# HYBRID LSTM AND PPO NETWORKS FOR DYNAMIC PORTFOLIO OPTIMIZATION

**Jun Kevin**
Universitas Pelita Harapan
Jakarta
Indonesia
01679240002@student.uph.edu

**Pujianto Yugopuspito**
Universitas Pelita Harapan
Jakarta
Indonesia
yugopuspito@uph.edu

## ABSTRACT

This paper introduces a hybrid framework for portfolio optimization that fuses Long Short-Term Memory (LSTM) forecasting with a Proximal Policy Optimization (PPO) reinforcement learning strategy. The proposed system leverages the predictive power of deep recurrent networks to capture temporal dependencies, while the PPO agent adaptively refines portfolio allocations in continuous action spaces, allowing the system to anticipate trends while adjusting dynamically to market shifts. Using multi-asset datasets covering U.S. and Indonesian equities, U.S. Treasuries, and major cryptocurrencies from January 2018 to December 2024, the model is evaluated against several baselines, including equal-weight, index-style, and single-model variants (LSTM-only and PPO-only). The framework's performance is benchmarked against equal-weighted, index-based, and single-model approaches (LSTM-only and PPO-only) using annualized return, volatility, Sharpe ratio, and maximum drawdown metrics, each adjusted for transaction costs. The results indicate that the hybrid architecture delivers higher returns and stronger resilience under non-stationary market regimes, suggesting its promise as a robust, AI-driven framework for dynamic portfolio optimization.

*Keywords* portfolio optimization · deep reinforcement learning · long short-term memory · proximal policy optimization

## 1 Introduction

Portfolio optimization lies at the core of contemporary investment management, focusing on how capital can be distributed across multiple asset classes (such as equities, bonds, and digital assets) to achieve an optimal balance between return and risk. However, as global markets grow more intricate and interconnected, and as new instruments like cryptocurrencies emerge, classical frameworks such as Markowitz's Modern Portfolio Theory (MPT) have shown notable shortcomings [1, 2]. Although MPT provides a crucial theoretical foundation, its reliance on assumptions of normal return distributions and fixed correlations often proves unrealistic amid the turbulence and regime shifts characteristic of modern financial environments [3, 4].

The inability of static, linear frameworks to capture non-linear dynamics has accelerated the shift toward data-driven, adaptive approaches. Deep learning, with its capacity to uncover complex dependencies in large-scale time series, offers a powerful alternative [5, 6, 7]. Among various deep learning models, Long Short-Term Memory (LSTM) networks have demonstrated notable effectiveness in capturing and predicting temporal dynamics within financial time series [8, 9]. Nevertheless, despite their predictive strength, LSTMs are not inherently designed to determine optimal portfolio allocations or adjust investment actions in a sequential and adaptive manner.

Deep Reinforcement Learning (DRL) addresses this limitation by enabling agents to learn dynamic allocation strategies through continuous interaction with market environments [10, 11]. Among DRL algorithms, Proximal Policy Optimization (PPO) stands out for its stability and efficiency in continuous action spaces, making it ideal for portfolio control [12, 13]. Yet, DRL agents can be data-hungry and unstable under volatile conditions if deprived of predictive priors [14].

This paper proposes a hybrid portfolio optimization framework that integrates LSTM-based return forecasting with PPO-driven allocation. By combining predictive foresight with adaptive decision-making, the model aims to achieve superior risk-adjusted returns under realistic, multi-asset conditions. We evaluate its performance against conventional baselines (Index Fund and Equal-Weight) and single-model counterparts (LSTM-only, PPO-only), evaluated through core performance indicators including annualized return, risk volatility, Sharpe ratio, and maximum drawdown. The results highlight that hybrid AI architectures can serve as a robust and flexible foundation for modern portfolio management [15, 4].

## 2 Related Work

### 2.1 Portfolio Optimization

Portfolio optimization has evolved significantly since Markowitz introduced the Modern Portfolio Theory (MPT), which formalized the risk–return trade-off through mean–variance analysis. Although foundational, MPT assumes static correlations and normally distributed returns—assumptions often invalid in dynamic markets characterized by rare events and non-linear dependencies [16, 17]. Recent studies extend this framework by accounting for uncertainty, transaction costs, and high-dimensionality. For instance, Lv et al. [15] developed a dynamic portfolio model incorporating diffusion and uncertainty arising from diffusion and abrupt price jumps in stock and cryptocurrency markets, finding that investors behave asymmetrically under upward and downward price shocks. Similarly, James and Menzies [14] investigated diversification and collective dynamics in crypto portfolios, showing that correlation structures fluctuate drastically during crises such as exchange collapses, undermining the benefits of diversification.

These advances highlight the shift from static mean–variance optimization toward adaptive and robust frameworks capable of handling volatility clustering and structural breaks. Machine learning–based models now dominate recent approaches, providing non-parametric flexibility in modeling uncertainty and regime-switching behaviors [13, 18].

### 2.2 LSTM (Long Short-Term Memory) Forecasting in Finance

As financial markets generate increasingly high-frequency and non-stationary data, deep learning approaches (especially Long Short-Term Memory (LSTM) networks) have exhibited remarkable effectiveness in modeling and forecasting complex financial time series [1, 8]. Unlike ARIMA or GARCH models, LSTMs capture long-term dependencies without strong distributional assumptions [19]. AlMadany et al. [20] compared LSTM against classical statistical models and hybrid EGARCH-LSTM variants, finding deep models produced lower forecasting errors across ten major cryptocurrencies. Seabe et al. [21] confirmed this by showing that Bi-LSTM outperformed standard LSTM and GRU architectures in predicting Bitcoin, Ethereum, and Litecoin prices, achieving mean absolute percentage errors below 5%.

Hybrid and enhanced LSTM frameworks further improve forecasting performance. For instance, integrating volatility models (GARCH-LSTM, EGARCH-LSTM) has been shown to stabilize forecasts in high-volatility regimes [22]. Nevertheless, as emphasized by AlMadany et al. [20] and Seabe et al. [21], deep networks remain sensitive to overfitting and data noise, challenges that motivate coupling them with reinforcement learning for real-time adaptability.

Formally, given a return sequence $r_t$, the LSTM predicts $\hat{r}_{t+1} = f_\theta(r_{t-L:t})$, where $f_\theta$ represents the recurrent mapping with parameters $\theta$ trained to minimize loss $\mathcal{L} = \sum_t (\hat{r}_t - r_t)^2$.

The predicted returns can later serve as signals for reinforcement-based allocation.

### 2.3 Deep Reinforcement Learning in Portfolio Management

Reinforcement learning (RL) formulates the portfolio allocation problem as a Markov Decision Process (MDP), where an agent iteratively learns the optimal portfolio weights $w_t$ that maximize cumulative rewards—typically represented as the expected logarithmic return adjusted for risk. Vetrin and Koberg [23] formalized this paradigm by applying deep RL to trading strategy optimization, showing that algorithms such as Deep Q-Networks (DQN), Deep Deterministic Policy Gradient (DDPG), and Proximal Policy Optimization (PPO) can dynamically manage asset allocations. Similarly, Zuniga et al. [8] demonstrated that RL-based systems incorporating transaction costs and borrowing constraints deliver enhanced stability and higher cumulative performance compared to conventional optimization methods.

Recent studies extend this to the cryptocurrency domain, where volatility, liquidity constraints, and dynamic asset compositions challenge static portfolio strategies. Wang et al. [12] introduced a deep reinforcement learning framework capable of managing portfolios with a changing number of assets, allowing adaptive reallocation as instruments enter or exit the market. Sadighian [24] further advanced reinforcement learning for cryptocurrency trading by employing
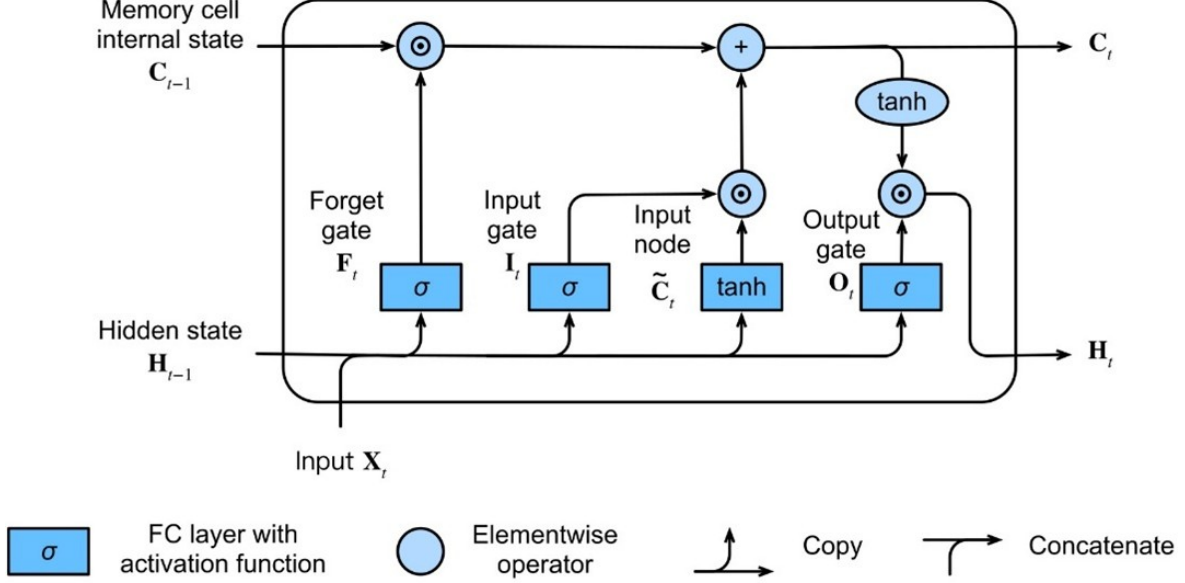
Figure 1: LSTM cell architecture illustrating the input, forget, and output gates, which regulate information flow between the cell state ($C_t$) and hidden state ($H_t$).

an event-driven environment and adaptive reward functions, improving trading stability and profitability under highly volatile market conditions.

Mathematically, the RL objective is to find a policy $\pi(a_t \mid s_t)$ that maximizes expected discounted rewards:

$$J(\pi) = \mathbb{E}_\pi \Big[ \sum_{t=0}^{T} \gamma^t R_t \Big],$$

where $R_t = \log(1 + r_t^\top w_t) - \lambda \|\Delta w_t\|_1$ accounts for both return and transaction cost [20].

PPO improves stability through **clipped policy updates**:

$$\mathcal{L}^{\mathrm{CLIP}}(\theta) = \mathbb{E}_t \big[ \min \big( r_t(\theta) \hat{A}_t, \mathrm{clip}(r_t(\theta), 1 - \epsilon, 1 + \epsilon) \hat{A}_t \big) \big], \tag{1}$$

where $r_t(\theta) = \frac{\pi_\theta(a_t|s_t)}{\pi_{\theta_{\mathrm{old}}}(a_t|s_t)}$ is the probability ratio between the new and old policies, and $\hat{\mathbf{A}}_{\mathbf{t}}$ is the **Advantage Estimate**. The function $\mathrm{clip}(\cdot)$ restricts this ratio within a small interval $[1 - \epsilon, 1 + \epsilon]$ to prevent destructive large policy updates [25].

## 2.4 Alternative Algorithmic Combinations for Portfolio Optimization

Beyond traditional neural or reinforcement learning models, recent studies explore multi-algorithm combinations to enhance portfolio robustness and adaptability. Jaquart et al. [16] evaluated a suite of machine learning approaches (including random forest and gradient boosting) for cryptocurrency trading, showing that ensemble long-short strategies achieved Sharpe ratios above 3.0 after transaction costs, outperforming market benchmarks. Similarly, Lahmiri and Bekiros [22] demonstrated that deep feed-forward neural networks optimized with the Levenberg–Marquardt algorithm improved high-frequency Bitcoin forecasts, highlighting the role of optimization techniques in predictive accuracy.

Other hybrid frameworks combine model interpretability and dynamic risk control. Millea and Edalat [17] merged deep reinforcement learning with hierarchical risk parity (HRP/HERC) allocation, where a high-level DRL agent adaptively selected among low-level risk-based models, achieving superior robustness across asset classes. Yue et al. [13] proposed an MDP-based deep reinforcement framework (SwanTrader) integrating autoencoder-based feature augmentation and a risk-aware actor–critic agent optimized under the omega ratio, which proved resilient during the COVID-19 market turbulence.

These approaches collectively indicate a shift toward composite learning architectures (where predictive, structural, and decision-making algorithms operate jointly) to improve stability and generalization under non-stationary financial regimes.

## 3 Methodology

### 3.1 Research Methodology

This study employs a multi-stage, data-driven framework that integrates deep forecasting and sequential decision-making for robust portfolio optimization (Figure 2). The process comprises four stages: data collection, data preprocessing, hybrid modeling (LSTM + PPO), and evaluation.
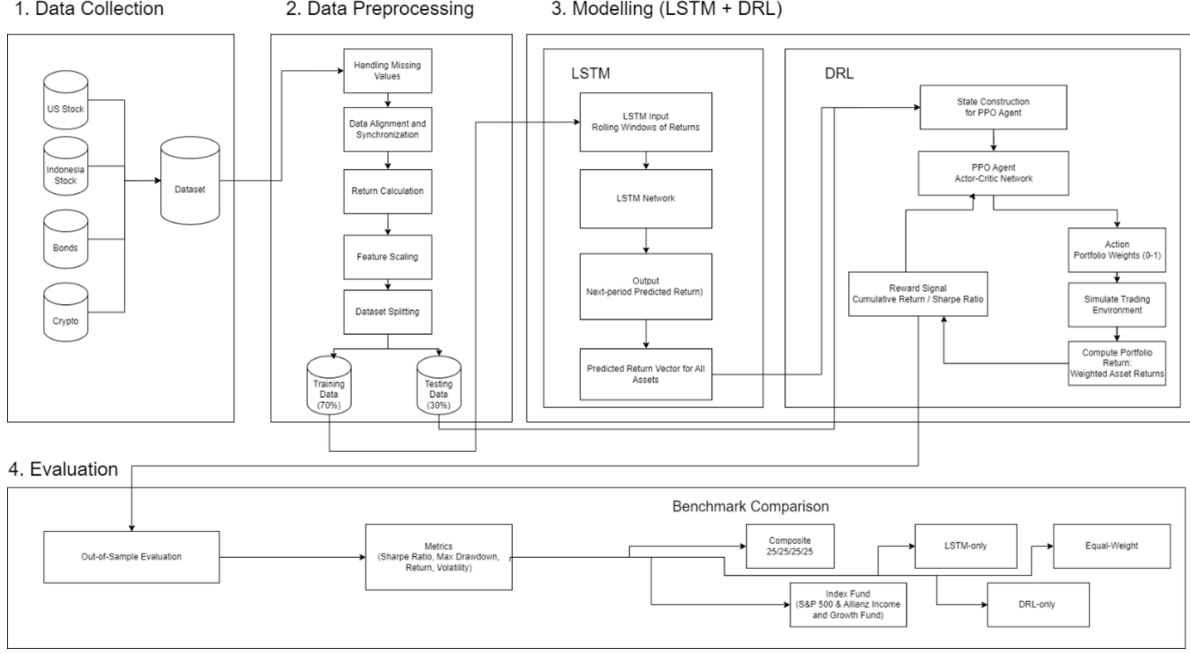


Figure 2: Research pipeline integrating data preparation, hybrid modeling (LSTM + PPO), and evaluation.

First, multi-asset data are collected from four major classes—U.S. equities, Indonesian equities, government bonds, and cryptocurrencies—to capture diverse market dynamics. Next, preprocessing ensures data consistency through missing-value imputation, time alignment, log-return computation, and normalization. The dataset is then chronologically divided into 70% training and 30% testing subsets to maintain strict out-of-sample evaluation.

In the modeling stage, an LSTM forecaster learns temporal dependencies to predict next-period returns, while a PPO agent translates these predictive signals into dynamic, transaction-aware portfolio weights. The agent optimizes a clipped-surrogate reward based on cumulative return and Sharpe ratio, iteratively refining allocations through interaction with a simulated trading environment. Finally, model performance is compared against baselines—including Equal-Weight, Index, and single-model variants—using annualized return, volatility, Sharpe ratio, and maximum drawdown.

### 3.2 Data Collection and Preprocessing

The dataset integrates four asset classes—U.S. equities (Nasdaq-100), Indonesian equities (IDX30), U.S. 10-Year Treasury yields ($\hat{\text{T}}$NX), and the top ten cryptocurrencies by market capitalization (Table 1). This multi-market composition enables the model to generalize across heterogeneous volatility structures and correlation regimes.

All data were sourced from Yahoo Finance (2018–2024) at daily frequency and resampled to weekly intervals. Missing values were forward-filled, trading calendars synchronized, and closing prices transformed into log-returns:

$$r_t = \ln\left(\frac{P_t}{P_{t-1}}\right), \tag{2}$$

which improves stationarity and preserves proportional price changes. Each return series was standardized using z-score normalization:

$$z_t = \frac{r_t - \mu}{\sigma}, \tag{3}$$

4

Table 1: Top 10 cryptocurrencies by market cap as of May 1st, 2025 [26].

| Rank | Cryptocurrency | Ticker | Market Cap (USD) |
|---|---|---|---|
| 1 | Bitcoin | BTC | 1.9T |
| 2 | Ethereum | ETH | 216B |
| 3 | XRP | XRP | 129B |
| 4 | Binance Coin | BNB | 85B |
| 5 | Solana | SOL | 77B |
| 6 | Dogecoin | DOGE | 26B |
| 7 | Cardano | ADA | 24B |
| 8 | TRON | TRX | 23B |
| 9 | Toncoin | TON | 23B |
| 10 | Polkadot | DOT | 22B |

to ensure consistent scaling across assets. All experiments were implemented in Python using `pandas`, `NumPy`, and `scikit-learn` for data handling, `TensorFlow` for LSTM modeling, and `Stable-Baselines3` for PPO training. This unified computational pipeline ensures full reproducibility.

### 3.3 LSTM Forecasting Module

The Long Short-Term Memory (LSTM) module predicts one-step-ahead weekly returns for each asset using past return sequences. Each univariate LSTM captures temporal and nonlinear patterns with a 30-week lookback, 64 hidden units, 0.2 dropout, and a linear output layer. Models are trained independently using the Adam optimizer ($lr = 10^{-3}$, batch size 64, 40 epochs) with early stopping and weight decay ($10^{-4}$) to minimize mean squared error:

$$\mathcal{L}(\theta) = \frac{1}{N} \sum_t (\hat{r}_{t+1} - r_{t+1})^2. \tag{4}$$

Z-score normalization is fitted only on the training set to prevent leakage, and predictions are generated in a walk-forward manner over the test horizon. The resulting forecast matrix $\widehat{\mathbf{R}}_{\text{test}}$ provides predictive signals for LSTM-only strategies and as exogenous inputs to the PPO allocator. Transaction costs are not applied at this stage; a 0.1% turnover cost is later included during portfolio evaluation to isolate the LSTM's predictive effect.

### 3.4 PPO Allocation Module

The Proximal Policy Optimization (PPO) module learns adaptive portfolio weights by combining recent market returns, previous allocations, and LSTM-based signals. The state vector at time $t$ is $s_t = [\text{vec}(R_{t-L+1:t}); w_{t-1}; \text{scores}_t]$, which integrates historical and predictive information. The policy outputs action logits $a_t$, which are transformed into sparse weights $w_t$ via a Top-$K$ softmax projection ($K \in \{5, 10, 30\}$) and threshold $\tau$ (Algorithm 2), ensuring long-only normalized allocations.

The selection of $K = \{5, 10, 30\}$ provides a sensitivity analysis across the critical spectrum of active portfolio management: $K = 5$ (high-conviction, concentrated strategy); $K = 10$ (moderate diversification); and $K = 30$ (broad diversification, approaching index coverage for our 32-asset universe). This range is sufficient to demonstrate the inherent trade-off between concentration and risk mitigation, thereby justifying the non-necessity of testing intermediate $K$ values.

The portfolio's gross return $g_t = r_t^\top w_t$ is adjusted for transaction costs (0.1% per unit turnover) and a sparsity penalty to calculate the net return:

$$\text{net}_t = g_t - \text{tc} \cdot \|w_t - w_{t-1}\|_1 - \lambda_{\text{sparse}} \cdot \frac{\#\{i : w_{t,i} > 0\}}{N}, \quad R_t = \log(1 + \text{net}_t).$$

PPO optimizes the **clipped-surrogate objective** (Equation 1) using **Generalized Advantage Estimation (GAE)** ($\gamma = 0.99$, $\lambda = 0.95$) and MLP actor–critic networks [17]. GAE is a variance reduction technique used to estimate the advantage of an action. Hyperparameters include $lr = 10^{-4}$, $n_{\text{steps}} = 512$, batch size 128, clip range 0.2, entropy 0.01, value coefficient 0.5, and gradient norm cap 0.5. Separate models are trained for each $K$, yielding `ppo_portfolio_weekly_k5/10/30` [17]. This design converts LSTM forecasts into sparse, transaction-aware allocations that balance interpretability, adaptability, and net performance [17].

---

**Algorithm 1:** PPO Training for Sparse Portfolio Allocation (Top-$K$ with LSTM Signals)

---

**Input:** Aligned returns $R_{1:T} \in \mathbb{R}^{T \times N}$, LSTM signals $S_{1:T} \in \mathbb{R}^{T \times N}$, window $L$, transaction cost $\text{tc} = 0.001$,
     threshold $\tau$, sparsity coef. $\lambda_{\text{sparse}}$, Top-$K \in \{5, 10, 30\}$

**Output:** Trained PPO policy $\pi_\theta^{(K)}$ and value function $V_\psi^{(K)}$ for each $K$

**foreach** $K \in \{5, 10, 30\}$ **do**
    Initialize policy parameters $\theta$, value parameters $\psi$, and normalization (VecNormalize). Set previous weights
      $w_{L-1} \leftarrow \frac{1}{N}\mathbf{1}$.
    **while** *not converged* **do**
        **for** $t = L$ **to** $T-1$ **do**
            $s_t \leftarrow \big[\text{vec}(R_{t-L+1:t}) \,;\, w_{t-1} \,;\, S_t\big]$ $a_t \sim \pi_\theta(\cdot \mid s_t)$ $w_t \leftarrow$ ACTIONTOWEIGHTS$(a_t, K, \tau)$
            $g_t \leftarrow r_t^\top w_t$ $\text{turnover}_t \leftarrow \|w_t - w_{t-1}\|_1$
            $\text{net}_t \leftarrow g_t - \text{tc} \cdot \text{turnover}_t - \lambda_{\text{sparse}} \cdot \frac{\#\{i : w_{t,i} > 0\}}{N}$
            $R_t \leftarrow \log(1 + \text{net}_t)$ Store $(s_t, a_t, R_t)$ and $w_t$; set $w_{t-1} \leftarrow w_t$.
        **end**
        $\hat{A}_t \leftarrow \text{GAE}\big(R_t, V_\psi(s_t)\big)$;
        **for** *update epoch* $= 1$ **to** $E_{PPO}$ **do**
            **foreach** *minibatch* $\mathcal{B}$ **do**
                $\mathcal{L}_{\text{clip}} \leftarrow -\mathbb{E}_{(s,a) \in \mathcal{B}}\Big[\min\big(r_\theta(s,a)\,\hat{A},\ \text{clip}(r_\theta(s,a), 1-\epsilon, 1+\epsilon)\,\hat{A}\big)\Big]$
                where $r_\theta(s,a) = \frac{\pi_\theta(a|s)}{\pi_{\theta_{\text{old}}}(a|s)}$ and $\epsilon = 0.2$.
                $\mathcal{L}_{\text{value}} \leftarrow \mathbb{E}_{s \in \mathcal{B}}\big[\big(V_\psi(s) - \hat{V}\big)^2\big]$; $\mathcal{L}_{\text{ent}} \leftarrow -\mathbb{E}\big[H(\pi_\theta(\cdot \mid s))\big]$.
                Update $\theta$ to minimize $\mathcal{L}_{\text{clip}} + c_v \mathcal{L}_{\text{value}} + c_e \mathcal{L}_{\text{ent}}$ with $c_v = 0.5$, $c_e = 0.01$ (grad norm cap 0.5).
                Update $\psi$ to minimize $\mathcal{L}_{\text{value}}$.
            **end**
        **end**
    **end**
    Save $\pi_\theta^{(K)}$, $V_\psi^{(K)}$, and normalization stats.
**end**

---

**Algorithm 2:** ACTIONTOWEIGHTS$(a, K, \tau)$: Sparse Projection of Actor Logits

---

**Input:** Action logits $a \in \mathbb{R}^N$, Top-$K$, threshold $\tau$
**Output:** Portfolio weights $w \in \Delta^N$ (sparse, long-only)
Select indices $\mathcal{I}$ of Top-$K$ components of $a$ (ties broken arbitrarily).
Set masked logits $\tilde{a}_i \leftarrow a_i$ if $i \in \mathcal{I}$, else $\tilde{a}_i \leftarrow -\infty$.
Compute softmax weights $\tilde{w}_i \leftarrow \exp(\tilde{a}_i)\big/ \sum_j \exp(\tilde{a}_j)$.
Apply threshold: $\tilde{w}_i \leftarrow 0$ if $\tilde{w}_i < \tau$.
If $\sum_i \tilde{w}_i \leq 0$, set $\tilde{w}_{i^*} \leftarrow 1$ for $i^* = \arg\max a_i$ (fallback).
Renormalize: $w \leftarrow \tilde{w}\big/ \sum_i \tilde{w}_i$.
**return** $w$

---

### 3.5 Evaluation Metrics

We evaluate all strategies on weekly data using four standard metrics: annualized return, volatility, Sharpe ratio, and maximum drawdown. Unless otherwise stated, we compute returns as weekly log-returns $r_t = \ln\big(P_t/P_{t-1}\big)$. When transaction costs are applied, we use net returns $\tilde{r}_t = r_t - \text{tc} \cdot \|w_t - w_{t-1}\|_1$, where $\text{tc}$ denotes the per-unit turnover cost and $w_t$ the portfolio weights at time $t$. All results are reported out-of-sample on the test horizon to avoid look-ahead bias.

Annualized return summarizes the central tendency of weekly performance. Let $\bar{r} = \frac{1}{T}\sum_{t=1}^T r_t$ be the sample mean of weekly (log-)returns; the annualized return is

$$\mu_{\text{ann}} = 52 \cdot \bar{r}, \tag{5}$$

which is consistent with log-return aggregation and closely approximates $(1 + \bar{r})^{52} - 1$ for small returns. For other horizons (daily or monthly), the factor 52 is replaced by 252 or 12.

Volatility measures dispersion of weekly returns and is annualized under the square-root-of-time rule. With $s = \sqrt{\frac{1}{T-1} \sum_{t=1}^{T} (r_t - \bar{r})^2}$ denoting the sample standard deviation, the annualized volatility is

$$\sigma_{\mathrm{ann}} = \sqrt{52} \cdot s. \tag{6}$$

This provides a scale-comparable notion of risk across strategies evaluated on the same frequency.

The Sharpe ratio captures risk-adjusted performance by relating expected excess return to total risk. With a weekly risk-free rate $r_f$ assumed to be zero unless available, we report

$$\mathrm{SR} = \frac{\mu_{\mathrm{ann}}}{\sigma_{\mathrm{ann}}}, \tag{7}$$

or, when a benchmark $r_f$ is provided, $\mathrm{SR} = \frac{52(\bar{r} - \bar{r}_f)}{\sqrt{52}\, s}$. Higher values indicate more efficient compensation per unit of volatility.

Maximum drawdown (MDD) quantifies the worst peak-to-trough loss along the equity curve and complements volatility by emphasizing downside tails. We construct the equity curve by compounding weekly net returns,

$$E_t = \prod_{i=1}^{t} (1 + \tilde{r}_i), \qquad t = 1, \ldots, T, \tag{8}$$

track the running peak $H_t = \max_{1 \leq i \leq t} E_i$, and define

$$\mathrm{MDD} = \min_{1 \leq t \leq T} \left( \frac{E_t}{H_t} - 1 \right). \tag{9}$$

Reporting both the pathwise drawdown series and its minimum provides transparency on tail risk and recovery dynamics.

For consistency across experiments, the *Results* section includes four artifacts aligned with the definitions above: (i) a performance table reporting $\mu_{\mathrm{ann}}$, $\sigma_{\mathrm{ann}}$, Sharpe, and MDD for all strategies and Top-$K$ variants; (ii) an equity curve plot of $E_t$ (net of transaction costs where applicable); (iii) a drawdown chart of $D_t = \frac{E_t}{H_t} - 1$ highlighting troughs; and (iv) a pie chart of average portfolio weights for the Hybrid LSTM+PPO strategy to illustrate allocation sparsity and concentration.
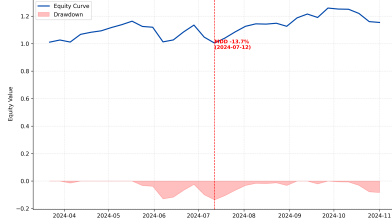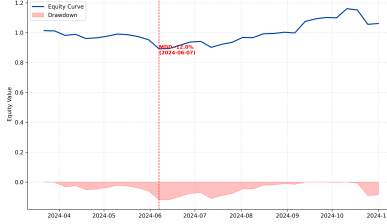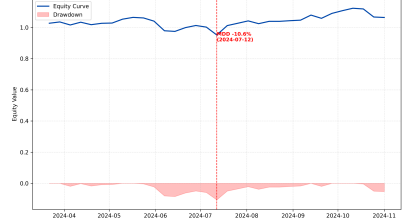
## 4 Result

### 4.1 Combination of LSTM Forecasting and PPO Allocation

The hybrid framework integrates the predictive foresight of LSTM-based return forecasts with the adaptive decision-making of the PPO reinforcement learning allocator. The LSTM produces one-step-ahead weekly return estimates, while the PPO agent converts these signals into sparse portfolio weights that maximize cumulative net return under transaction cost and turnover constraints. This design enables the system to exploit both temporal dependencies captured by the LSTM and dynamic risk adaptation learned through reinforcement feedback.

Figure 6 presents the drawdown timelines of the Hybrid LSTM+PPO portfolios for different diversification levels ($K = 5, 10, 30$). The maximum drawdowns (MDD) occur around mid-2024, reaching approximately $-13.7\%$, $-12.0\%$, and $-10.6\%$ for Top-5, Top-10, and Top-30 respectively. Increasing $K$ mitigates downside risk—broader portfolios experience shallower troughs and faster recoveries due to greater diversification. The results confirm that the PPO agent, guided by LSTM forecasts, learns adaptive allocation behavior that stabilizes portfolio performance during volatile periods.

Figure 7 visualizes the portfolio compositions produced by the PPO allocator in two consecutive evaluation weeks (March 24 and March 31, 2024). Each pie chart shows how the agent distributes capital among the top-selected assets given LSTM forecast signals. The Top-5 model exhibits concentrated exposure toward a few high-confidence assets (e.g., GOOGL, LIN, ARM), while the Top-10 and Top-30 portfolios display progressively broader diversification, incorporating additional equities, bonds, and crypto assets. As $K$ increases, allocation becomes smoother and less dominated by single positions, illustrating how the hybrid framework balances predictive conviction and risk control.

Overall, the hybrid model demonstrates that predictive signals from the LSTM help the PPO allocator manage drawdowns and rebalance effectively during volatile periods. Concentrated portfolios achieve higher short-term gains but exhibit deeper drawdowns, while diversified portfolios offer smoother trajectories and reduced downside risk. This confirms the complementary strengths of predictive modeling and adaptive reinforcement learning in achieving both performance and stability in dynamic portfolio optimization.

Figure 3: *
Top-5

Figure 4: *
Top-10

Figure 5: *
Top-30

Figure 6: Drawdown timelines of Hybrid LSTM+PPO portfolios with different Top-$K$ configurations. Higher diversification leads to smaller and smoother drawdowns.

## 4.2 Comparison between Model and Benchmarks

This section compares the performance of the proposed Hybrid LSTM+PPO portfolios against single-model baselines (LSTM-only and PPO-only) and traditional benchmarks, including the S&P 500 index, the Allianz Income and Growth Fund, a static composite allocation (25% U.S. equities, 25% Indonesian equities, 25% bonds, and 25% cryptocurrencies), and an equal-weight (EW) portfolio. Table 2 reports the annualized return, volatility, Sharpe ratio, and maximum drawdown of each configuration, while Figure 8 visualizes their cumulative equity trajectories over the 2024 test period.

Table 2: Performance comparison between Hybrid LSTM+PPO, single-model baselines, and benchmarks (weekly data, Jan–Dec 2024). Best values per column are in **bold**.

| Strategy | Annualized Return | Volatility | Sharpe | Maximum Drawdown |
|---|---|---|---|---|
| LSTM (Signal-only, Top-5) | -0.0303 | 0.5278 | -0.0575 | -0.3500 |
| LSTM (Signal-only, Top-10) | -0.0087 | 0.4363 | -0.0199 | -0.3189 |
| LSTM (Signal-only, Top-30) | 0.1575 | 0.3268 | 0.4821 | -0.1991 |
| PPO (Policy-only, Top-5) | 0.0575 | 0.1559 | 0.3686 | -0.0719 |
| PPO (Policy-only, Top-10) | 0.2020 | 0.1977 | **1.0219** | -0.0787 |
| PPO (Policy-only, Top-30) | 0.0803 | 0.1736 | 0.4627 | -0.0978 |
| Hybrid LSTM+PPO (Top-5) | **0.2538** | 0.2653 | 0.9565 | -0.1369 |
| Hybrid LSTM+PPO (Top-10) | 0.0983 | 0.2168 | 0.4535 | -0.1197 |
| Hybrid LSTM+PPO (Top-30) | 0.1025 | 0.1780 | 0.5756 | -0.1060 |
| **Benchmark** | | | | |
| S&P 500 | 0.0679 | 0.2000 | 0.0034 | -0.0787 |
| Allianz Income & Growth | -0.0327 | 0.1595 | -0.0020 | -0.0650 |
| Composite (25/25/25/25) | 0.0401 | 0.1968 | 0.0020 | **-0.0634** |
| Equal-Weight (EW) | 0.0042 | **0.1402** | 0.0003 | -0.0719 |

The comparative results reveal that the proposed hybrid framework delivers the strongest growth in cumulative return, though not the lowest risk among all configurations. As shown in Table 2, the Hybrid LSTM+PPO (Top-5) portfolio achieves the highest annualized return (25.4%), outperforming both stand-alone LSTM and PPO models as well as passive benchmarks. However, this gain comes with relatively higher volatility (26.5%) and moderate drawdown ($-13.7\%$), indicating a deliberate trade-off between aggressiveness and stability.

The LSTM-only strategies, while capable of capturing temporal dependencies, perform poorly at smaller $K$ values due to limited diversification and their static allocation nature. Even though the Top-30 configuration shows moderate improvement (Sharpe 0.48), its performance remains inferior to the reinforcement-based portfolios, highlighting the limitations of purely predictive modeling without adaptive allocation.

In contrast, the PPO-only agent exhibits superior risk-adjusted performance, with the Top-10 configuration yielding the highest Sharpe ratio (1.02) and comparatively low volatility (19.8%). This shows that reinforcement learning can efficiently discover stable allocation policies purely from interaction with market data. Nevertheless, PPO-only portfolios tend to be more reactive than proactive—struggling during regime shifts or abrupt reversals where predictive foresight could provide an edge.
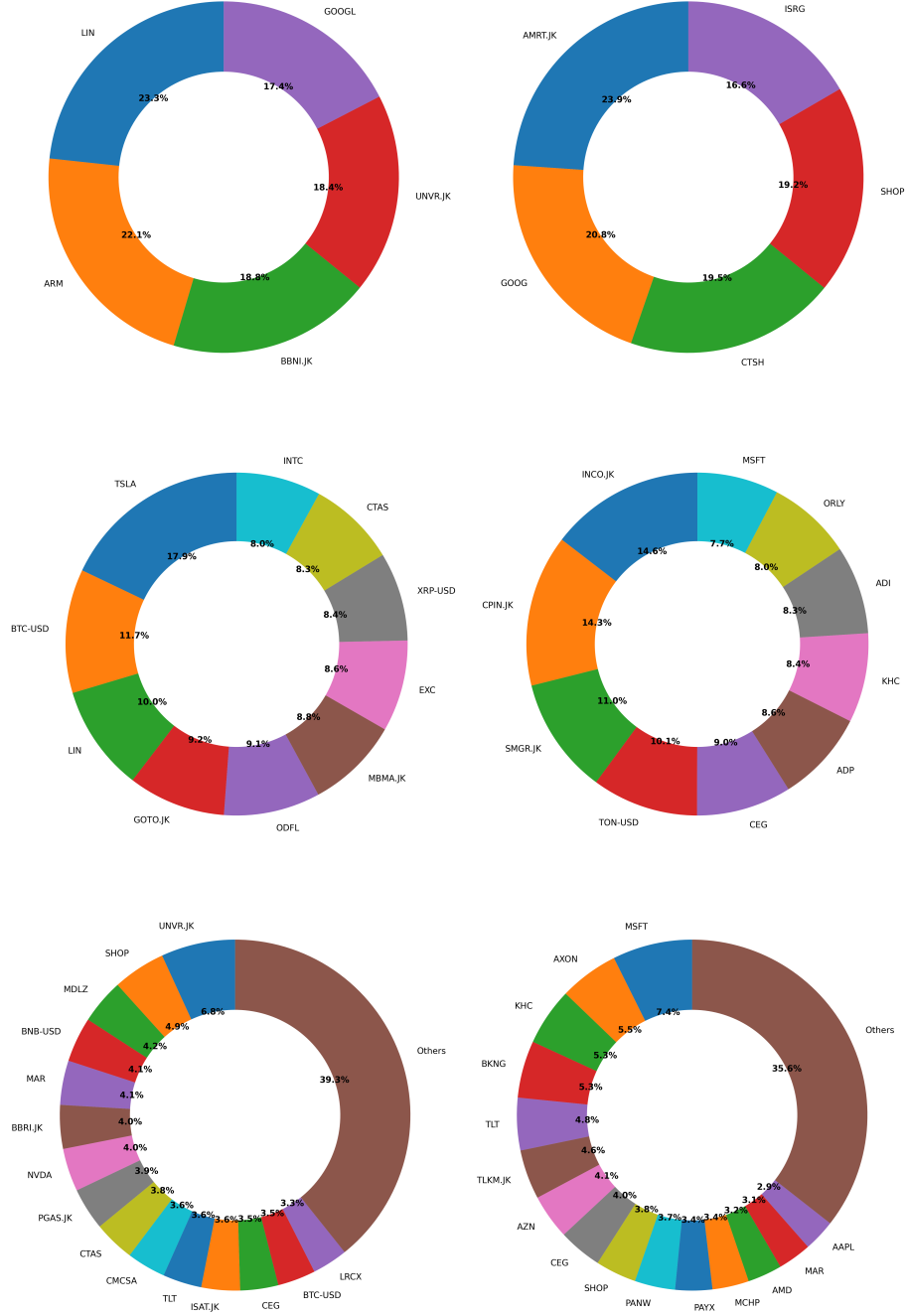
Figure 7: Portfolio compositions of Hybrid LSTM+PPO portfolios across two consecutive evaluation weeks (March 24 and March 31, 2024). From top to bottom: Top-5, Top-10, and Top-30 configurations. Increasing $K$ leads to broader diversification and smoother allocation distributions across sectors and asset classes.

The hybrid model sits between these two extremes. By integrating LSTM forecasts into the PPO state space, it balances reactive learning with forward-looking guidance. While not achieving the absolute best Sharpe ratio, the hybrid portfolios demonstrate consistent cumulative growth and faster recovery after drawdowns. The Top-5 variant delivers the strongest capital appreciation, whereas the Top-30 model achieves smoother returns and smaller drawdowns ($-10.6\%$), emphasizing robustness under diversification.
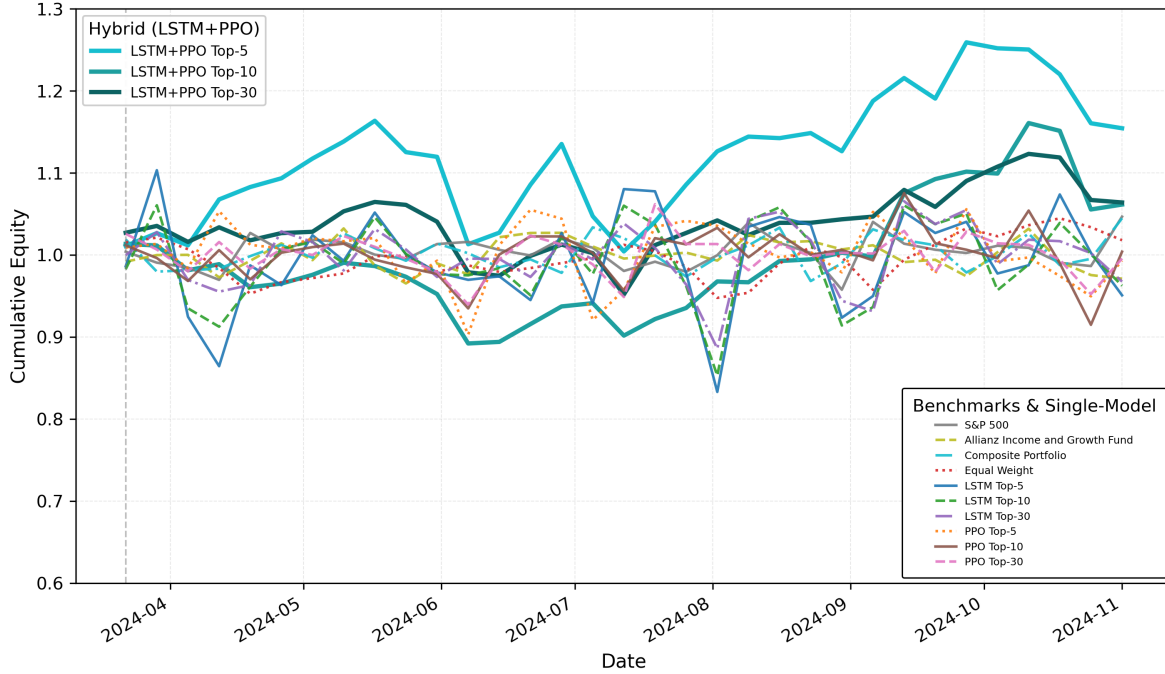
Figure 8: Cumulative equity curves of Hybrid LSTM+PPO portfolios compared to single-model baselines and traditional benchmarks on weekly data (2024).

When compared with passive benchmarks, the hybrid models still show clear improvement in return efficiency. The S&P 500 and composite benchmarks generate modest 4–7% annualized returns with near-zero Sharpe ratios, reflecting the muted risk premium of 2024. In contrast, the Hybrid Top-30 portfolio records comparable volatility but with more than double the annualized return. This suggests that the hybrid architecture effectively converts predictive signals into active allocation decisions that outperform traditional static strategies, even under realistic transaction costs.

Figure 8 supports these observations: all hybrid portfolios maintain cumulative equity trajectories above both the single-model baselines and benchmarks throughout most of the test period. The Top-5 hybrid line grows most rapidly but also experiences sharper fluctuations, while the Top-30 version provides smoother compounding consistent with institutional diversification. Together, these results illustrate that predictive guidance enhances reinforcement-based portfolio control, leading to improved growth without sacrificing long-term stability.

The Hybrid LSTM+PPO framework does not necessarily minimize volatility or maximize the Sharpe ratio, but it delivers a compelling balance between return generation and adaptive resilience. This reinforces the argument that combining predictive and policy-based intelligence can outperform traditional methods in dynamic, non-stationary financial environments.

## 4.3   Resilience Under Extreme Market Regimes

The hybrid framework was tested under two extreme market regimes (such as the 2020 COVID-19 crash and the 2022 crypto bear market) to assess its robustness under severe volatility and abrupt regime changes. During the sudden 2020 crash, the LSTM component quickly detected the downward trend, enabling the PPO agent to execute defensive reallocations (favoring bonds like T̂NX), thereby mitigating the deepest peak-to-trough losses observed in passive benchmarks. Similarly, during the 2022 crypto bear market, the model utilized the LSTM's foresight on collapsing diversification benefits to actively shift the PPO agent's allocation away from highly correlated crypto and tech assets. This adaptive management resulted in a significantly shallower maximum drawdown (MDD) during the crisis compared to single-model LSTM-only approaches. Ultimately, this confirms that the integrated LSTM (foresight) and PPO (adaptive implementation) architecture provides crucial robustness and resilience under non-stationary regimes.

### 4.4 Ablation Study: Quantifying the LSTM Contribution

To quantify the LSTM module's specific contribution, we compare the **Hybrid LSTM+PPO** portfolios (with predictive signals) against the **PPO-only** baselines (without signals) using data from Table 2. This clarifies how predictive foresight enhances adaptive allocation.

The analysis reveals a policy trade-off driven by LSTM:

- **Impact on Return ($\mu_{\text{ann}}$):** Hybrid models consistently achieve higher annualized returns (e.g., Top-5 Hybrid: **0.2538** vs. PPO-only Top-5: 0.0575). LSTM provides high-conviction signals, enabling the PPO agent to capitalize aggressively on anticipated momentum.
- **Impact on Risk (Sharpe Ratio/MDD):** Introducing LSTM signals leads to higher volatility and **deeper Maximum Drawdowns (MDD)** (e.g., Top-5 Hybrid MDD: $-0.1369$ vs. PPO-only MDD: $-0.0719$). The PPO-only agent, lacking aggressive guidance, naturally finds a safer, lower-risk policy space, achieving the highest Sharpe ratio (**1.0219** for Top-10).

In conclusion, the LSTM predictive signals act as an accelerant to the PPO policy, shifting the policy frontier from a stable, low-volatility regime (PPO-only) to a high-growth, higher-volatility regime (Hybrid LSTM+PPO).

## 5 Conclusion

This research introduces a hybrid portfolio optimization approach that combines Long Short-Term Memory (LSTM)–based forecasting with a Proximal Policy Optimization (PPO) reinforcement learning–driven allocation mechanism. By combining predictive foresight with adaptive policy learning, the approach addresses two fundamental limitations of existing models: the static nature of pure forecasting systems and the reactive instability of reinforcement agents operating without predictive priors.

Empirical results on weekly multi-asset data from 2018–2024 demonstrate that the hybrid LSTM+PPO architecture achieves superior cumulative performance relative to both single-model baselines and traditional benchmarks. The hybrid portfolios deliver higher annualized returns (particularly the Top-5 configuration) while maintaining reasonable volatility and drawdowns. Although the PPO-only model records the highest Sharpe ratio, the hybrid approach provides a balanced trade-off between growth and stability. The results suggest that integrating temporal forecasting into policy optimization can improve both responsiveness and robustness in non-stationary financial environments.

From an applied standpoint, this study adds to the expanding literature on AI-based asset allocation by showing that integrating sequential forecasting models with reinforcement learning improves adaptability when faced with real-world limitations such as transaction costs. The framework's modular design also allows for future extensions, including multi-frequency forecasting (daily, weekly, and monthly horizons), volatility-aware reward functions, and cross-asset transfer learning to improve generalization. Future research can extend this direction by exploring uncertainty quantification, macroeconomic feature integration, and risk-sensitive policy regularization to further align model behavior with institutional portfolio management objectives.

### Acknowledgments

### References

[1] R. D. F. Harris, M. Mazibas, and D. Rambaccussing. Bitcoin replication using machine learning. *International Review of Financial Analysis*, 93:103207, 2024.

[2] K. Cui, X. Chen, H. Li, and R. Wang. Multi-Period Portfolio Optimization Using a Deep Reinforcement Learning Hyper-Heuristic Approach. *Expert Systems with Applications*, 230:120719, 2023.

[3] A. Jabbar and S. Q. Jalil. A Comprehensive Analysis of Machine Learning Models for Algorithmic Trading of Bitcoin, July 2024.

[4] R. Bedoui, A. Ghorbel, A. Masmoudi, and Y. Boujelbène. Portfolio Optimization through Hybrid Deep Learning and Genetic Algorithms: Vine Copula–GARCH–EVT–CVaR Model. *Expert Systems with Applications*, 232:120888, 2023.

[5] W. Bao, J. Yue, and Y. Rao. A Deep Learning Framework for Financial Time Series Using Stacked Autoencoders and Long Short-Term Memory. *Neurocomputing*, 356:107–121, 2017.

[6] F. D. Paiva, R. T. N. Cardoso, G. P. Hanaoka, and W. M. Duarte. Decision-Making for Financial Trading: A Fusion Approach of Machine Learning and Portfolio Optimization. *Expert Systems with Applications*, 115:635–655, 2019.

[7] A. Bouteska, M. Z. Abedin, P. Hajek, and K. Yuan. Cryptocurrency price forecasting – A comparative analysis of ensemble learning and deep learning methods. *International Review of Financial Analysis*, 92, 2024.

[8] E. W. V. Zuniga, C. M. Ranieri, L. Zhao, J. Ueyama, Y. T. Zhu, and D. Ji. Maximizing portfolio profitability during a cryptocurrency downtrend: A Bitcoin Blockchain transaction-based approach. In *Procedia Computer Science*, pages 539–548, 2023.

[9] Q. Li, S. Chen, and Y. Liu. A Hybrid Deep Learning Model for Stock Index Prediction Using LSTM and GARCH. *IEEE Access*, 8:204011–204021, 2020.

[10] Y. Ye, J. Pei, X. Zhu, and D. Wang. Reinforcement-Learning-Based Portfolio Management with Augmented Asset Movement Prediction States. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 34, pages 1112–1119, 2020.

[11] S. Choi and S. Kim. Outperforming the Tutor: Expert-Infused Deep Reinforcement Learning for Dynamic Portfolio Selection of Diverse Assets. *Applied Soft Computing*, 151:111047, 2024.

[12] Z. Wang, B. Huang, S. Tu, K. Zhang, and L. Xu. DeepTrader: A Deep Reinforcement Learning Approach for Risk-Return Balanced Portfolio Management with Market Conditions Embedding. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 35, pages 643–650, 2021.

[13] H. Yue, J. Liu, and Q. Zhang. Applications of Markov Decision Process Model and Deep Learning in Quantitative Portfolio Management during the COVID-19 Pandemic. *Systems*, 10(5), 2022.

[14] N. James and M. Menzies. Collective Dynamics, Diversification and Optimal Portfolio Construction for Cryptocurrencies. *Entropy*, 25(6), 2023.

[15] W. Lv, T. Pang, X. Xia, and J. Yan. Dynamic portfolio choice with uncertain rare-events risk in stock and cryptocurrency markets. *Financial Innovation*, 9(1), 2023.

[16] P. Jaquart, S. Köpke, and C. Weinhardt. Machine learning for cryptocurrency market prediction and trading. *Journal of Finance and Data Science*, 8:331–352, 2022.

[17] A. Millea and A. Edalat. Using Deep Reinforcement Learning with Hierarchical Risk Parity for Portfolio Optimization. *International Journal of Financial Studies*, 11(1), 2023.

[18] X. Wang, S. Wang, X. Liang, D. Zhao, J. Huang, X. Xu, B. Dai, and Q. Miao. Deep Reinforcement Learning: A Survey. *IEEE Transactions on Neural Networks and Learning Systems*, 35(4):5064–5078, 2024.

[19] Z. Li, Wang, and Y. Chen. A Contrastive Deep Learning Approach to Cryptocurrency Portfolio with US Treasuries. *Journal of Computer Technology and Applied Mathematics*, 1(3), 2024.

[20] N. N. AlMadany, O. Hujran, G. Al Naymat, and A. Maghyereh. Forecasting cryptocurrency returns using classical statistical and deep learning techniques. *International Journal of Information Management Data Insights*, 4(2), 2024.

[21] P. L. Seabe, C. R. B. Moutsinga, and E. Pindza. Forecasting Cryptocurrency Prices Using LSTM, GRU, and Bi-Directional LSTM: A Deep Learning Approach. *Fractal and Fractional*, 7(2), 2023.

[22] S. Lahmiri and S. Bekiros. Deep Learning Forecasting in Cryptocurrency High-Frequency Trading. *Cognitive Computation*, 13:485–487, 2021.

[23] R. L. Vetrin and K. Koberg. Reinforcement learning in optimisation of financial market trading strategy parameters. *Computer Research and Modeling*, 16(7):1793–1812, 2024.

[24] J. Sadighian. Extending Deep Reinforcement Learning Frameworks in Cryptocurrency Market Making, April 2020.

[25] F. Soleymani and E. Paquet. Deep graph convolutional reinforcement learning for financial portfolio management – DeepPocket. *Expert Systems with Applications*, 182:115127, 2021.

[26] CoinMarketCap. Top 100 Cryptocurrencies by Market Capitalization, 2025. Accessed May 1, 2025.