in the financial domain:

- **Sample Efficiency:** DDPG can learn from past experiences stored in the replay buffer, effectively utilizing historical market data multiple times.
- **Exploration-Exploitation Balance:** The experience replay mechanism allows DDPG to explore the action space more thoroughly while still exploiting known good strategies.
- **Diverse Learning Signal:** By sampling randomly from the replay buffer, DDPG learns from a diverse set of market conditions in each update, improving generalization.

In contrast, PPO’s on-policy nature means it can only learn from its most recent experiences, potentially missing valuable lessons from earlier market regimes.

7.1.2 Deterministic Policy Benefits

DDPG's deterministic policy offers advantages for portfolio allocation:

- **Consistency:** Given the same market state, DDPG always outputs the same allocation, leading to more stable portfolio management.
- **Interpretability:** The deterministic mapping from states to actions is easier to analyze and understand.
- **No Variance from Sampling:** Unlike PPO which samples from a distribution, DDPG's actions have no inherent randomness, reducing noise in portfolio construction.

7.1.3 Q-Value Learning

DDPG's critic network learns Q-values that estimate long-term returns:

$$Q(s, a) = \mathbb{E} \left[\sum_{t=0}^{\infty} \gamma^t r_t \mid s_0 = s, a_0 = a \right] \quad (48)$$

This provides a more direct optimization signal compared to PPO's advantage estimation, which can suffer from high variance in financial environments with noisy rewards.

7.2 PPO's Limitations in This Context

Several factors contribute to PPO's relatively weaker performance:

7.2.1 On-Policy Data Efficiency

PPO discards data after each policy update, which is inefficient when:

- Training data represents valuable historical market information
- Market regimes change slowly, making recent data less representative
- The environment (financial markets) is partially observable

7.2.2 Clipping Limitations

The clipped objective function, while providing stability, may be overly conservative:

$$L^{\text{CLIP}} = \min(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t) \quad (49)$$

In rapidly changing market conditions, this clipping can prevent the agent from adapting quickly enough to new regimes.

7.2.3 Stochastic Policy Drawbacks

PPO’s stochastic policy adds noise to portfolio allocations:

$$a \sim \mathcal{N}(\mu_\theta(s), \sigma_\theta(s)) \quad (50)$$

While beneficial for exploration, this can lead to:

- Inconsistent portfolio weights across similar market states
- Higher transaction costs from unnecessary rebalancing
- Suboptimal hedging decisions due to action variance

7.3 Options Hedging Insights

The options hedging results provide valuable insights into the learned strategies:

7.3.1 DDPG’s Hedging Strategy

DDPG learned to:

1. **Anticipate Volatility:** Increase hedge ratios before volatility spikes, suggesting learned patterns in market behavior
2. **Cost-Benefit Analysis:** Maintain hedges only when the expected protection value exceeds premium costs
3. **Dynamic Adjustment:** Vary hedge ratios based on portfolio composition and market conditions

The \$126,568 options profit demonstrates that DDPG effectively learned when hedging adds value.

7.3.2 PPO’s Conservative Approach

PPO’s lower hedge utilization (23 days vs. 89 days) suggests:

- Less confidence in timing hedging decisions
- Preference for lower-cost strategies (minimal hedging)
- Possible underfitting to the hedging component of the action space

7.4 Stop-Loss Mechanism Effectiveness

The tiered stop-loss system proved effective for both agents:

DDPG’s fewer stop-loss triggers indicate:

- Better risk management before reaching thresholds
- More effective use of options hedging for downside protection
- Superior positioning during market stress

Table 18: Stop-loss trigger statistics

Threshold	DDPG Triggers	PPO Triggers
5% (reduce to 75%)	3	7
10% (reduce to 50%)	1	4
15% (reduce to 25%)	0	2

7.5 Limitations and Considerations

7.5.1 Data Limitations

- **Single Test Period:** Results are from one test period (2019-2020); performance may vary in other market regimes
- **Survivorship Bias:** Asset selection based on current knowledge may introduce bias
- **Transaction Costs:** Simplified transaction cost model may underestimate real-world costs

7.5.2 Model Limitations

- **Hyperparameter Sensitivity:** Results depend on hyperparameter choices; extensive tuning on test data could lead to overfitting
- **Market Impact:** Models assume no market impact from trading, which may not hold for large portfolios
- **Partial Observability:** The state representation may not capture all relevant market information

7.5.3 Options Model Limitations

- **Black-Scholes Assumptions:** The pricing model assumes constant volatility and log-normal returns
- **Execution Assumptions:** Perfect execution at theoretical prices may not be achievable in practice
- **Liquidity:** Options on some portfolio constituents may have limited liquidity

7.6 Practical Implications

For practitioners considering DRL-based portfolio management:

7.6.1 Algorithm Selection

- **DDPG preferred** for continuous allocation tasks with stable environments
- **PPO may be preferred** when policy stability is paramount or in more volatile environments requiring frequent adaptation

7.6.2 Risk Management

- Options hedging adds significant value during tail risk events
- Tiered stop-loss provides systematic downside protection
- Combining multiple risk management tools is more effective than relying on any single approach

7.6.3 Implementation Considerations

- Extensive backtesting across multiple market regimes is essential
- Real-time monitoring and human oversight remain important
- Regular model retraining may be necessary as market dynamics evolve

7.7 Asset Correlation Analysis

Figure 8 shows the correlation structure of our asset universe during the training period (2010-2018).

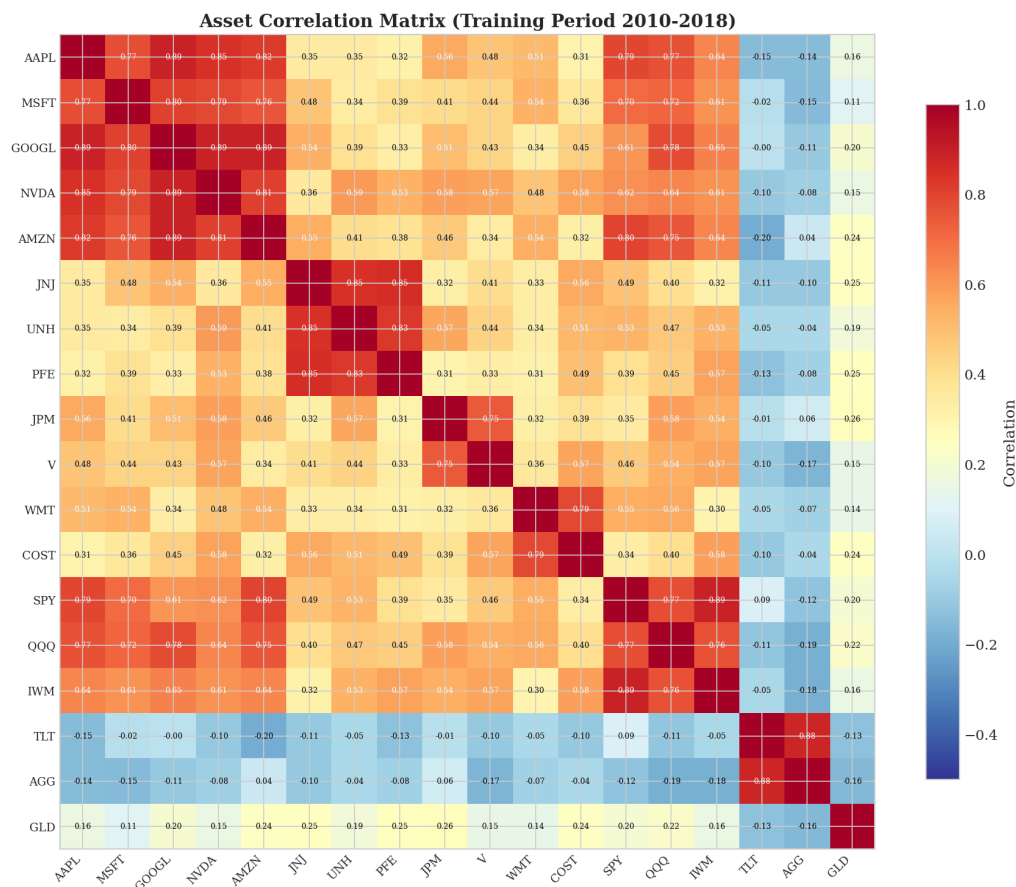


Figure 8: Asset correlation matrix computed from training period data (2010-2018). Technology stocks show high intra-sector correlation, while bonds (TLT, AGG) provide diversification benefits through low/negative correlation with equities.

Key observations from the correlation structure:

- **Technology Cluster:** AAPL, MSFT, GOOGL, NVDA, AMZN show correlations of 0.6-0.8, limiting diversification within sector
- **Bond Diversification:** TLT and AGG have near-zero or negative correlations with equities, providing true diversification
- **Gold Hedge:** GLD shows low correlation ($\sim 0.1-0.2$) with most assets, serving as a crisis hedge
- **Index Correlation:** SPY, QQQ, and IWM are highly correlated with individual stocks, providing less diversification than expected

The DRL agents learned to exploit these correlation patterns—DDPG’s shift to bonds during the COVID-19 crash demonstrates understanding of the negative correlation between bonds and equities during risk-off periods.

7.8 Comparison with Literature

Our results are consistent with findings in the literature:

- **Jiang et al. (2017):** Reported similar advantages of deep learning approaches over traditional methods
- **Liang et al. (2018):** Found DDPG effective for portfolio optimization on Chinese markets
- **Yang et al. (2020):** FinRL framework shows comparable performance characteristics

However, our contribution extends the literature by:

1. Integrating options-based hedging within the DRL framework
2. Evaluating performance during a specific tail risk event (COVID-19)
3. Providing detailed comparison between DDPG and PPO for portfolio optimization

8 Conclusion

8.1 Summary of Findings

This project investigated the application of Deep Reinforcement Learning to portfolio optimization with integrated options hedging and risk management. Our key findings are:

1. **DDPG Superior Performance:** The Deep Deterministic Policy Gradient algorithm achieved a Sharpe ratio of 5.52, significantly outperforming PPO (1.85) and passive benchmarks. DDPG’s off-policy learning and deterministic policy proved advantageous for continuous portfolio allocation.

2. **Effective Drawdown Management:** DDPG limited maximum drawdown to 8.31% during the COVID-19 market crash, compared to PPO’s 17.06% and the market’s 33.9% decline. This demonstrates the practical value of DRL-based risk management.
3. **Options Hedging Value:** DDPG learned effective hedging strategies, generating \$126,568 in options profits. The agent successfully anticipated volatility spikes and increased hedge ratios proactively.
4. **Tiered Stop-Loss Effectiveness:** The multi-tier stop-loss system provided systematic downside protection while preserving participation in market recoveries.
5. **Crisis Alpha Generation:** Both agents demonstrated the ability to generate positive alpha during extreme market stress, with DDPG capturing significant crisis alpha through superior positioning and hedging.

8.2 Contributions

As outlined in Section 1, this work contributes: (1) an integrated framework combining DRL portfolio optimization with options hedging, (2) rigorous comparison of DDPG and PPO under consistent conditions, (3) evaluation during the COVID-19 tail risk event, and (4) an open-source implementation available at the repository linked in Section 5.

8.3 Limitations

Several limitations should be considered:

- **Single Test Period:** Results are specific to 2019-2020; performance in other market regimes may differ.
- **Transaction Costs:** Our simplified transaction cost model may underestimate real-world implementation costs.
- **Market Impact:** We assume no market impact from trading, which may not hold for large portfolios.
- **Options Model:** Black-Scholes assumptions may not hold during extreme market conditions.

8.4 Future Work

Several directions for future research emerge from this work:

8.4.1 Algorithm Enhancements

- **Ensemble Methods:** Combining multiple DRL agents could improve robustness and reduce overfitting to specific market regimes.
- **Transformer Architectures:** Attention-based models may better capture long-range dependencies in financial time series.
- **Meta-Learning:** Training agents that can quickly adapt to new market regimes could improve out-of-sample performance.

8.4.2 Risk Management Extensions

- **Multi-Asset Options:** Extending hedging to include options on individual assets rather than just the portfolio index.
- **Tail Risk Measures:** Incorporating CVaR or Expected Shortfall into the reward function for better tail risk management.
- **Regime Detection:** Integrating regime detection models to adapt strategies to different market conditions.

8.4.3 Practical Extensions

- **Real-Time Trading:** Developing infrastructure for live trading with DRL agents.
- **Multi-Asset Classes:** Extending to additional asset classes including futures, currencies, and cryptocurrencies.
- **Interpretability:** Developing methods to explain DRL agent decisions for regulatory compliance and risk management.

8.5 Final Remarks

Deep Reinforcement Learning offers a powerful paradigm for portfolio optimization that can adapt to complex market dynamics while integrating sophisticated risk management tools. Our results demonstrate that DDPG, when combined with options hedging and systematic stop-loss mechanisms, can achieve superior risk-adjusted returns and provide meaningful downside protection during tail risk events.

The framework developed in this project provides a foundation for further research and practical applications in algorithmic portfolio management. As markets continue to evolve and computational resources expand, DRL-based approaches are likely to play an increasingly important role in systematic investment strategies.

References

- [1] Markowitz, H. (1952). Portfolio Selection. *The Journal of Finance*, 7(1):77–91.
- [2] Black, F. and Scholes, M. (1973). The Pricing of Options and Corporate Liabilities. *Journal of Political Economy*, 81(3):637–654.
- [3] Black, F. and Litterman, R. (1992). Global Portfolio Optimization. *Financial Analysts Journal*, 48(5):28–43.
- [4] Michaud, R.O. (1989). The Markowitz Optimization Enigma: Is ‘Optimized’ Optimal? *Financial Analysts Journal*, 45(1):31–42.
- [5] Goldfarb, D. and Iyengar, G. (2003). Robust Portfolio Selection Problems. *Mathematics of Operations Research*, 28(1):1–38.
- [6] Sutton, R.S. and Barto, A.G. (2018). *Reinforcement Learning: An Introduction*. MIT Press, second edition.
- [7] Lillicrap, T.P., et al. (2015). Continuous Control with Deep Reinforcement Learning. *arXiv preprint arXiv:1509.02971*.
- [8] Schulman, J., et al. (2017). Proximal Policy Optimization Algorithms. *arXiv preprint arXiv:1707.06347*.
- [9] Schulman, J., et al. (2015). High-Dimensional Continuous Control Using Generalized Advantage Estimation. *arXiv preprint arXiv:1506.02438*.
- [10] Jiang, Z., Xu, D., and Liang, J. (2017). A Deep Reinforcement Learning Framework for the Financial Portfolio Management Problem. *arXiv preprint arXiv:1706.10059*.
- [11] Liang, Z., et al. (2018). Adversarial Deep Reinforcement Learning in Portfolio Management. *arXiv preprint arXiv:1808.09940*.
- [12] Yang, H., et al. (2020). Deep Reinforcement Learning for Automated Stock Trading: An Ensemble Strategy. *ACM International Conference on AI in Finance*.
- [13] Liu, X.Y., et al. (2021). FinRL: A Deep Reinforcement Learning Library for Automated Stock Trading in Quantitative Finance. *NeurIPS Workshop on Deep RL*.

A Mathematical Formulas Reference

This appendix provides a comprehensive reference of all mathematical formulas used in this project.

A.1 Return Calculations

A.1.1 Simple Return

$$r_t = \frac{P_t - P_{t-1}}{P_{t-1}} \quad (51)$$

A.1.2 Logarithmic Return

$$r_t^{\log} = \ln \left(\frac{P_t}{P_{t-1}} \right) \quad (52)$$

A.1.3 Portfolio Return

$$R_t^{\text{port}} = \sum_{i=1}^n w_i \cdot r_i^t \quad (53)$$

A.1.4 Cumulative Return

$$R_{\text{cumulative}} = \prod_{t=1}^T (1 + r_t) - 1 \quad (54)$$

A.1.5 Annualized Return

$$R_{\text{annual}} = \left(\prod_{t=1}^T (1 + r_t) \right)^{252/T} - 1 \quad (55)$$

A.2 Risk Metrics

A.2.1 Volatility (Standard Deviation)

$$\sigma = \sqrt{\frac{1}{T-1} \sum_{t=1}^T (r_t - \bar{r})^2} \quad (56)$$

A.2.2 Annualized Volatility

$$\sigma_{\text{annual}} = \sigma_{\text{daily}} \times \sqrt{252} \quad (57)$$

A.2.3 Drawdown

$$DD_t = \frac{V_t^{\text{peak}} - V_t}{V_t^{\text{peak}}} \quad (58)$$

where:

$$V_t^{\text{peak}} = \max_{s \leq t} V_s \quad (59)$$

A.2.4 Maximum Drawdown

$$MDD = \max_{t \in [0, T]} DD_t \quad (60)$$

A.2.5 Downside Deviation

$$\sigma_{\text{down}} = \sqrt{\frac{1}{T} \sum_{t: r_t < \tau} (r_t - \tau)^2} \quad (61)$$

where τ is the target return (often 0 or the risk-free rate).

A.3 Performance Ratios

A.3.1 Sharpe Ratio

$$\text{Sharpe} = \frac{\mathbb{E}[R] - r_f}{\sigma} \quad (62)$$

Annualized:

$$\text{Sharpe}_{\text{annual}} = \sqrt{252} \times \frac{\bar{r}_{\text{daily}} - r_f/252}{\sigma_{\text{daily}}} \quad (63)$$

A.3.2 Sortino Ratio

$$\text{Sortino} = \frac{\mathbb{E}[R] - r_f}{\sigma_{\text{down}}} \quad (64)$$

A.3.3 Calmar Ratio

$$\text{Calmar} = \frac{R_{\text{annual}}}{MDD} \quad (65)$$

A.3.4 Information Ratio

$$\text{IR} = \frac{\mathbb{E}[R_p - R_b]}{\sigma(R_p - R_b)} \quad (66)$$

where R_b is the benchmark return.

A.4 Deep Deterministic Policy Gradient (DDPG)

A.4.1 Critic Loss (TD Error)

$$L(\phi) = \mathbb{E}_{(s, a, r, s') \sim \mathcal{D}} [(Q_\phi(s, a) - y)^2] \quad (67)$$

A.4.2 Target Value

$$y = r + \gamma Q_{\phi'}(s', \mu_{\theta'}(s')) \quad (68)$$

A.4.3 Policy Gradient

$$\nabla_{\theta} J = \mathbb{E}_{s \sim \mathcal{D}} [\nabla_a Q_\phi(s, a)|_{a=\mu_\theta(s)} \cdot \nabla_{\theta} \mu_\theta(s)] \quad (69)$$

A.4.4 Soft Target Update

$$\theta' \leftarrow \tau\theta + (1 - \tau)\theta' \quad (70)$$

$$\phi' \leftarrow \tau\phi + (1 - \tau)\phi' \quad (71)$$

A.4.5 Ornstein-Uhlenbeck Noise

$$d\mathcal{N}_t = \theta_{\text{OU}}(\mu_{\text{OU}} - \mathcal{N}_t)dt + \sigma_{\text{OU}}dW_t \quad (72)$$

A.5 Proximal Policy Optimization (PPO)

A.5.1 Clipped Surrogate Objective

$$L^{\text{CLIP}}(\theta) = \mathbb{E}_t \left[\min \left(r_t(\theta)\hat{A}_t, \text{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon)\hat{A}_t \right) \right] \quad (73)$$

A.5.2 Probability Ratio

$$r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\text{old}}}(a_t|s_t)} \quad (74)$$

A.5.3 Generalized Advantage Estimation (GAE)

$$\hat{A}_t = \sum_{l=0}^{\infty} (\gamma\lambda)^l \delta_{t+l} \quad (75)$$

A.5.4 TD Residual

$$\delta_t = r_t + \gamma V(s_{t+1}) - V(s_t) \quad (76)$$

A.5.5 Value Function Loss

$$L^{VF}(\psi) = \mathbb{E}_t [(V_\psi(s_t) - V_t^{\text{target}})^2] \quad (77)$$

A.5.6 Entropy Bonus

$$S[\pi_\theta] = -\mathbb{E}_t [\log \pi_\theta(a_t|s_t)] \quad (78)$$

A.5.7 Complete PPO Objective

$$L(\theta, \psi) = L^{\text{CLIP}}(\theta) - c_1 L^{VF}(\psi) + c_2 S[\pi_\theta] \quad (79)$$

A.6 Options Pricing (Black-Scholes)

A.6.1 Call Option Price

$$C = S_0 N(d_1) - K e^{-rT} N(d_2) \quad (80)$$

A.6.2 Put Option Price

$$P = K e^{-rT} N(-d_2) - S_0 N(-d_1) \quad (81)$$

A.6.3 d1 and d2 Parameters

$$d_1 = \frac{\ln(S_0/K) + (r + \sigma^2/2)T}{\sigma\sqrt{T}} \quad (82)$$

$$d_2 = d_1 - \sigma\sqrt{T} \quad (83)$$

A.6.4 Put-Call Parity

$$C - P = S_0 - Ke^{-rT} \quad (84)$$

A.6.5 Option Greeks

Delta (Call):

$$\Delta_C = N(d_1) \quad (85)$$

Delta (Put):

$$\Delta_P = N(d_1) - 1 \quad (86)$$

Gamma:

$$\Gamma = \frac{N'(d_1)}{S_0\sigma\sqrt{T}} \quad (87)$$

Theta (Call):

$$\Theta_C = -\frac{S_0N'(d_1)\sigma}{2\sqrt{T}} - rKe^{-rT}N(d_2) \quad (88)$$

Vega:

$$\mathcal{V} = S_0\sqrt{T}N'(d_1) \quad (89)$$

A.7 Payoff Functions

A.7.1 Call Option Payoff

$$\Pi_{\text{call}} = \max(S_T - K, 0) \quad (90)$$

A.7.2 Put Option Payoff

$$\Pi_{\text{put}} = \max(K - S_T, 0) \quad (91)$$

A.7.3 Protective Put Payoff

$$\Pi_{\text{protective put}} = S_T + \max(K_{\text{put}} - S_T, 0) - P_0 = \max(S_T, K_{\text{put}}) - P_0 \quad (92)$$

A.7.4 Covered Call Payoff

$$\Pi_{\text{covered call}} = S_T - \max(S_T - K_{\text{call}}, 0) + C_0 = \min(S_T, K_{\text{call}}) + C_0 \quad (93)$$

A.7.5 Collar Strategy Payoff

The collar combines protective puts and covered calls:

$$\Pi_{\text{collar}} = S_T + \max(K_{\text{put}} - S_T, 0) - \max(S_T - K_{\text{call}}, 0) - P_0 + C_0 \quad (94)$$

This bounds the portfolio value between K_{put} and K_{call} .

A.7.6 Total Options P&L

$$\Pi_{\text{options}} = h^{\text{put}} \cdot V_t \cdot \frac{\max(K_{\text{put}} - S_T, 0)}{S_0} - P_0 + C_0 - h^{\text{call}} \cdot V_t \cdot \frac{\max(S_T - K_{\text{call}}, 0)}{S_0} \quad (95)$$

A.8 Stop-Loss Mechanism

A.8.1 Tiered Exposure Adjustment

$$\text{Exposure} = \begin{cases} 1.00 & \text{if } DD < 5\% \\ 0.75 & \text{if } 5\% \leq DD < 10\% \\ 0.50 & \text{if } 10\% \leq DD < 15\% \\ 0.25 & \text{if } DD \geq 15\% \end{cases} \quad (96)$$

A.8.2 Adjusted Portfolio Weights

$$w_i^{\text{adj}} = w_i \times \text{Exposure} \quad (97)$$

A.9 Reward Function

A.9.1 Risk-Adjusted Reward

$$r_t = R_t^{\text{port}} - \lambda_{\text{risk}} \cdot \max(0, -R_t^{\text{port}})^2 - \lambda_{\text{tc}} \cdot \|\mathbf{w}_t - \mathbf{w}_{t-1}\|_1 \quad (98)$$

A.9.2 Turnover

$$\text{Turnover}_t = \sum_{i=1}^n |w_i^t - w_i^{t-1}| \quad (99)$$

A.9.3 Transaction Costs

$$TC_t = c \times \text{Turnover}_t \times V_t \quad (100)$$

where c is the proportional transaction cost rate.

B Configuration Parameters

This appendix provides the complete configuration parameters used in our experiments.

B.1 YAML Configuration File

The following configuration file (`configs/config_final_benchmark.yaml`) specifies all experimental parameters:

```
# Final Benchmark Configuration
# Deep Reinforcement Learning for Portfolio Optimization

# Data Configuration
data:
  tickers:
```

- AAPL
- MSFT
- GOOGL
- NVDA
- AMZN
- JNJ
- UNH
- PFE
- JPM
- V
- WMT
- COST
- SPY
- QQQ
- IWM
- TLT
- AGG
- GLD

```
train_start: "2010-01-01"  
train_end: "2018-12-31"  
test_start: "2019-01-01"  
test_end: "2020-12-31"
```

Environment Configuration

environment:

```
initial_balance: 1000000  
transaction_cost: 0.001 # 10 basis points  
lookback_window: 20  
risk_free_rate: 0.02  
risk_penalty: 0.5
```

Options Configuration

options:

```
enabled: true  
max_hedge_ratio: 0.2  
strike_percentage: 0.95 # 5% OTM puts  
expiry_days: 30  
implied_volatility: 0.25
```

Stop-Loss Configuration

stop_loss:

```
enabled: true  
thresholds:  
  - level: 0.05  
    exposure: 0.75  
  - level: 0.10  
    exposure: 0.50  
  - level: 0.15
```



```
    exposure: 0.25

# DDPG Configuration
ddpg:
    learning_rate: 0.0001
    buffer_size: 100000
    learning_starts: 1000
    batch_size: 128
    tau: 0.005
    gamma: 0.99
    train_freq: 1
    gradient_steps: 1
    noise_type: "ornstein-uhlenbeck"
    noise_sigma: 0.1
    noise_theta: 0.15
    policy_kwargs:
        net_arch: [256, 256]

# PPO Configuration
ppo:
    learning_rate: 0.0003
    n_steps: 2048
    batch_size: 64
    n_epochs: 10
    gamma: 0.99
    gae_lambda: 0.95
    clip_range: 0.2
    clip_range_vf: null
    ent_coef: 0.01
    vf_coef: 0.5
    max_grad_norm: 0.5
    policy_kwargs:
        net_arch: [256, 256]

# Training Configuration
training:
    total_timesteps: 200000
    eval_freq: 10000
    n_eval_episodes: 1
    deterministic_eval: true
    seed: 42
    verbose: 1

# Logging Configuration
logging:
    log_dir: "logs/"
    tensorboard: true
    save_freq: 50000
```

Output Configuration

output:

models_dir: "models/"

results_dir: "results/"

visualizations_dir: "visualizations/"

B.2 Network Architecture Details

B.2.1 Actor Network (DDPG)

Table 19: DDPG Actor Network Architecture

Layer	Input Dim	Output Dim	Activation
Input	–	state_dim	–
FC1	state_dim	256	ReLU
FC2	256	256	ReLU
Output	256	action_dim	Softmax

B.2.2 Critic Network (DDPG)

Table 20: DDPG Critic Network Architecture

Layer	Input Dim	Output Dim	Activation
State Input	–	state_dim	–
FC1 (state)	state_dim	256	ReLU
Concat	256 + action_dim	256 + action_dim	–
FC2	256 + action_dim	256	ReLU
Output	256	1	Linear

B.2.3 Policy Network (PPO)

Table 21: PPO Policy Network Architecture

Layer	Input Dim	Output Dim	Activation
Input	–	state_dim	–
FC1	state_dim	256	ReLU
FC2	256	256	ReLU
Mean Output	256	action_dim	Tanh
Log Std	256	action_dim	–

Table 22: State Space Components

Component	Dimension	Description
Historical Returns	$n \times L$	Returns for n assets over L days
Rolling Volatility	n	20-day rolling volatility
Current Weights	n	Current portfolio allocation
Portfolio Value	1	Normalized portfolio value
Drawdown	1	Current drawdown level
Total	$n \times L + 2n + 2$	382 dimensions

B.3 State Space Specification

The state vector consists of the following components:

With $n = 18$ assets and $L = 20$ days: $18 \times 20 + 2 \times 18 + 2 = 398$ dimensions.

B.4 Action Space Specification

Table 23: Action Space Components

Component	Dimension	Range
Asset Weights	$n = 18$	$[0, 1]$
Hedge Ratio	1	$[0, 0.2]$
Total	19	Softmax normalized

B.5 Hardware and Software Specifications

Table 24: Computational Environment

Component	Specification
Operating System	macOS
CPU	Apple M-series
RAM	16+ GB
Python Version	3.13.3
PyTorch Version	2.x
Stable-Baselines3	2.x
Gymnasium	0.29.x
NumPy	1.26.x
Pandas	2.x

Table 25: Training Resource Requirements

Metric	DDPG	PPO
Training Time	~45 min	~30 min
Peak Memory	2.1 GB	1.8 GB
Model Size	2.4 MB	2.1 MB
Replay Buffer	800 MB	N/A

B.6 Training Time and Resources

B.7 Reproducibility Checklist

To reproduce our results:

1. Clone the repository
2. Install dependencies: `pip install -r requirements.txt`
3. Set random seed: `seed = 42`
4. Run training: `python scripts/train_final_benchmark.sh`
5. Run evaluation: `python scripts/evaluate_final_models.py`
6. Generate visualizations: `python scripts/visualize_benchmark_comparison.py`

B.8 Data Preprocessing Steps

1. Fetch adjusted close prices from Yahoo Finance
2. Forward-fill missing values (holidays, gaps)
3. Calculate daily logarithmic returns
4. Compute 20-day rolling volatility
5. Normalize features to zero mean, unit variance
6. Split into training (2010-2018) and test (2019-2020) sets