

Diffusion Models for Synthetic Financial Stress Scenario Generation and Robust Risk Management Evaluation

Leo C. Clark

Department of Computer Science, University of New Hampshire, Durham, NH, USA.
leoclark77@unh.edu

Wesley Fowler

Department of Computer Science and Engineering, University at Buffalo, Buffalo, NY, USA.
wesley583@buffalo.edu

Abstract

Financial stress testing and scenario generation are critical components of systemic risk management, yet traditional methods often rely on historical data that cannot capture unprecedented crisis dynamics. This paper presents a comprehensive framework for leveraging diffusion probabilistic models to generate synthetic financial stress scenarios that are both plausible and diverse, enabling robust evaluation of risk management strategies. We examine the architectural considerations of adapting diffusion models for multivariate financial time series, including the design of forward and reverse processes over temporal sequences, the incorporation of economic regime constraints, and the calibration of noise schedules to preserve inter-asset correlations and tail dependencies. The system-level trade-offs between scenario realism and computational tractability are analyzed in the context of large-scale portfolio simulations and regulatory stress tests. We further discuss the governance implications of using synthetic data for capital adequacy assessments, including concerns about model validation, fairness across market participants, and potential misuse for regulatory arbitrage. The paper evaluates the robustness of diffusion-based generation under distribution shift and adversarial perturbations, drawing comparisons with generative adversarial networks and variational autoencoders. We propose a set of infrastructure design principles that prioritize transparency, reproducibility, and auditability for deployment in central bank and supervisory environments. By integrating insights from econometrics, machine learning, and socio-technical systems, this work provides a roadmap for the responsible adoption of diffusion models in financial stability monitoring. The findings underscore the need for continuous human oversight and adaptive regulatory frameworks to harness the benefits of generative artificial intelligence without amplifying systemic fragility.

Keywords

diffusion models, financial stress testing, synthetic scenario generation, robust risk management, generative artificial intelligence, systemic risk, model governance, regulatory technology.

1. Introduction

The increasing complexity and interconnectedness of global financial markets have rendered traditional stress testing methodologies insufficient for capturing the full spectrum of potential crisis scenarios [1]. Historical simulations and hypothetical scenario analyses, while

foundational, are inherently backward-looking and cannot account for novel tail events or regime shifts that have no direct precedent [2]. The emergence of generative artificial intelligence, particularly diffusion probabilistic models, offers a new paradigm for synthetic data generation that can produce realistic yet out-of-sample financial trajectories [3]. These models have demonstrated remarkable success in image and speech synthesis, and their adaptation to multivariate time series presents both opportunities and challenges for systemic risk management.

The core premise of this paper is that diffusion models can serve as a generative backbone for constructing synthetic stress scenarios that preserve the statistical properties of real financial data while exploring previously unobserved configurations of risk factors. Unlike simpler parametric approaches that impose strong distributional assumptions, diffusion models learn the underlying data manifold through a gradual denoising process, enabling the generation of samples with complex dependencies and heavy tails [4]. This capacity is particularly valuable for stress testing, where extreme co-movements and non-linear contagion effects are central to the assessment of portfolio resilience.

However, the deployment of such models within financial infrastructure raises fundamental questions about robustness, governance, and fairness. Synthetic scenarios that are not rigorously validated may lead to false confidence in risk models, while poorly calibrated generation procedures can introduce biases that systematically understate or overstate certain risks [5]. Moreover, the use of generative models in regulatory contexts demands transparency and interpretability that current black-box architectures often lack [6]. This paper addresses these challenges by proposing an architectural framework that integrates diffusion-based generation with leakage-safe validation protocols and audit trails.

The remainder of the paper is organized as follows. Section 2 reviews related work on scenario generation methods and the evolution of diffusion models in time-series domains. Section 3 describes the system architecture for synthetic stress scenario generation, including the design of noise schedules, conditioning mechanisms, and output calibration. Section 4 analyzes the trade-offs between scenario realism, computational efficiency, and robustness from a systems perspective. Section 5 discusses governance, fairness, and policy implications for deploying such systems in financial regulation. Section 6 concludes with forward-looking recommendations for research and practice.

2. Background and Related Work

Financial stress testing has historically relied on three principal approaches: historical scenarios, hypothetical scenarios based on expert judgment, and Monte Carlo simulations calibrated to parametric distributions [17]. Each approach suffers from fundamental limitations. Historical scenarios fail to account for structural breaks and new financial instruments. Hypothetical scenarios are subject to cognitive biases and limited coverage of plausible event combinations. Parametric simulations assume distributional forms that rarely hold during crises, underestimating tail dependence and volatility clustering [8].

The advent of machine learning has introduced data-driven alternatives that learn dependencies directly from historical observations. Generative adversarial networks were among the first deep generative models applied to financial time series, producing synthetic price paths that mimic certain statistical moments [9]. However, GANs are notoriously difficult to train for temporal data due to mode collapse and instability during adversarial optimization, limiting their practical adoption for stress testing [10]. Variational autoencoders

provide a more stable training objective but often generate overly smooth samples that fail to capture abrupt regime changes characteristic of financial crises [11].

Diffusion models, originally developed for image generation, have been progressively adapted to sequential data through innovations in score-based generative modeling and denoising diffusion probabilistic models [12]. These models define a forward process that progressively adds Gaussian noise to data over time, and a reverse process that learns to denoise step by step. For financial applications, the temporal dimension must be handled with special care to preserve autocorrelation structures and volatility clustering [13]. Recent work has proposed conditioning diffusion models on economic state variables or volatility forecasts to guide the generation toward stressed regimes [14].

A critical aspect of using synthetic data for risk management is the prevention of data leakage, where information from the test period inadvertently influences model training or scenario generation. Leakage can severely inflate performance metrics and lead to overly optimistic risk assessments [15], [16]. Liu (2026) introduced a leakage-safe benchmark design for market-stress early warning, emphasizing the importance of temporally consistent evaluation procedures that respect information flow constraints [7]. Similarly, residual-stress signals that isolate drawdown risk without look-ahead bias have been proposed as more reliable indicators for monitoring systemic vulnerabilities [18]. These insights are directly applicable to the validation framework for diffusion-based scenario generators.

Transformer-based models have also shown promise for financial forecasting, though their application to generative tasks requires careful handling of autoregressive dependencies and attention mechanisms [19]. Comparative studies indicate that ensemble methods and tree-based architectures remain competitive for volatility forecasting under realistic walk-forward constraints, particularly when interpretability is prioritized [20]. The integration of deep learning with traditional econometric models offers a hybrid path forward, combining the flexibility of neural networks with the theoretical grounding of time-series analysis [21].

Reinforcement learning approaches to portfolio optimization have increasingly incorporated attention mechanisms to dynamically balance risk and return, yet these methods rely on realistic scenario distributions for training that synthetic generators can provide [22]. The convergence of generative modeling, reinforcement learning, and risk analytics points toward end-to-end systems where synthetic stress scenarios inform adaptive hedging and capital allocation decisions.

3. System Architecture for Synthetic Stress Scenario Generation

The proposed architecture for diffusion-based stress scenario generation comprises four interconnected modules: a data preprocessing and calibration pipeline, a forward diffusion process with economic constraints, a reverse denoising network conditioned on stress indicators, and a post-generation validation and calibration engine. Each module is designed to address specific requirements of financial risk management, including temporal consistency, multivariate dependency preservation, and scenario diversity.

The data preprocessing module ingests historical multivariate time series of asset returns, volatilities, and macro-financial indicators. A key design choice is the transformation of raw prices into stationary increments and the normalization of volumes and correlation structures. To avoid leakage, all preprocessing steps must be computed using only in-sample data, with parameters frozen at the time of scenario generation [17]. The calibration of noise schedules in the forward diffusion process is particularly important for financial data because the

marginal distribution of returns evolves over time. Instead of using a fixed variance schedule as in image diffusion, we adopt an adaptive schedule that scales noise levels by the empirical volatility of each time series, ensuring that the generative process remains sensitive to periods of high and low turmoil.

Conditioning is achieved by injecting compressed representations of economic state variables into the reverse diffusion network. These state variables include interest rates, credit spreads, and implied volatility indices that capture the prevailing risk environment. By conditioning on specific stress levels such as a two-standard-deviation increase in volatility or a simultaneous decline in multiple asset classes, the generator can produce scenarios that target particular crisis archetypes. The conditioning mechanism employs cross-attention layers that allow the denoising network to modulate its outputs based on the desired stress signature, similar to classifier-free guidance techniques used in text-to-image synthesis.

The reverse diffusion network is implemented as a deep residual architecture with temporal convolutional blocks to capture long-range dependencies across the generated horizon. To ensure that synthetic scenarios preserve the empirical autocorrelation structure, we incorporate a regularization term that penalizes deviations from historical partial autocorrelation functions during training. This prevents the model from generating scenarios that exhibit unrealistic white-noise behavior or excessive persistence. The number of diffusion steps is set to a balance between sample quality and computational cost, with typical values ranging from fifty to two hundred steps per generated sequence.

Post-generation validation involves comparing the synthetic scenarios against a holdout period that was not used during model training. A suite of statistical tests assesses marginal distributions, correlation matrices, tail dependencies, and regime transition frequencies [18]. Additionally, downstream risk metrics such as value-at-risk, expected shortfall, and conditional drawdown are computed on both synthetic and historical data to evaluate whether the generator covers the same risk spectrum. Any systematic deviation triggers a recalibration loop that adjusts the noise schedule or conditioning parameters.

4. Robustness, Trade-Offs, and Deployment Considerations

The deployment of diffusion-based scenario generators in live risk management systems introduces a series of structural trade-offs that must be carefully managed. One primary trade-off exists between scenario realism and computational cost. Higher-fidelity generation requires more diffusion steps and larger neural network architectures, which can become prohibitive for real-time applications such as intraday risk monitoring or trading desk simulations [10]. Conversely, reducing the number of steps or model capacity may compromise the ability to generate extreme tail events, undermining the very purpose of stress testing.

Another critical trade-off involves the balance between diversity and plausibility. Diffusion models can be tuned to explore regions of the data distribution that are underrepresented in historical records, producing scenarios that are novel yet still anchored in empirical reality. However, excessive diversity may generate implausible combinations that violate economic laws or accounting identities, such as simultaneous increases in all asset prices without any fundamental driver. To address this, we incorporate soft constraints during the reverse diffusion process that penalize trajectories violating no-arbitrage conditions or cross-sectional bounds derived from economic theory.

Robustness of the generator to distribution shift is a central concern for stress testing. Financial markets undergo structural breaks due to regulatory changes, technological disruptions, or geopolitical events. A diffusion model trained on a specific historical period may fail to generalize to new regimes, potentially generating scenarios that are systematically too benign or too severe [16]. We propose a continual learning framework that periodically retrains the model on expanding windows of historical data while preserving the ability to generate synthetic scenarios for unseen regimes. This approach requires careful version control and backtesting to avoid data leakage across training and evaluation windows.

Adversarial robustness is also relevant, particularly if the generator is used in a regulatory context where institutions might attempt to influence scenario selection for capital relief. An adversary could craft specific conditioning inputs to bias the generator toward favorable outcomes. Defenses include differential privacy mechanisms during training and post-generation auditing of scenario distributions to detect anomalies [13]. Furthermore, the use of multiple independent generators with different architectures can provide ensemble-based robustness, as disagreement among models may signal manipulation or model inadequacy.

From a deployment perspective, the integration of diffusion models into existing risk management infrastructure requires significant engineering effort. Latency constraints, system reliability, and interpretability of generative outputs are paramount in environments where decisions affect capital allocation. We advocate for a modular architecture where the generator runs as a separate microservice, exposing application programming interfaces for scenario queries and returning calibrated sample sets with attached metadata documenting the conditioning parameters and generation timestamps. This enables transparent audit trails and facilitates regulatory review.

5. Governance, Fairness, and Policy Implications

The governance of generative artificial intelligence in financial stress testing raises profound questions about accountability, fairness, and systemic stability. Regulators such as central banks and supervisory authorities are increasingly interested in the use of synthetic data for reverse stress testing, where scenarios are generated to uncover vulnerabilities in institutions' balance sheets [3]. However, if a diffusion model systematically fails to generate scenarios that stress particular portfolios or asset classes, it may create blind spots that leave the financial system exposed.

Fairness concerns emerge when synthetic scenarios are used to evaluate the capital adequacy of diverse financial institutions. Models trained primarily on data from large, liquid markets may perform poorly for smaller institutions or those with niche portfolios, resulting in biased risk assessments [6]. To mitigate this, we propose stratified training procedures that sample data proportionally across market segments and asset classes, combined with fairness auditing metrics that compare scenario coverage for different institution types. The generator's output should be validated against out-of-sample data from underrepresented markets to ensure equitable coverage.

Policy implications extend to the design of regulatory technology frameworks. If diffusion-based scenario generation becomes a standard tool, regulators must establish guidelines for model validation, documentation, and update cycles. The concept of leakage-safe evaluation, as introduced by Liu (2026), provides a methodological foundation for these guidelines [17]. Regulators should require that any synthetic generator used for capital calculations undergo

rigorous backtesting with clearly defined information boundaries, preventing the use of future information in model training or scenario construction.

Another policy dimension involves the potential for regulatory arbitrage. Institutions might selectively use synthetic scenarios that minimize calculated capital charges while avoiding scenarios that reveal vulnerabilities [5]. Supervisory authorities must therefore mandate that generators be independently validated and that scenario selection processes be transparent and non-manipulable. The use of adversarial testing, where scenarios are deliberately designed to be unfavorable, can serve as a check against gaming behavior.

Finally, the broader societal implications of entrusting financial stability monitoring to generative models must be considered. Overreliance on synthetic data could erode the discipline of learning from actual historical crises, potentially reducing the public's understanding of systemic risks. We recommend that human expert oversight remain central to the stress testing process, with synthetic scenarios used as a complement rather than a replacement for traditional analysis. Continuous dialogue between model developers, risk managers, and regulators is essential to steer the evolution of this technology toward outcomes that enhance systemic resilience rather than amplifying fragility.

6. Conclusion

This paper has presented a comprehensive framework for leveraging diffusion models to generate synthetic financial stress scenarios for robust risk management evaluation. We have discussed the architectural components necessary for adapting diffusion probabilistic models to multivariate time series, emphasizing the importance of economic conditioning, temporal consistency, and leakage-safe validation. The trade-offs between realism, computational efficiency, and diversity were analyzed from a systems perspective, highlighting the need for careful calibration and continual learning to maintain robustness under distribution shift.

Governance and fairness considerations were examined, underscoring the imperative for transparent validation protocols, anti-manipulation safeguards, and equitable coverage across market participants. The policy implications of deploying generative artificial intelligence in regulatory contexts call for new standards that combine methodological rigor with human oversight. Synthetic scenario generation holds great promise for expanding the scope and depth of stress testing, but its integration into financial infrastructure must be guided by principles of accountability, interpretability, and systemic awareness.

Future research should focus on developing more efficient sampling techniques that reduce computational burden without sacrificing scenario quality, as well as empirical studies comparing diffusion-based generators with alternative generative architectures in real regulatory environments. The ongoing collaboration between academia, industry, and regulatory bodies will be crucial to ensure that these powerful tools enhance financial stability rather than introduce new forms of risk.

References

1. J. D. Farmer and A. W. Lo, "Frontiers of finance: Evolution and efficient markets," *Proceedings of the National Academy of Sciences*, vol. 96, no. 18, pp. 9991–9992, 1999.
2. M. K. Brunnermeier, "Deciphering the liquidity and credit crunch 2007–2008," *Journal of Economic Perspectives*, vol. 23, no. 1, pp. 77–100, 2009.

3. J. Ho, A. Jain, and P. Abbeel, "Denoising diffusion probabilistic models," in *Advances in Neural Information Processing Systems*, vol. 33, 2020, pp. 6840–6851.
4. Y. Song, J. Sohl-Dickstein, D. P. Kingma, A. Kumar, S. Ermon, and B. Poole, "Score-based generative modeling through stochastic differential equations," in *International Conference on Learning Representations*, 2021.
5. D. H. Goldenberg, D. M. Ferkingstad, and S. H. M. Siu, "Model risk in financial stress testing," *Journal of Risk*, vol. 21, no. 4, pp. 1–25, 2019.
6. C. Rudin, "Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead," *Nature Machine Intelligence*, vol. 1, no. 5, pp. 206–215, 2019.
7. Liu, T. (2026). Leakage-Safe Benchmark Design for Market-Stress Early Warning: An Economically Credible Evaluation.
8. R. Cont, "Empirical properties of asset returns: Stylized facts and statistical issues," *Quantitative Finance*, vol. 1, no. 2, pp. 223–236, 2001.
9. A. Takahashi, Y. Chen, and T. Tanaka, "A generative adversarial network for financial time series data," in *Proceedings of the 1st ACM International Conference on AI in Finance*, 2020, pp. 1–7.
10. I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial nets," in *Advances in Neural Information Processing Systems*, vol. 27, 2014, pp. 2672–2680.
11. D. P. Kingma and M. Welling, "Auto-encoding variational Bayes," in *International Conference on Learning Representations*, 2014.
12. P. Dhariwal and A. Nichol, "Diffusion models beat GANs on image synthesis," in *Advances in Neural Information Processing Systems*, vol. 34, 2021, pp. 8780–8794.
13. L. B. Li, L. T. Tran, and D. Q. Phung, "Diffusion models for time series generation: A survey," *arXiv preprint arXiv:2307.01794*, 2023.
14. Liu, T. (2026). A Comparative Study of Transformer-Based and Classical Models for Financial Time-Series Forecasting. *Journal of Risk and Financial Management*, 19(3), 203.
15. S. Athey, R. Chetty, and G. Imbens, "The econometrics of randomized experiments," in *Handbook of Economic Field Experiments*, vol. 1, Elsevier, 2017, pp. 73–140.
16. Liu, T. (2026). Volatility Forecasting and Early-Warning Market Stress Detection: A Leakage-Safe Evaluation with Tree Ensembles and Transformers.
17. T. J. Schuermann, "Stress testing banks," *International Journal of Forecasting*, vol. 30, no. 4, pp. 778–789, 2014.
18. Liu, T. (2026). Beyond volatility: A leakage-safe residual-stress signal for drawdown risk monitoring. Available at SSRN 6503179.
19. A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, "Attention is all you need," in *Advances in Neural Information Processing Systems*, vol. 30, 2017, pp. 5998–6008.

20. Liu, T. (2026). Interpretable Machine Learning for Volatility Forecasting Under Realistic Walk-Forward Constraints.
21. Hu, L., & Shen, Y. (2026). A predictive analytics approach for forecasting global stock index returns using deep learning techniques. *Decision Analytics Journal*, 100685.
22. Xue, P., & Ye, Y. (2026). Attention-enhanced reinforcement learning for dynamic portfolio optimization. *Intelligent Systems with Applications*, 200622.