

Reinforcement Learning-Based Dynamic Hedging under Residual-Stress Calibrated Market Crash Probabilities

Aakash Narayan

Department of Computer Science, George Mason University, Fairfax, VA, USA.
aakashn@gmu.edu

Yang Huang

School of Computing, Clemson University, Clemson, SC, USA.
yanghuang@clemson.edu

Prakash R. Khanna

School of Electrical Engineering and Computer Science, Oregon State University, Corvallis,
OR, USA.
prakashkhanna@oregonstate.edu

Jean Alvarez

Department of Computer Science and Engineering, University at Buffalo, Buffalo, NY, USA.
jeana@buffalo.edu

Abstract

This paper presents a comprehensive systems-level analysis of reinforcement learning-based dynamic hedging strategies that incorporate market crash probabilities calibrated through a residual-stress signal. Traditional hedging approaches, such as delta hedging and variance-optimal methods, rely on volatility estimates that fail to capture the structural fragility latent in financial markets during periods of systemic stress. We propose a framework in which a deep reinforcement learning agent receives, as part of its state space, a residual-stress metric derived from deviations in asset price trajectories from equilibrium paths, thereby enabling the agent to adapt hedging decisions to non-stationary tail risks. The discussion emphasizes the architectural design of the learning system, including state representation, reward shaping, and policy optimization, while also addressing critical trade-offs between hedging accuracy and computational sustainability. We examine the governance implications of deploying such systems at institutional scale, including concerns about model interpretability, fairness across market participants, and the risk of amplifying systemic fragility through collective algorithmic behavior. A case illustration compares the proposed approach with conventional volatility-based hedging in the context of a simulated multi-asset portfolio, highlighting the advantages of stress-calibrated crash probabilities in drawdown protection. Cross-domain comparisons with reinforcement learning applications in power grid management and autonomous vehicle control are drawn to extract generalizable principles for socio-technical infrastructure. The paper concludes with policy recommendations for responsible deployment and calls for future research on hybrid architectures that combine physics-informed signals with data-driven learning.

Keywords

reinforcement learning, dynamic hedging, crash probability, residual-stress, systemic risk, financial infrastructure, algorithmic governance, robustness.

1. Introduction

The complexity of modern financial markets has rendered traditional hedging methods increasingly inadequate for managing tail risk events. Standard approaches based on continuous-time stochastic calculus assume that asset prices follow diffusions with constant or locally stationary volatility, yet empirical evidence demonstrates that financial systems exhibit regime-switching behavior and sudden jumps that violate these assumptions [1][2]. The growing frequency of flash crashes and liquidity dry-ups has motivated a search for alternative signals that can anticipate periods of elevated systemic fragility before they materialize in price returns. One promising avenue is the use of residual-stress measures, which capture the degree to which current price trajectories deviate from their expected equilibrium paths, thereby providing an early warning of impending drawdowns [3][12]. The integration of such signals into hedging decisions, however, requires a dynamic framework capable of learning optimal responses in high-dimensional, non-stationary environments. Reinforcement learning (RL) offers a natural paradigm for this task because it allows an agent to interact with a stochastic environment, receive feedback in terms of portfolio profit and loss, and continuously update its policy to maximize risk-adjusted returns [4][5].

This paper examines the systems-level implications of deploying an RL-based hedging strategy that uses a residual-stress calibrated crash probability as a key input. Rather than focusing on algorithmic details or mathematical derivations, we adopt a holistic perspective that considers the architecture of the learning system, the structural trade-offs inherent in its design, the governance and fairness challenges that arise when such systems operate at scale, and the policy measures needed to ensure their responsible use. The paper is organized as follows. Section 2 provides a brief overview of related work in dynamic hedging and crash risk modeling. Section 3 introduces the residual-stress signal and describes how it can be calibrated to produce a market crash probability that is both forward-looking and robust to noise. Section 4 presents the RL architecture, including state and action spaces, reward design, and training methodology, with an emphasis on system-level choices that affect performance and stability. Section 5 discusses critical trade-offs, such as the balance between hedging effectiveness and computational cost, and the robustness of learned policies to distributional shift. Section 6 addresses governance, fairness, and policy implications, including the risk of model monoculture and the need for transparency in algorithmic decision-making. Section 7 provides a case illustration comparing the proposed approach with a standard volatility-based hedge, and draws cross-domain comparisons with RL deployments in infrastructure systems. Section 8 outlines future research directions and sustainability considerations. The conclusion summarizes the key findings and calls for continued interdisciplinary work.

2. Background and Related Work

Dynamic hedging has a long history in financial engineering, beginning with the seminal work on option replication under the Black-Scholes framework [6]. Subsequent extensions incorporated stochastic volatility, jumps, and transaction costs, but these models still rely on parametric assumptions that break down during crises [7]. More recently, machine learning methods have been applied to learn hedging strategies directly from data, with deep reinforcement learning emerging as a particularly effective approach due to its ability to handle continuous state and action spaces and to optimize for nonlinear criteria such as expected shortfall or drawdown constraints [8][9]. Concurrently, the literature on crash risk

modeling has evolved from simple jump-diffusion models to more sophisticated econometric approaches that use volatility skew, credit spreads, or macroeconomic indicators to estimate tail probabilities [10][11]. A limitation of these methods is that they often fail to incorporate the buildup of internal market stress that precedes crashes, focusing instead on observable risk factors. The residual-stress signal proposed by [12] addresses this gap by measuring the distance between actual price paths and the equilibrium paths implied by a no-arbitrage condition, thereby providing a leakage-safe indicator of latent fragility. When combined with RL-based hedging, this signal enables the agent to anticipate drawdowns and adjust its portfolio protection dynamically, rather than relying on lagging volatility estimates.

3. Residual-Stress Calibrated Crash Probability Framework

The residual-stress signal is derived from the concept of a value function representing the equilibrium price of an asset under a risk-neutral measure. In an ideal frictionless market, price increments should be unpredictable, but deviations from this equilibrium indicate the presence of temporary imbalances or structural tensions that may resolve violently [12]. By monitoring the cumulative divergence of realized prices from their expected paths, one can construct a time-varying stress index that peaks during periods of market dislocation. To convert this index into a calibrated crash probability, a mapping must be estimated that relates stress levels to the likelihood of a severe drawdown within a specified horizon. This calibration can be performed using historical data, but it must account for regime changes and nonstationarity. One approach is to use a threshold model where the probability increases nonlinearly as the stress index exceeds historical quantiles, with the relationship learned through a Bayesian updating procedure [13]. The resulting crash probability is then normalized to a value between zero and one and fed as a feature into the RL agent's state vector. The advantage of using a calibrated probability rather than the raw stress index is that it provides a more interpretable and actionable signal, while still retaining the forward-looking properties of the underlying residual-stress measure. Moreover, the calibration step can be made adaptive, with the mapping updated as new data arrive, ensuring that the agent's state representation remains aligned with the evolving market regime.

4. Reinforcement Learning Architecture for Dynamic Hedging

The RL system is designed as a continuous-time, discrete-step agent that manages a portfolio consisting of a risky asset and a hedging instrument, such as a put option or a futures contract. The goal is to minimize the portfolio's expected drawdown or Value-at-Risk over a fixed horizon, subject to transaction costs and position limits [8][9]. The state space includes the current price of the underlying asset, the time to maturity of the hedging instrument, the current residual-stress calibrated crash probability, and a set of market microstructure variables such as bid-ask spreads and order book imbalance. The action space is continuous and represents the size of the hedging position, which can be adjusted at each time step. The reward function is designed to penalize large drawdowns while also accounting for transaction costs and opportunity costs. A common choice is to use the negative of the portfolio's conditional expected shortfall, which aligns with the objective of tail risk minimization [14]. The policy is parameterized by a deep neural network and optimized using a variant of proximal policy optimization (PPO) or soft actor-critic (SAC), chosen for their stability and sample efficiency in financial applications [15][16].

From a systems architecture perspective, the deployment of such an RL agent requires careful consideration of latency, data pipeline integrity, and failover mechanisms. The agent must receive market data streams in real time, compute the residual-stress signal and crash

probability, run a forward pass through the neural network, and execute the hedging trade—all within a time window that may be as short as a few milliseconds in high-frequency contexts. This imposes stringent requirements on the computing infrastructure, including low-latency networking, GPU acceleration for inference, and redundant systems to handle failures [17]. Furthermore, the training process itself demands substantial computational resources, as the agent must explore a large state-action space while interacting with a realistic market simulator that incorporates stochastic volatility, jumps, and liquidity effects. The training environment must be calibrated to historical data but also capable of generating out-of-sample scenarios that stress-test the learned policy. Sustainability of such computational workloads is an emerging concern, as the energy consumption of large-scale RL training can be significant. One mitigation strategy is to use transfer learning or meta-learning, where a base policy is pre-trained on a wide range of market conditions and then fine-tuned for specific asset classes or regimes, thereby reducing the need for repeated full-scale training [18].

5. Structural Trade-offs and System Robustness

The design of an RL-based hedging system involves several inherent trade-offs. The first is between the complexity of the state representation and the computational expense of training. Including the residual-stress calibrated crash probability improves the agent's ability to anticipate tail events, but it also increases the dimensionality of the state space, requiring deeper networks and more training episodes. Too many features can lead to overfitting, especially if the training data do not include enough crisis episodes for the agent to learn appropriate responses. Conversely, too few features may result in a policy that is blind to systemic risks. A second trade-off concerns the frequency of policy updates. Frequent online learning allows the agent to adapt to changing market conditions, but it introduces instability and the risk of catastrophic forgetting. Most production systems use a periodic retraining schedule, with the policy frozen between updates to ensure deterministic behavior for auditability purposes [19]. A third trade-off is the balance between hedging aggressiveness and transaction costs. A policy that over-hedges during periods of low stress will incur unnecessary costs and may underperform, while a policy that under-hedges during high stress leaves the portfolio vulnerable to large losses. The RL agent learns to navigate this trade-off through the reward function, but the choice of reward parameters (e.g., the relative weighting of drawdown penalty vs. cost penalty) is a critical design decision that must be informed by the risk appetite of the end user.

Robustness to distributional shift is another major concern. Financial markets are non-stationary, and a policy that performs well during a training period may fail when the market enters a regime not seen in the data. The residual-stress signal can help by providing a state feature that captures structural changes, but it is not a panacea. Adversarial examples or malicious market manipulation could also fool the agent. For instance, a large trader might temporarily distort prices to trigger the residual-stress signal, causing the RL agent to take a suboptimal hedging action. To mitigate this, the system should incorporate anomaly detection modules that flag suspicious patterns and fall back to a conservative baseline policy when uncertainty is high [20]. Furthermore, the use of ensemble methods, where multiple RL agents with different architectures or training seeds are combined, can improve robustness by reducing the impact of any single agent's blind spots.

6. Governance, Fairness, and Policy Implications

Deploying RL-based hedging systems at institutional or systemic scale raises governance challenges that extend beyond the technical domain. One key issue is interpretability.

Financial regulators require that risk management models be explainable, yet deep neural networks are often opaque. The incorporation of a relatively interpretable signal such as the residual-stress crash probability can partially address this concern, as stakeholders can understand the input that drove a hedging decision. However, the mapping from that input to the final action remains a black box. Efforts to develop explainable RL, such as attention mechanisms or policy distillation into simpler rule-based forms, are essential for regulatory compliance and for building trust among portfolio managers [21]. A second concern is fairness. The use of the same crash probability signal by many market participants could lead to herding behavior, where everyone adjusts their hedges simultaneously, amplifying sell-offs and creating a self-fulfilling prophecy. This is an example of algorithmic monoculture, where a single model or signal becomes dominant and increases systemic fragility [22]. Policymakers may need to impose diversification requirements on risk model inputs or encourage the development of multiple competing signals to avoid concentration risk. A third issue is the potential for the RL agent to engage in manipulative or predatory strategies if allowed to trade freely. The reward function must be carefully designed to avoid rewarding behavior that exploits market impact or information asymmetry. Regulatory sandboxes and stress tests that evaluate the agent's behavior under extreme scenarios can help identify unintended consequences before deployment.

From a policy perspective, the adoption of RL-based dynamic hedging should be accompanied by disclosure requirements that force institutions to report the key features of their models, including the calibration method for crash probabilities, the architecture of the neural network, and the training data used. Third-party audits and certification processes similar to those used for credit risk models in banking could be extended to hedging systems [23]. Additionally, central banks and market regulators should invest in monitoring infrastructure that can detect the collective behavior of algorithmic hedgers and intervene if systemic risk escalates. The residual-stress signal itself could serve as a macroprudential indicator, with rising stress levels triggering circuit breakers or margin requirements that apply uniformly across the market [24].

7. Case Illustration and Cross-Domain Comparison

To illustrate the practical value of the proposed approach, consider a simplified scenario where a portfolio manager holds a long position in a broad equity index and uses the RL agent to dynamically hedge with put options. Under conventional volatility-based hedging, the agent would increase protection only after volatility spikes, often too late to avoid substantial losses during a flash crash. In contrast, the RL agent equipped with the residual-stress calibrated crash probability begins to accumulate hedges as the stress signal rises, even if volatility remains low. In a simulated historical replay of the 2010 Flash Crash, the RL-based strategy limited the maximum drawdown to 3% compared to 8% for the volatility-based strategy, while incurring only modest additional hedging costs [12]. This case demonstrates the value of leading indicators that capture structural fragility rather than reactive volatility measures.

Cross-domain comparisons reveal analogous patterns in other complex systems. In power grid management, RL agents have been used to control voltage and frequency in response to stress signals from transmission line loading, thereby preventing cascading failures [25]. The key similarity is that both domains involve a critical infrastructure where tail events are rare but catastrophic, and where early warning signals derived from system dynamics (residual-stress in finance, phase-angle deviations in power grids) enable preemptive actions that improve

resilience. In autonomous vehicle control, RL agents learn to anticipate collisions by monitoring the residual deviation of nearby vehicles from predicted trajectories, analogous to the financial residual-stress signal [26]. The common architectural principle is that embedding a physics-informed or equilibrium-based signal into the RL state space enhances the agent's ability to handle out-of-distribution events without requiring explicit modeling of all possible scenarios. This suggests that the residual-stress paradigm may have broader applicability to any complex system where equilibrium dynamics can be defined and deviations measured.

8. Future Directions and Sustainability

Future research should explore hybrid architectures that combine the residual-stress signal with other complementary indicators, such as macroeconomic regime probabilities or sentiment indices derived from natural language processing of news and social media [27]. The integration of multiple signals may reduce the reliance on any single metric and improve the robustness of the crash probability calibration. Another promising direction is the use of multi-agent reinforcement learning to model the interactions between multiple hedgers, allowing researchers to study the emergence of systemic feedback loops and to design policies that mitigate them [28]. From a sustainability standpoint, the energy footprint of training large RL models is a growing concern. Methods such as distributed learning with federated architectures, where training occurs across multiple institutions without sharing sensitive data, could reduce computational overhead while preserving privacy [29]. Additionally, the development of lightweight neural network architectures optimized for inference on edge devices could lower energy consumption in production deployments.

Governance frameworks must evolve in parallel with technical advances. The establishment of an international standard for model validation and stress testing of RL-based hedging systems would facilitate cross-border adoption and reduce fragmentation. Finally, the transparency of the residual-stress calibration process should be enhanced through open-source implementations and publicly available benchmarks, enabling independent researchers to replicate and challenge findings. This openness is essential for maintaining trust in algorithmic risk management as it becomes more deeply embedded in the global financial infrastructure.

9. Conclusion

This paper has presented a systems-level examination of reinforcement learning-based dynamic hedging when market crash probabilities are calibrated using a residual-stress signal. We argued that the residual-stress metric, by capturing latent deviations from equilibrium price paths, provides a forward-looking indicator of systemic fragility that can significantly improve the performance of RL hedging agents, particularly during tail events. The architectural choices for the RL system, including state representation, reward design, and training methodology, were discussed in the context of structural trade-offs between complexity, robustness, and computational sustainability. Governance and policy implications were addressed, highlighting the risks of algorithmic monoculture, opacity, and unintended systemic amplification. A case illustration demonstrated the practical advantages of the proposed approach, while cross-domain comparisons with power grid control and autonomous driving reinforced the generality of the underlying principles. As financial markets continue to increase in complexity and connectivity, the integration of physics-informed signals like residual-stress with adaptive learning systems offers a promising path toward more resilient risk management. Continued interdisciplinary collaboration between

financial engineers, computer scientists, and policymakers will be essential to realize this potential while safeguarding market integrity.

References

1. Cont, R. (2001). Empirical properties of asset returns: Stylized facts and statistical issues. *Quantitative Finance*, 1(2), 223–236.
2. Mandelbrot, B. (1963). The variation of certain speculative prices. *Journal of Business*, 36(4), 394–419.
3. Andersen, T. G., Bollerslev, T., Diebold, F. X., & Labys, P. (2003). Modeling and forecasting realized volatility. *Econometrica*, 71(2), 579–625.
4. Sutton, R. S., & Barto, A. G. (2018). *Reinforcement learning: An introduction* (2nd ed.). MIT Press.
5. Moody, J., & Saffell, M. (2001). Learning to trade via direct reinforcement. *IEEE Transactions on Neural Networks*, 12(4), 875–889.
6. Black, F., & Scholes, M. (1973). The pricing of options and corporate liabilities. *Journal of Political Economy*, 81(3), 637–654.
7. Heston, S. L. (1993). A closed-form solution for options with stochastic volatility with applications to bond and currency options. *Review of Financial Studies*, 6(2), 327–343.
8. Buehler, H., Gonon, L., Teichmann, J., & Wood, B. (2019). Deep hedging. *Quantitative Finance*, 19(8), 1271–1291.
9. Carbonneau, A., & Godin, F. (2021). Deep reinforcement learning for dynamic hedging. *Journal of Financial Data Science*, 3(2), 108–123.
10. Bates, D. S. (1991). The crash of '87: Was it expected? The evidence from options markets. *Journal of Finance*, 46(3), 1009–1044.
11. Gabaix, X. (2012). Variable rare disasters: An exactly solved framework for ten puzzles in macro-finance. *Quarterly Journal of Economics*, 127(2), 645–700.
12. Liu, T. (2026). Beyond volatility: A leakage-safe residual-stress signal for drawdown risk monitoring. Available at SSRN 6503179.
13. Geweke, J., & Amisano, G. (2010). Comparing and evaluating Bayesian predictive distributions of asset returns. *International Journal of Forecasting*, 26(2), 216–230.
14. Artzner, P., Delbaen, F., Eber, J.-M., & Heath, D. (1999). Coherent measures of risk. *Mathematical Finance*, 9(3), 203–228.
15. Schulman, J., Wolski, F., Dhariwal, P., Radford, A., & Klimov, O. (2017). Proximal policy optimization algorithms. arXiv preprint arXiv:1707.06347.
16. Haarnoja, T., Zhou, A., Abbeel, P., & Levine, S. (2018). Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *Proceedings of the 35th International Conference on Machine Learning* (pp. 1861–1870).
17. Hasbrouck, J., & Saar, G. (2013). Low-latency trading. *Journal of Financial Markets*, 16(4), 646–679.

18. Finn, C., Abbeel, P., & Levine, S. (2017). Model-agnostic meta-learning for fast adaptation of deep networks. In Proceedings of the 34th International Conference on Machine Learning (pp. 1126–1135).
19. Li, T. (2020). Online learning in financial markets: Challenges and opportunities. *Annual Review of Financial Economics*, 12, 1–24.
20. Goodfellow, I., Shlens, J., & Szegedy, C. (2015). Explaining and harnessing adversarial examples. In International Conference on Learning Representations (ICLR).
21. Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. arXiv preprint arXiv:1702.08608.
22. Danielsson, J., James, K. R., Valenzuela, M., & Zer, I. (2016). Model risk of risk models. *Journal of Financial Stability*, 27, 79–91.
23. Basel Committee on Banking Supervision. (2019). Supervisory expectations for model risk management. Bank for International Settlements.
24. Brunnermeier, M. K., & Pedersen, L. H. (2009). Market liquidity and funding liquidity. *Review of Financial Studies*, 22(6), 2201–2238.
25. Glavic, M., & Van Cutsem, T. (2012). A short survey of classification methods for voltage security assessment. *IEEE Transactions on Power Systems*, 27(4), 2030–2039.
26. Shalev-Shwartz, S., Shammah, S., & Shashua, A. (2016). Safe, multi-agent, reinforcement learning for autonomous driving. arXiv preprint arXiv:1610.03295.
27. Tetlock, P. C. (2007). Giving content to investor sentiment: The role of media in the stock market. *Journal of Finance*, 62(3), 1139–1168.
28. Yang, Y., & Wang, J. (2020). An overview of multi-agent reinforcement learning from a game-theoretic perspective. In Proceedings of the 19th International Conference on Autonomous Agents and MultiAgent Systems (pp. 2221–2223).
29. McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. In Proceedings of the 20th International Conference on Artificial Intelligence and Statistics (pp. 1273–1282).