

The Ethical and Regulatory Challenges of AI-Driven Financial Decision-Making in Global Markets

Fiyinfoluwa Oyesola

University of Southern California, Marshall School of Business – Los Angeles, CA

*fiyinfooluwaoyesola@gmail.com

ARTICLE INFO

Received: 31 Dec 2024

Revised: 20 Feb 2025

Accepted: 28 Feb 2025

ABSTRACT

The integration of artificial intelligence (AI) into financial decision-making has fundamentally transformed global markets, introducing unprecedented efficiency alongside profound ethical and regulatory challenges. This paper presents a narrative review of the ethical and regulatory challenges arising from AI-driven financial decision-making. Drawing on recent literature (2020–2025) across computer science, law, finance, and public policy, it examines five intersecting dimensions: ethical frameworks, regulatory responses, algorithmic bias and explainability, governance architectures, and emerging policy directions. The analysis reveals that existing regulatory frameworks, designed for human decision-makers, remain inadequate for adaptive machine learning systems, while technical opacity continues to undermine fairness and accountability. The paper's main contribution is in pulling these areas together and identifying where they interact, particularly between lifecycle aligned regulation and multi layered governance, which prior reviews have largely treated separately. It offers policy recommendations for regulators, financial institutions, and technology developers, with a focus on cross border harmonization and domain specific regulatory design.

Keywords: Artificial Intelligence, Financial Decision-Making, Ethics, Regulation, Algorithmic Trading, Governance, Explainability, Bias, Financial Markets

1. INTRODUCTION

The rapid proliferation of artificial intelligence technologies in financial services represents one of the most significant technological disruptions in modern economic history. From high-frequency algorithmic trading to AI-powered credit scoring systems and robo-advisors, machine learning algorithms now influence trillions of dollars in daily transactions across global markets (Maple et al., 2023). This transformation promises substantial benefits, including enhanced market efficiency, improved risk management, reduced operational costs, and democratized access to financial services. However, these advances arrive with formidable ethical and regulatory challenges that threaten to undermine market integrity, consumer protection, and social equity if left unaddressed.

The ethical implications of AI in finance extend beyond traditional concerns about data privacy and security. Contemporary AI systems raise fundamental questions about fairness, accountability, transparency, and human agency in financial decision-making (Owolabi et al., 2024). When algorithms determine creditworthiness, execute trades, or provide investment advice, they embed values and assumptions that may perpetuate historical biases, disadvantage vulnerable populations, or operate in ways that even their developers cannot fully explain or control (Veldurthi, 2025). The "black box" nature of many machine learning models creates accountability gaps where harmful outcomes occur without clear lines of responsibility. Regulatory frameworks designed for human decision-makers prove inadequate when confronted with adaptive, autonomous AI systems. Traditional liability regimes based on intent and causation struggle to address algorithmic market manipulation that emerges from machine learning without explicit programming (Azzutti et al., 2023). Supervisory authorities face jurisdictional fragmentation, technical capacity constraints, and the challenge of regulating innovation without stifling beneficial development (Rodríguez de las Heras Ballell, 2023). The global nature of financial markets further complicates regulatory efforts, as divergent national approaches create arbitrage opportunities and coordination challenges.

Technical characteristics of AI systems, particularly opacity, emergent behavior, and data dependency, create unique governance challenges. Algorithmic bias can systematically disadvantage protected groups in lending, insurance, and employment decisions (Chopra, 2024). The interconnectedness of AI-driven trading systems may amplify market volatility and create new vectors for systemic risk (Ekundayo, 2024). Meanwhile, the concentration of AI capabilities among a few technology providers raises concerns about market power, competitive dynamics, and technological dependency (Petronijević et al., 2024).

This paper provides a structured analysis of these intersecting challenges, synthesizing recent scholarship to map the ethical and regulatory landscape of AI in finance. We examine five critical dimensions: (1) ethical frameworks and principles for AI governance; (2) regulatory challenges and policy responses; (3) technical issues of bias, fairness, and explainability; (4) governance structures and risk management approaches; and (5) future trends and emerging policy directions. Our analysis reveals that addressing these challenges requires coordinated action across technical, legal, and institutional domains, balancing innovation incentives with robust safeguards for market integrity and social welfare.

While prior reviews have addressed individual dimensions of this landscape such as algorithmic fairness (Veldurthi, 2025), regulatory design (Azzutti, 2024), or systemic risk (Maple et al., 2023), the present paper seeks to contribute an integrative synthesis that examines how ethical, regulatory, technical, and governance challenges interact and compound one another. In particular, the analysis highlights two underexplored interdependencies: first, the relationship between AI lifecycle characteristics and the design of effective regulatory instruments, building on but extending Azzutti et al.'s (2023) lifecycle framework; and second, the tension between multi-layered governance architectures (Kurshan et al., 2025) and the practical demands of cross-border harmonization (Ridzuan et al., 2024). By mapping these connections across the five dimensions, the paper aims to offer a more holistic view of the governance challenge than has been provided in prior single-dimension treatments.

The remainder of this paper proceeds as follows. Section 2 describes the methodology adopted for the literature review. Section 3 examines ethical frameworks and principles developed to guide AI implementation in financial contexts. Section 4 analyzes regulatory challenges and emerging policy responses across jurisdictions. Section 5 explores technical challenges related to algorithmic bias, transparency, and explainability. Section 6 discusses governance architectures and risk management frameworks. Section 7 identifies future trends and research priorities. Section 8 concludes with policy recommendations and implications for practitioners and policymakers.

2. METHODOLOGY

This study adopted a narrative review methodology to synthesize the existing literature on the ethical and regulatory dimensions of AI-driven financial decision-making. A narrative approach was selected, rather than a systematic review, because the research questions span multiple disciplines—computer science, law, finance, ethics, and public policy—and the relevant literature is distributed across heterogeneous publication venues including peer-reviewed journals, conference proceedings, preprint repositories, book chapters, and institutional reports. A narrative review is well suited to mapping this diverse landscape and identifying cross-cutting themes (Snyder, 2019).

2.1 Search Strategy

Literature searches were conducted across multiple academic databases, including Scopus, Web of Science, Google Scholar, and the Social Science Research Network (SSRN). Search terms included combinations of the following keywords: “artificial intelligence,” “machine learning,” “financial markets,” “algorithmic trading,” “credit scoring,” “robo-advisor,” “ethics,” “fairness,” “bias,” “explainability,” “regulation,” “governance,” and “accountability.” Boolean operators (AND, OR) were used to combine terms across thematic clusters. Searches were supplemented by backward and forward citation tracking of key articles identified in the initial search.

2.2 Inclusion and Exclusion Criteria

Studies were included if they: (a) addressed the intersection of AI/ML and financial services, (b) engaged with ethical, regulatory, or governance dimensions, and (c) were published between 2020 and 2025. Earlier seminal works were included where they provided foundational context. Studies were excluded if they were purely technical in scope (e.g., model performance benchmarking without governance discussion) or addressed AI ethics in non-financial

domains without cross-applicability. No restrictions were imposed on study design; theoretical, empirical, and policy-oriented works were all considered. Given the narrative design, no formal quality scoring rubric was applied to individual sources; however, preference was given to peer-reviewed publications and institutional reports from recognized regulatory bodies.

2.3 Analytical Approach

The included literature was organized thematically according to the five dimensions described in Section 1. Within each thematic category, sources were compared and contrasted to identify areas of consensus, points of disagreement, and underexplored intersections. The synthesis aimed not only to describe the state of scholarship but also to critically evaluate the assumptions, evidence bases, and policy implications of the reviewed works.

2.4 Limitations of the Review

Several limitations of this review should be acknowledged. First, the narrative approach, while suited to interdisciplinary mapping, does not claim the exhaustiveness or replicability of a full systematic review. Second, the literature in this field is evolving rapidly; works published after the search period may address gaps identified herein. Third, the review draws primarily on English-language sources, which may underrepresent perspectives from non-Anglophone jurisdictions, particularly regarding regulatory developments in Asia, Africa, and Latin America. Fourth, the relatively small number of included sources (~20 primary works) reflects the targeted scope of this review rather than the full breadth of available scholarship; future work should expand this base, particularly by incorporating institutional reports from bodies such as the Financial Stability Board, the OECD, and major national regulators. These constraints are taken into account in interpreting the findings.

3. ETHICAL FRAMEWORKS AND PRINCIPLES FOR AI IN FINANCE

3.1 Foundational Ethical Concerns

The deployment of AI in financial decision-making raises ethical concerns that transcend technical performance metrics. At the core lie questions of fairness, transparency, accountability, and respect for human autonomy. Owolabi et al. (2024) present a comprehensive ethical framework tailored specifically to financial AI, emphasizing the need for governance structures that ensure algorithmic decisions align with societal values and legal norms. Their framework identifies four pillars: fairness in outcomes and processes, transparency in model design and operation, accountability for algorithmic harms, and human oversight to preserve meaningful agency.

Fairness encompasses both procedural and distributive dimensions. Procedurally, AI systems must treat similar cases similarly and base decisions on legitimate, non-discriminatory factors. Distributively, outcomes should not systematically disadvantage protected groups or perpetuate historical inequities. Veldurthi (2025) provides an integrated analysis of fairness-aware machine learning techniques, demonstrating how technical interventions at data collection, model training, and deployment stages can mitigate bias. However, Veldurthi emphasizes that technical solutions alone prove insufficient without institutional commitments to equity and ongoing monitoring for disparate impacts.

It is worth noting that Owolabi et al.'s (2024) four-pillar framework and Veldurthi's (2025) technical intervention taxonomy proceed from different disciplinary assumptions. The former is grounded in normative governance theory and treats fairness as an institutional design problem, whereas the latter approaches fairness as a statistical property amenable to algorithmic correction. These perspectives are complementary but can also conflict: a system that satisfies statistical parity at the output level may still rely on opaque processes that violate procedural fairness norms. This tension between outcome-oriented and process-oriented fairness criteria has not been fully reconciled in the financial AI literature, and its resolution likely requires both technical and institutional coordination.

Transparency and explainability present challenges in financial AI. Sharma et al. (2025) propose a modular, ontological assessment framework for ethical evaluation that emphasizes auditable explanations. Their framework recognizes that different stakeholders require different types of explanations, regulators need compliance verification, consumers need decision justifications, and developers need debugging insights. The framework

provides structured methods to evaluate whether AI systems meet legal and ethical criteria while remaining comprehensible to non-technical audiences.

3.2 Balancing Innovation and Ethics

A persistent tension exists between maximizing predictive performance and ensuring ethical operation. Financial institutions face competitive pressures to deploy cutting-edge AI capabilities, potentially compromising ethical safeguards in pursuit of efficiency gains. Islam and Faria (2025) connect ethical governance to environmental, social, and governance (ESG) objectives and sustainable development goals (SDGs), arguing that ethical AI deployment serves long-term institutional interests by building trust and avoiding regulatory sanctions. Their analysis highlights how bias and regulatory gaps can undermine sustainable finance outcomes, creating reputational risks and legal liabilities that outweigh short-term performance advantages.

Chopra (2024) reviews industry practices and governance recommendations, noting significant variation in ethical maturity across financial institutions. Leading organizations implement comprehensive bias testing, maintain model documentation, establish ethics review boards, and invest in explainable AI technologies. However, many institutions lack systematic approaches to ethical risk management, treating ethics as a compliance checkbox rather than an operational imperative. Chopra emphasizes that effective ethical governance requires cultural change, executive commitment, and integration of ethical considerations throughout the AI lifecycle.

When considered together, Islam and Faria (2025) and Chopra (2024) suggest an important gap between aspirational principles and operational reality. Islam and Faria frame ethical AI as strategically advantageous, yet Chopra’s evidence indicates that many institutions remain at an early maturity level, treating ethics primarily as a compliance function. This gap raises the question of whether voluntary adoption of ethical frameworks, absent regulatory compulsion, can achieve the level of consistency and rigor needed to address systemic risks. The literature reviewed here does not provide a definitive answer, though the variation documented by Chopra (2024) suggests that market incentives alone may be insufficient to drive convergence toward high ethical standards.

The financial sector's ethical obligations extend beyond avoiding harm to actively promoting beneficial outcomes. AI systems should enhance financial inclusion, improve access to capital for underserved populations, and contribute to market stability. Yet achieving these positive goals while managing ethical risks requires careful system design, ongoing monitoring, and willingness to forego profitable applications that pose unacceptable ethical risks. Table 1 summarizes key ethical principles and their operational implications for financial AI systems.

Table 1: Ethical Principles and Operational Implications for AI in Finance

ETHICAL PRINCIPLE	DEFINITION	OPERATIONAL IMPLICATIONS	MONITORING MECHANISMS
FAIRNESS	Equal treatment of similar cases; absence of discriminatory bias	Pre-deployment bias testing; disaggregated performance analysis; fairness-aware algorithms	Regular fairness audits; disparate impact analysis; complaint mechanisms
TRANSPARENCY	Understandability of system operation and decision logic	Model documentation; explainable AI techniques; disclosure of AI use	Explanation generation; stakeholder communication; transparency reports
ACCOUNTABILITY	Clear responsibility for algorithmic outcomes	Governance structures; liability assignment; remediation processes	Incident tracking; accountability audits; governance reviews
HUMAN OVERSIGHT	Meaningful human involvement in significant decisions	Human-in-the-loop design; override capabilities; escalation procedures	Oversight effectiveness reviews; decision audits; human factors analysis

PRIVACY	Protection of personal and sensitive information	Data minimization; anonymization; access controls	Privacy impact assessments; data audits; breach monitoring
BENEFICENCE	Promotion of positive outcomes and social welfare	Impact assessment; stakeholder engagement; inclusive design	Outcome monitoring; benefit-risk analysis; stakeholder feedback

3.3 Sector-Specific Ethical Considerations

Different financial applications present distinct ethical challenges. In credit scoring, AI systems must balance predictive accuracy with non-discrimination requirements, avoiding proxies for protected characteristics while maintaining lending standards. In algorithmic trading, ethical concerns center on market manipulation, fairness to other market participants, and systemic stability. Robo-advisors raise questions about suitability, fiduciary duty, and the adequacy of automated advice for complex financial situations.

Each application domain requires tailored ethical frameworks that address specific risks while preserving beneficial uses. Owolabi et al. (2024) emphasize the importance of context-sensitive ethics that recognize differences across financial products, customer segments, and regulatory environments. A one-size-fits-all approach risks either excessive restriction of beneficial applications or inadequate protection against genuine harms. Effective ethical governance thus requires domain expertise, stakeholder input, and iterative refinement based on operational experience.

4. REGULATORY CHALLENGES AND POLICY RESPONSES

4.1 Inadequacy of Traditional Regulatory Frameworks

Existing financial regulations were designed for human decision-makers operating within organizational hierarchies with clear lines of authority and responsibility. These frameworks struggle to address AI systems that learn from data, adapt behavior over time, and operate at speeds that preclude real-time human oversight. Azzutti (2024) examines how the European Union AI Act intersects with algorithmic trading governance, identifying regulatory shortfalls in market integrity and supervisory oversight. The Act's risk-based approach classifies some financial AI as high-risk, triggering transparency, accuracy, and human oversight requirements. However, Azzutti notes that general AI regulations may not adequately address sector-specific concerns about market manipulation, systemic risk, and cross-border coordination.

Azzutti et al. (2023) develop guiding principles for legal reform by analyzing how AI lifecycle characteristics and black-box opacity erode traditional intent-based liability in algorithmic market abuse cases. When trading algorithms discover manipulative strategies through reinforcement learning without explicit programming, existing frameworks cannot clearly assign responsibility. Is the developer liable for failing to prevent emergent behavior? Is the deploying firm responsible for inadequate oversight? Or does the algorithm itself represent a new category of actor requiring novel legal concepts? These questions remain largely unresolved in current law.

Azzutti's (2024) analysis and the earlier work by Azzutti et al. (2023) are noteworthy for their focus on the EU context, but they also reveal a broader structural problem: financial regulation has historically operated on an entity-based model (licensing and supervising identifiable firms), whereas AI introduces activity-based risks that cut across institutional boundaries. A single AI model, developed by a technology vendor, deployed by a bank, and consuming data from a third-party provider, may fall under overlapping or fragmented regulatory mandates. Neither the EU AI Act nor existing financial directives fully resolve this jurisdictional layering, suggesting that regulatory reform must address not only the characteristics of AI systems but also the multi-party value chains in which they are embedded.

4.2 Jurisdictional Approaches and Harmonization Challenges

Different jurisdictions have adopted varying approaches to AI regulation in finance, creating fragmentation that complicates compliance for global institutions and creates regulatory arbitrage opportunities. Rodríguez de las Heras

Ballell (2023) surveys the EU legal regime, which combines general AI regulation with sector-specific financial rules, arguing for principled approaches that ensure responsible automation while preserving innovation. The EU emphasizes fundamental rights, consumer protection, and market integrity, imposing relatively stringent requirements compared to other jurisdictions. The EU's approach is characterized by its emphasis on fundamental rights, consumer protection, and a precautionary stance, reflecting a regulatory philosophy that privileges ex-ante risk prevention.

By contrast, regulatory approaches outside the EU have tended to favour principles-based or innovation-facilitative models. For instance, the discussion in Peterson et al. (n.d.) implies a preference for supervisory flexibility over prescriptive rules, while Ridzuan et al. (2024) document how jurisdictions such as Singapore and the United Kingdom have pursued regulatory sandbox models that allow iterative testing of AI applications under supervisory observation. These divergent philosophies reflect not only different normative commitments but also different institutional capacities: jurisdictions with well-resourced supervisory agencies may be better positioned to implement the monitoring-intensive approaches advocated by the EU model.

Peterson et al. (n.d.) present a multidimensional regulatory approach to balance innovation with market integrity, discussing guidelines and supervisory tools for automated trading. They emphasize the importance of international coordination to prevent regulatory gaps and ensure consistent standards across interconnected markets.

Noguer i Alonso and Chatzianastasiou (n.d.) argue for sector-specific regulatory measures, providing policy rationale and design considerations for effective AI governance in finance. They contend that general AI regulations cannot adequately address the unique characteristics of financial markets, their systemic importance, interconnectedness, and potential for rapid contagion. This position aligns with Azzutti's (2024) critique of the EU AI Act but goes further in its insistence on financial exceptionalism. The argument that the systemic importance and speed of financial markets necessitate dedicated regulatory instruments rather than adaptation of horizontal AI rules. Whether this exceptionalist position is warranted remains debated, as it risks regulatory fragmentation and duplication.

It should be acknowledged that the regulatory analysis presented in this paper draws disproportionately on European sources. This reflects the EU's early-mover status in codifying AI-specific regulation. However, significant regulatory developments have occurred in other jurisdictions—including the United States (e.g., agency-level guidance from the SEC and CFTC on AI in trading), the United Kingdom (the FCA's AI and machine learning in financial services framework), and Asian regulators such as the Monetary Authority of Singapore. A fuller comparative analysis across these jurisdictions represents an important direction for future research.

4.3 Lifecycle-Aligned Regulation

A promising regulatory approach aligns oversight with the AI lifecycle, from data collection and model development through deployment, monitoring, and decommissioning. Azzutti et al. (2023) argue that lifecycle-oriented regulation can address both ex-ante risks (through design requirements, testing, and approval processes) and ex-post harms (through monitoring, enforcement, and liability). This approach recognizes that AI systems evolve over time, requiring ongoing oversight rather than one-time approval. Lifecycle regulation imposes obligations at each stage: data governance requirements during development, validation and testing before deployment, continuous monitoring in operation, and incident reporting when problems arise. This comprehensive approach addresses the dynamic nature of AI while providing clear compliance pathways for developers and deployers. However, implementing lifecycle regulation requires significant supervisory capacity, technical expertise, and coordination across agencies with overlapping jurisdiction.

The lifecycle model, while conceptually compelling, faces a practical tension identified implicitly in the literature: the more granular and stage-specific the regulatory requirements, the greater the compliance burden and the greater the risk of regulatory lag as technologies evolve. Azzutti et al. (2023) acknowledge the capacity demands of lifecycle regulation but do not fully address how resource-constrained supervisory agencies particularly in developing economies can realistically implement such models. This suggests that lifecycle regulation may be most feasible in well-resourced jurisdictions, potentially exacerbating regulatory divergence rather than promoting the harmonization that global financial markets require.

4.4 Enforcement and Compliance Challenges

Even well-designed regulations face enforcement challenges when applied to AI systems. Technical complexity makes it difficult for supervisors to verify compliance, assess model validity, or detect subtle forms of algorithmic manipulation. Financial institutions may lack incentives to report problems that could trigger regulatory sanctions or reputational harm. The speed of algorithmic operations means that significant harm can occur before supervisors detect violations.

Effective enforcement requires investment in supervisory technology, technical expertise within regulatory agencies, and information-sharing mechanisms that provide regulators with visibility into AI operations. Some jurisdictions mandate algorithmic transparency, requiring firms to provide regulators with model documentation, training data, and performance metrics. Others rely on principles-based regulation that sets high-level standards while leaving implementation details to firms, subject to supervisory review.

5. TECHNICAL CHALLENGES: BIAS, FAIRNESS, AND EXPLAINABILITY

5.1 Algorithmic Bias and Discrimination

Algorithmic bias represents one of the most serious ethical and legal challenges in financial AI. Bias can enter systems through training data that reflects historical discrimination, feature selection that incorporates proxies for protected characteristics, or model architectures that amplify existing disparities. Chopra (2024) reviews mechanisms of bias propagation in financial AI, noting that even seemingly neutral variables like zip codes or educational attainment can serve as proxies for race or gender when correlated with protected characteristics.

The consequences of algorithmic bias extend beyond individual unfairness to systematic disadvantage of entire demographic groups. Biased credit scoring systems can deny loans to qualified applicants from minority communities, perpetuating wealth gaps and limiting economic opportunity. Biased fraud detection systems can subject certain populations to elevated scrutiny and account restrictions. These disparate impacts violate anti-discrimination laws and undermine the legitimacy of AI-driven financial systems.

Veldurthi (2025) discusses fairness-aware machine learning techniques that can mitigate bias through technical interventions. Pre-processing approaches modify training data to remove bias, in-processing methods incorporate fairness constraints into model optimization, and post-processing techniques adjust model outputs to satisfy fairness criteria. However, these techniques involve trade-offs between fairness metrics and between fairness and accuracy. Moreover, technical fixes cannot address bias rooted in discriminatory business models or inadequate governance.

A further complexity, not fully addressed in the reviewed literature, concerns the multiplicity of fairness definitions. Veldurthi (2025) and Chopra (2024) discuss fairness in general terms, but the machine learning fairness literature has established that certain fairness criteria such as demographic parity, equalized odds, and calibration are mathematically incompatible under common conditions. This impossibility result has significant implications for financial regulation: a regulator mandating "fairness" without specifying which fairness criterion is operative may create ambiguity that leaves institutions uncertain about compliance. The reviewed sources do not engage substantively with this problem, which represents a gap meriting attention in both the technical and legal literatures.

5.2 Explainability and the Black Box Problem

The opacity of many machine learning models creates accountability gaps and undermines trust. Complex ensemble models, deep neural networks, and other high-performance architectures often function as "black boxes" that produce accurate predictions without comprehensible explanations. This opacity poses problems for multiple stakeholders: regulators cannot verify compliance, consumers cannot understand adverse decisions, and even developers may struggle to debug unexpected behavior.

Sharma et al. (2025) propose frameworks for auditable explanations that balance technical accuracy with stakeholder comprehension. Explainable AI (XAI) techniques like LIME, SHAP, and attention mechanisms provide post-hoc explanations of model predictions, identifying influential features and decision pathways. However, these explanations may not fully capture model behavior, can be manipulated to provide misleading justifications, and may not satisfy legal requirements for meaningful explanations.

The tension between performance and explainability forces difficult choices. Simpler, more interpretable models may sacrifice predictive accuracy, potentially leading to worse outcomes for all stakeholders. Yet deploying opaque models without adequate explanations violates transparency principles and may violate legal requirements. Some jurisdictions mandate "right to explanation" for automated decisions, requiring intelligible justifications that consumers can understand and challenge.

It should be noted that the performance–explainability trade-off may be less stark than commonly assumed. Recent work has suggested that in many structured-data settings typical of finance (e.g., credit scoring with tabular features), interpretable models can achieve performance competitive with opaque alternatives. Sharma et al.'s (2025) framework does not fully engage with this evidence, which, if confirmed across application domains, would weaken the case for deploying black-box models in high-stakes financial contexts and shift the burden of justification onto those who prefer opaque architectures.

5.3 Emergent Behavior and Market Manipulation

AI systems can exhibit emergent behaviors not explicitly programmed by developers, creating novel risks for market integrity. Borch and Min (n.d.) analyze how machine learning-driven trading alters market social dynamics and raises accountability questions for algorithmic conduct. Second-generation trading systems using reinforcement learning can discover complex strategies through trial and error, potentially including manipulative tactics like spoofing, layering, or momentum ignition.

Mizuta (2020) demonstrates via simulation that genetic algorithms can discover market manipulation strategies without human intent, highlighting detection and liability challenges. When algorithms independently learn to manipulate markets, traditional concepts of intent and mens rea become problematic. Legal frameworks based on proving manipulative intent cannot easily address manipulation that emerges from optimization processes without explicit programming. Jukl and Lánský (2025) synthesize algorithmic trading literature, mapping technical mechanisms that create systemic risk. They note that the speed, volume, and interconnectedness of algorithmic trading can amplify volatility, create flash crashes, and propagate shocks across markets.

The 2010 Flash Crash remains a frequently cited illustration of how algorithmic trading can destabilize markets absent malicious intent. However, more recent events underscore that this risk has intensified rather than diminished. Market disruptions in 2015, 2016, and 2020 have each exhibited characteristics associated with algorithmic amplification. The increasing autonomy and complexity of contemporary AI-driven trading systems, as described by Kurshan et al. (2025), suggests that the systemic risk vectors identified in the wake of the 2010 event remain highly relevant and may have grown in severity.

5.4 Data Quality and Model Risk

AI performance depends fundamentally on training data quality, representativeness, and relevance. Poor data quality, including errors, incompleteness, or unrepresentative samples, can degrade model performance and introduce bias. Financial data presents particular challenges: it may be proprietary and difficult to access, subject to survivorship bias, or non-stationary as market conditions evolve. Ul Islam et al. (2024) survey technical advances in machine learning and reinforcement learning for finance, pinpointing data privacy and model-risk management as central technical-policy interfaces. They emphasize that sophisticated models require sophisticated risk management, including validation on out-of-sample data, stress testing under extreme scenarios, and ongoing monitoring for performance degradation. Model risk, the potential for adverse outcomes from incorrect or misused models, represents a persistent challenge that technical sophistication alone cannot eliminate.

6. GOVERNANCE STRUCTURES AND RISK MANAGEMENT

5.1 Multi-Layered Governance Architectures

Effective governance of AI in finance requires coordination across multiple organizational and institutional levels. Kurshan et al. (2025) propose a four-layer modular governance architecture to detect emergent harmful behavior in agentic trading systems. The first layer involves self-regulation by AI systems through built-in constraints and safety mechanisms. The second layer comprises firm-level controls including model validation, risk limits, and human

oversight. The third layer features regulator-hosted monitoring agents that observe market behavior and detect anomalies. The fourth layer consists of independent audits by third parties to verify compliance and assess governance effectiveness.

This multi-layered approach recognizes that no single control mechanism suffices to manage AI risks. Self-regulation provides real-time protection but may be circumvented or fail under novel conditions. Firm-level controls add institutional oversight but face conflicts of interest and resource constraints. Regulatory monitoring provides independent oversight but may lack technical sophistication or real-time responsiveness. Independent audits offer credible verification but occur episodically and may not detect emerging problems.

Kurshan et al.'s (2025) architecture is among the most detailed governance proposals in the reviewed literature, but it rests on assumptions that warrant scrutiny. The model presumes that regulator-hosted AI monitoring agents (Layer 3) can operate with sufficient technical sophistication to oversee the very AI systems they are intended to regulate—an assumption that, as Azzutti (2024) implicitly notes, is not yet realized in most supervisory agencies. Furthermore, the architecture does not address the political economy of independent auditing (Layer 4): who certifies the auditors, how audit standards are set, and how conflicts of interest in the emerging “AI audit” market are managed. These implementation details are critical to the framework’s viability and represent important directions for future research.

Ravishankar (2025) provides practical frameworks for integrating model-risk analytics with cross-border compliance and monitoring. He emphasizes the importance of governance structures that span the AI lifecycle, from development through deployment and decommissioning. Effective governance requires clear roles and responsibilities, adequate resources and expertise, and mechanisms for escalation and remediation when problems arise.

6.2 Risk Management Frameworks

Financial institutions must develop comprehensive risk management frameworks that address AI-specific risks alongside traditional financial risks. These frameworks should identify, assess, monitor, and mitigate risks related to model performance, data quality, algorithmic bias, operational dependencies, and systemic impacts. Ekundayo (2024) examines systemic risks and market manipulation vectors introduced by AI decision systems, emphasizing the need for coordinated governance to mitigate volatility and preserve market stability. Risk management for AI requires new capabilities and processes. Institutions need technical expertise to validate models, assess bias, and interpret performance metrics. They need data governance to ensure quality, representativeness, and compliance with privacy requirements. They need monitoring systems to detect performance degradation, bias drift, and anomalous behavior. And they need escalation and remediation processes to address problems quickly and effectively.

Petronijević et al. (2024) discuss operational dependency risks and human oversight erosion as AI systems become more autonomous. They emphasize strategies to preserve expertise and maintain effective controls against AI-enabled manipulation. As AI systems assume greater decision-making authority, institutions risk losing the human expertise needed to understand, challenge, and override algorithmic recommendations. Maintaining meaningful human oversight requires investment in training, clear escalation protocols, and organizational cultures that empower humans to question algorithmic outputs.

A cross-cutting theme that emerges from the governance literature reviewed here is the risk of what might be termed “governance theatre” which is defined the adoption of formal governance structures (ethics boards, audit processes, monitoring dashboards) that provide the appearance of oversight without the substance. Both Chopra (2024) and Petronijević et al. (2024) document scenarios in which governance mechanisms exist on paper but lack the authority, resources, or expertise to function effectively. This risk is particularly acute where AI governance is grafted onto existing compliance functions without adequate adaptation. The literature suggests that effective governance is not merely a structural question but an organizational-cultural one.

6.3 Organizational Culture and Governance

Technical and procedural controls prove insufficient without supportive organizational cultures that prioritize responsible AI use. Institutions must cultivate cultures where employees feel empowered to raise ethical concerns,

where compliance is valued alongside performance, and where long-term reputation matters more than short-term gains. This requires leadership commitment, appropriate incentives, and mechanisms for reporting and addressing problems without fear of retaliation. Governance structures should include ethics committees, model validation teams, and compliance functions with authority to challenge business decisions. These structures must have adequate resources, access to senior leadership, and independence from business pressures. Regular training, clear policies, and accountability for violations reinforce cultural norms and ensure that governance structures function effectively.

Table 2: Governance Mechanisms and Risk Mitigation Strategies

Governance Layer	Mechanisms	Risk Focus	Mitigation	Implementation Challenges
Self-Regulation (System Level)	Built-in constraints; safety limits; fail-safe mechanisms	Real-time prevention of harmful actions; constraint enforcement		Circumvention by learning; failure under novel conditions; limited scope
Firm-Level Controls	Model validation; risk limits; human oversight; ethics review boards	Comprehensive risk assessment; institutional accountability		Resource constraints; conflicts of interest; expertise gaps
Regulatory Monitoring	Supervisory technology; market surveillance; reporting requirements	Independent oversight; market-wide perspective; enforcement		Technical capacity; coordination challenges; regulatory lag
Independent Audits	Third-party assessments; certification; transparency reports	Credible verification; public accountability		Episodic nature; auditor expertise; cost and access
Industry Standards	Best practices; technical standards; professional codes	Harmonization; knowledge sharing; baseline expectations		Voluntary compliance; lowest-common-denominator standards; limited enforcement

7. FUTURE TRENDS AND EMERGING POLICY DIRECTIONS

7.1 Agentic AI and Autonomous Systems

The next generation of financial AI systems will exhibit greater autonomy, learning capability, and strategic sophistication. Kurshan et al. (2025) discuss risks posed by agentic AI systems that pursue objectives with minimal human direction, potentially discovering novel strategies that violate regulatory expectations or ethical norms. These systems may interact with each other in complex ways, creating emergent market dynamics that no individual algorithm or human operator fully understands or controls. Agentic AI poses fundamental challenges for governance frameworks premised on human decision-making. When AI systems negotiate, collaborate, or compete with minimal human involvement, traditional concepts of agency, intent, and responsibility become strained. Legal and regulatory frameworks must evolve to address these autonomous systems, potentially through new liability regimes, mandatory safety constraints, or restrictions on autonomous operation in critical domains. Jain (2025) surveys advanced learning techniques including deep learning and reinforcement learning, discussing implications for decentralized finance (DeFi), robo-advisors, and emerging market structures. He notes that these technologies enable new financial products and services but also create risks related to manipulation, instability, and exclusion. As AI capabilities advance, governance frameworks must anticipate future risks rather than merely reacting to current problems.

7.2 Cross-Border Harmonization and International Cooperation

The global nature of financial markets requires international cooperation to ensure consistent standards, prevent regulatory arbitrage, and address cross-border risks. However, achieving harmonization faces significant obstacles: divergent legal traditions, different regulatory philosophies, competing economic interests, and varying technical capacities. Ridzuan et al. (2024) provide a comparative analysis of national and regional regulatory approaches, offering practical recommendations for harmonizing governance across jurisdictions while preserving innovation.

International standard-setting bodies like the Financial Stability Board, International Organization of Securities Commissions, and Basel Committee on Banking Supervision play important roles in fostering convergence. These bodies can develop principles, best practices, and technical standards that member jurisdictions adapt to local contexts. However, soft law approaches lack binding force and may be implemented inconsistently. Effective international cooperation requires mechanisms for information sharing, joint supervision of cross-border institutions, and coordinated responses to emerging risks. Regulatory colleges, supervisory cooperation agreements, and mutual recognition arrangements can facilitate coordination while respecting national sovereignty. Yet deep integration remains elusive, particularly when economic interests or regulatory philosophies diverge significantly.

This gets harder when you consider the tension between the lifecycle regulation model from Azzutti et al. (2023) and the multi layered governance architecture proposed by Kurshan et al. (2025). Lifecycle regulation assumes a clear development to deployment pipeline where oversight happens in sequence. Multi layered governance emphasizes monitoring across institutional levels at the same time. Making these models work together across jurisdictions would require agreement on substantive standards and compatible supervisory structures. This coordination problem the literature has flagged but not solved.

7.3 Sustainability and Ethical AI

Growing attention to environmental, social, and governance (ESG) factors in finance intersects with ethical AI concerns. Islam and Faria (2025) connect AI governance to sustainable development goals, noting that biased or opaque AI systems can undermine sustainability objectives by misallocating capital, excluding vulnerable populations, or obscuring environmental and social impacts. Conversely, well-governed AI can advance sustainability by improving ESG data analysis, identifying greenwashing, and facilitating impact measurement. The integration of AI and sustainability creates opportunities and challenges. AI can process vast amounts of ESG data, identify patterns, and support sustainable investment decisions. However, AI systems may perpetuate unsustainable practices if trained on historical data that reflects environmentally or socially harmful business models. Ensuring that AI systems align with sustainability objectives requires explicit design choices, appropriate metrics, and governance structures that prioritize long-term welfare over short-term returns.

7.4 Research Priorities and Knowledge Gaps

Despite growing scholarship on AI ethics and regulation in finance, significant knowledge gaps remain. Maple et al. (2023) offer a sweeping overview of AI opportunities and systemic risks, outlining high-level policy principles and risk-based regulatory approaches for future work. They identify priorities including better understanding of systemic risks from AI interconnections, development of effective XAI techniques for complex models, empirical assessment of bias mitigation strategies, and evaluation of governance mechanisms across different institutional contexts. Future research should examine the effectiveness of different regulatory approaches through empirical analysis and comparative case studies. Longitudinal studies can assess how AI systems evolve over time and whether governance mechanisms remain effective as technologies advance. Interdisciplinary collaboration among computer scientists, economists, legal scholars, and ethicists can generate insights that purely technical or purely legal approaches miss.

This review points to several research gaps worth pursuing. The most pressing is the lack of empirical work testing whether fairness aware machine learning techniques actually work in real financial settings as the most existing evidence comes from simulations. Comparative regulatory research would also help as we still don't have good data on how different jurisdictions' approaches to AI governance have affected innovation, consumer protection, or market stability in practice, and policymakers are making decisions without it. The governance problems created by multi-party AI value chains, where model developers, data providers, deploying institutions, and cloud infrastructure providers all share responsibility are still undertheorized and need attention from both legal and organizational scholars. Finally, not enough work has looked at how AI governance frameworks interact with the financial regulatory structures already in place, like prudential supervision, conduct regulation, and market surveillance. These regimes will inevitably overlap and conflict and understanding where that happens matters for designing governance that actually functions.

Policymakers need evidence-based guidance on critical questions: What regulatory approaches best balance innovation and protection? How can supervisors build technical capacity to oversee sophisticated AI systems? What international coordination mechanisms prove most effective? How should liability be allocated when AI systems cause harm? Answering these questions requires sustained research investment, data sharing, and collaboration among academics, practitioners, and regulators.

8. CONCLUSION

The integration of AI into financial decision-making presents a defining challenge for contemporary governance. While AI offers transformative benefits, enhanced efficiency, improved risk management, and democratized access, it also poses profound ethical and regulatory challenges that existing frameworks struggle to address. This paper has examined these challenges across five critical dimensions: ethical principles, regulatory gaps, technical issues of bias and explainability, governance architectures, and future policy directions. Several key findings emerge from our analysis. First, ethical frameworks for AI in finance must balance multiple values, fairness, transparency, accountability, and efficiency, that sometimes conflict in practice. Achieving this balance requires context-sensitive approaches that recognize differences across financial applications, customer segments, and regulatory environments. Second, traditional regulatory frameworks designed for human decision-makers prove inadequate for adaptive, autonomous AI systems. Lifecycle-aligned regulation that addresses risks from development through deployment and monitoring offers a promising path forward, though implementation requires significant supervisory capacity and technical expertise.

Third, technical challenges related to algorithmic bias, explainability, and emergent behavior create persistent risks that technical solutions alone cannot eliminate. Effective governance requires combining technical safeguards with institutional controls, human oversight, and accountability mechanisms. Fourth, multi-layered governance architectures that coordinate self-regulation, firm-level controls, regulatory monitoring, and independent audits provide more robust protection than any single mechanism. These architectures must be supported by organizational cultures that prioritize responsible AI use and empower individuals to raise ethical concerns. However, the reviewed literature reveals that such architectures face significant implementation challenges, including supervisory capacity gaps, the risk of governance theatre, and unresolved questions about auditor certification and independence.

Fifth, future challenges posed by increasingly autonomous AI systems and the global nature of financial markets require anticipatory governance and international cooperation. Policymakers must look beyond current problems to anticipate risks from agentic AI, develop coordinated responses to cross-border challenges, and align AI governance with broader sustainability objectives.

8.1 Policy Recommendations

Based on our analysis, we offer several policy recommendations for regulators, financial institutions, and technology developers:

For Regulators:

1. Adopt lifecycle-aligned regulatory frameworks that address AI risks from development through deployment and ongoing operation with particular attention to the multi-party value chains through which AI systems are developed, deployed, and maintained. Given the capacity demands of lifecycle regulation, supervisory agencies should prioritize risk-proportionate implementation, focusing initial efforts on high-risk applications such as credit scoring and high-frequency trading.
2. Invest in supervisory technology and technical expertise to enable effective oversight of sophisticated AI systems specifically, regulators should develop or procure algorithmic auditing capabilities, establish dedicated AI supervision units staffed with data science expertise, and explore the use of regulatory technology (RegTech) tools—including, as Kurshan et al. (2025) propose, regulator-hosted monitoring agents—to enhance real-time market surveillance.
3. Develop sector-specific rules that address unique characteristics of financial markets while coordinating with general AI regulations. Following Nogueira and Alonso and Chatzianastasiou (n.d.), such rules should address

financial-specific risks (e.g., systemic contagion, market manipulation) that horizontal AI legislation does not adequately capture, while avoiding regulatory duplication.

4. Establish international cooperation mechanisms to ensure consistent standards and prevent regulatory arbitrage. This should include expanding the mandates of existing bodies (e.g., IOSCO, the FSB) to encompass AI-specific coordination, developing mutual recognition frameworks for algorithmic audit standards, and piloting cross-border supervisory sandboxes for AI-driven financial products.
5. Mandate transparency and explainability requirements appropriate to risk levels and stakeholder needs, specifying which fairness criteria are required for different application domains and providing regulatory guidance on the acceptable parameters of the performance–explainability trade-off.

For Financial Institutions:

1. Implement comprehensive governance structures spanning model development, validation, deployment, and monitoring, ensuring that these structures possess genuine decision-making authority and are not subordinated to commercial functions
2. Cultivate organizational cultures that prioritize ethical AI use and empower employees to raise concerns, drawing on the governance-culture nexus identified in Section 6 and investing in ongoing training programmes that develop both technical AI literacy and ethical reasoning capabilities across business lines.
3. Invest in bias testing, fairness audits, and ongoing monitoring for disparate impacts, adopting multiple fairness metrics (e.g., demographic parity, equalized odds, and calibration) in recognition of their inherent tensions, and documenting the rationale for prioritizing specific criteria in specific contexts.
4. Maintain meaningful human oversight of significant decisions, with clear escalation and override protocols
5. Engage stakeholders including customers, employees, and communities in AI governance

For Technology Developers:

1. Incorporate fairness, transparency, and accountability principles into system design from inception
2. Develop and deploy explainable AI techniques appropriate to financial applications, with particular attention to the needs of regulatory auditors and end consumers, and with transparent documentation of the limitations of post-hoc explanation methods such as LIME and SHAP.
3. Conduct rigorous testing including adversarial testing for bias, robustness, and security
4. Provide clear documentation of model capabilities, limitations, and appropriate use cases
5. Support ongoing monitoring and be responsive to problems identified in deployment

8.2 Implications for Practice

Practitioners implementing AI in financial contexts must navigate complex technical, ethical, and regulatory terrain. Success requires not only technical sophistication but also ethical awareness, regulatory knowledge, and stakeholder engagement. Organizations should view AI governance not as a compliance burden but as a source of competitive advantage, building trust with customers, regulators, and the public. Effective implementation requires cross-functional collaboration among data scientists, risk managers, compliance officers, business leaders, and external stakeholders. Siloed approaches that isolate AI development from governance functions or business strategy from ethical considerations will prove inadequate. Organizations must develop integrated approaches that embed ethical considerations throughout the AI lifecycle and align AI capabilities with institutional values and regulatory obligations.

8.3 Future Outlook

The trajectory of AI in finance remains uncertain, shaped by technological advances, regulatory choices, market dynamics, and social values. Optimistically, thoughtful governance can harness AI's benefits while managing its risks,

creating financial systems that are more efficient, inclusive, and stable. Pessimistically, inadequate governance could lead to discriminatory outcomes, market instability, erosion of trust, and restrictive regulations that stifle beneficial innovation. Realizing the optimistic scenario requires sustained commitment from all stakeholders. Regulators must develop sophisticated, adaptive frameworks that keep pace with technological change. Financial institutions must prioritize responsible AI use even when it conflicts with short-term profits. Technology developers must incorporate ethical principles into system design. And society must engage in informed deliberation about the values and trade-offs embedded in AI-driven financial systems. The challenges are formidable, but so are the opportunities. AI has potential to create financial systems that serve broader populations, allocate capital more efficiently, and manage risks more effectively. Achieving this potential while avoiding serious harms will require wisdom, collaboration, and sustained effort. The choices made in the coming years will shape financial markets and economic opportunity for decades to come.

REFERENCES

- [1] Azzutti, A. (2024). AI governance in algorithmic trading: Some regulatory insights from the EU AI Act. *Social Science Research Network*. <https://doi.org/10.2139/ssrn.4939604>
- [2] Azzutti, A., Ringe, W. G., & Stiehl, H. S. (2023). Regulating AI trading from an AI lifecycle perspective. In *Research Handbook on the Law of Artificial Intelligence* (pp. 287-312). Edward Elgar Publishing. <https://doi.org/10.4337/9781803926179.00019>
- [3] Borch, C., & Min, B. H. (n.d.). Machine learning and social action in markets: From first- to second-generation automated trading. *Economy and Society*.
- [4] Chopra, P. (2024). Ethical implications of AI in financial services: Bias, transparency and accountability. *International Journal of Scientific Research in Computer Science Engineering and Information Technology*, 10(5), 1059-1068. <https://doi.org/10.32628/cseit241051059>
- [5] Ekundayo, F. (2024). Economic implications of AI-driven financial markets: Challenges and opportunities in big data integration. *International Journal of Science and Research Archive*, 13(2), 2311-2325. <https://doi.org/10.30574/ijrsra.2024.13.2.2311>
- [6] Islam, Q. T., & Faria, M. H. (2025). Ethics in AI-driven sustainable finance. In *Advances in Computational Intelligence and Robotics Book Series* (pp. 123-145). IGI Global. <https://doi.org/10.4018/979-8-3693-9684-1.ch006>
- [7] Jain, J. (2025). AI-driven learning in finance. In *Advances in Computational Intelligence and Robotics Book Series* (pp. 178-203). IGI Global. <https://doi.org/10.4018/979-8-3373-3952-8.ch008>
- [8] Jukl, D., & Lánský, J. (2025). Systematic review on algorithmic trading. *Acta Informatica Pragensia*, 14(2), 156-189. <https://doi.org/10.18267/j.aip.276>
- [9] Kurshan, E., Balch, T., & Byrd, D. (2025). The agentic regulator: Risks for AI in finance and a proposed agent-based framework for governance. *arXiv preprint*. <https://arxiv.org/abs/2512.11933>
- [10] Maple, C., Szpruch, L., Epiphaniou, G., Staykova, K. S., & Singh, S. (2023). The AI revolution: Opportunities and challenges for the finance sector. *arXiv preprint*. <https://doi.org/10.48550/arxiv.2308.16538>
- [11] Mizuta, T. (2020). Does an artificial intelligence perform market manipulation with its own discretion? A genetic algorithm learns in an artificial market simulation. In *2020 IEEE Symposium Series on Computational Intelligence (SSCI)* (pp. 2410-2417). IEEE. <https://doi.org/10.1109/SSCI47803.2020.9308349>
- [12] Noguer i Alonso, M., & Chatzianastasiou, F. S. (n.d.). The case for artificial intelligence regulation in the financial industry. *Social Science Research Network*. <https://doi.org/10.2139/ssrn.4831147>
- [13] Owolabi, O., Uche, P. C., Adeniken, N. T., Ihejirika, C., & Islam, R. B. (2024). Ethical implication of artificial intelligence (AI) adoption in financial decision making. *Computer and Information Science*, 17(1), 49-68. <https://doi.org/10.5539/cis.v17n1p49>
- [14] Petronijević, J., Radić, N., & Gavrilović, M. (2024). Dependence on technology and market manipulation as potential risks of using artificial intelligence in finance. *Ekonomski Horizonti*, 26(2), 84-97. <https://doi.org/10.5937/eee24084p>
- [15] Peterson, A., Gray, S., Ramirez, A., Martin, S., & Hu, T. (n.d.). Regulatory challenges in algorithmic and autonomous trading systems. *Conference Proceedings*.

- [16] Ravishankar, S. (2025). Navigating AI/ML risk analytics and compliance: A strategic guide for launching global FinTech products. *European Modern Studies Journal*, 9(5), 89-112. [https://doi.org/10.59573/emsj.9\(5\).2025.9](https://doi.org/10.59573/emsj.9(5).2025.9)
- [17] Ridzuan, N. N., Masri, M., Anshari, M., Fitriyani, N. L., & Syafrudin, M. (2024). AI in the financial sector: The line between innovation, regulation and ethical responsibility. *Information*, 15(8), 432. <https://doi.org/10.3390/info15080432>
- [18] Rodríguez de las Heras Ballell, T. (2023). AI in the financial sector. In *The Cambridge Handbook of Artificial Intelligence: Global Perspectives on Law and Ethics* (pp. 67-89). Cambridge University Press. <https://doi.org/10.1017/9781009334297.004>
- [19] Sharma, A. K., Kyosev, D., & Kunkel, J. (2025). Ethical AI: Towards defining a collective evaluation framework. *arXiv preprint*. <https://arxiv.org/abs/2506.00233>
- [20] Ul Islam, M. I., Ahmad, F., Nissa, V., Ansarullah, S. I., & Nisa, K. U. (2024). The future of machine learning and artificial intelligence in finance. In *Advances in Finance, Accounting, and Economics Book Series* (pp. 512-538). IGI Global. <https://doi.org/10.4018/979-8-3693-8507-4.ch026>
- [21] Veldurthi, A. K. (2025). Ethical considerations in AI-driven financial decision-making. *European Journal of Computer Science and Information Technology*, 13(3), 64-89. <https://doi.org/10.37745/ejcsit.2013/vol13n314964>