# Standardized Threat Taxonomy for AI Security, Governance, and Regulatory Compliance

A Unified Taxonomy of Threat Vectors in Generative and Agentic AI and Machine Learning Systems

Prof. Hernan Huwyler, MBA CPA*

Department of Compliance, Control and Risk Management
IE Executive Education School & IE Law School
Maria de Molina, 15, 28006 Madrid, Spain

## Abstract

The accelerating deployment of artificial intelligence systems across regulated sectors has exposed critical fragmentation in risk assessment methodologies. A significant "language barrier" currently separates technical security teams, who focus on algorithmic vulnerabilities (e.g., MITRE ATLAS), from legal and compliance professionals, who address regulatory mandates (e.g., EU AI Act, NIST AI RMF). This disciplinary disconnect prevents the accurate translation of technical vulnerabilities into financial liability, leaving practitioners unable to answer fundamental economic questions regarding contingency reserves, control return-on-investment, and insurance exposure.

To bridge this gap, this research presents the **AI System Threat Vector Taxonomy**, a structured ontology designed explicitly for Quantitative Risk Assessment (QRA). The framework categorizes AI-specific risks into nine critical domains: Misuse, Poisoning, Privacy, Adversarial, Biases, Unreliable Outputs, Drift, Supply Chain, and IP Threat, integrating 53 operationally defined sub-threats. Uniquely, each domain maps technical vectors directly to business loss categories (Confidentiality, Integrity, Availability, Legal, Reputation), enabling the translation of abstract threats into measurable financial impact.

The taxonomy is empirically validated through an analysis of 133 documented AI incidents from 2025 (achieving 100% classification coverage) and reconciled against the main AI risk frameworks. Furthermore, it is explicitly aligned with ISO/IEC 42001 controls and NIST AI RMF functions to facilitate auditability. By providing the standardized inputs necessary for probabilistic modeling, specifically convolved Monte Carlo simulations, this framework allows organizations to transition from subjective, qualitative "heat maps" to rigorous financial exposure analysis. This establishes a unified language for AI risk communication, enabling evidence-based governance that satisfies both regulatory mandates and operational reality.

## 1 Introduction

Organizations deploying AI systems have responded to growing complexity by creating new governance roles, Chief AI Officers (CAIOs), AI risk managers, model auditors, and specialized red teams [28]. Unlike traditional IT security roles, these practitioners must assess threats spanning data integrity (poisoning), algorithmic fairness (bias), probabilistic outputs (hallucinations), and evolving data patterns

---

*Correspondence: hhuwyler@faculty.ie.edu; Tel: + 34 915 68 96 00.

(drift). However, a "Tower of Babel" problem persists: while engineering teams focus on technical metrics like "gradient descent manipulation" or "adversarial perturbations" [20], Board members and legal counsel require assessments in terms of financial liability, regulatory non-compliance, and reputation loss [1]. Without a shared lexicon and a structured quantification methodology, these stakeholders remain misaligned, leading to governance that is reactive rather than proactive.

## 1.1 The Regulatory Imperative

Regulatory developments have intensified the need for structured, evidence-based AI risk assessment. The European Union's AI Act, with enforcement phased in from mid-2025, classifies AI systems by risk level and mandates documented risk management systems for high-risk applications [4]. In the United States, the NIST AI Risk Management Framework (AI RMF) provides voluntary guidance but currently lacks operational threat taxonomies and specific quantification protocols [19]. Similarly, ISO/IEC 42001:2023, the first international AI management system standard, requires risk assessment under Clause 6.1 but defers the specific methodology to organizational discretion [12].

These frameworks mandate comprehensive risk assessment but provide limited operational guidance, leaving many organizations without structured methodologies to systematically identify and quantify AI-specific threats. Current industry resources provide either high-level governance principles [18] or qualitative, attack-centric threat lists like MITRE ATLAS [16]. While valuable for technical red teaming, these lists lack the "business translation" layer required for financial modeling. When asked to set contingency reserves for AI risks, estimate liability exposure, or justify the ROI of adversarial robustness testing, practitioners are left without structured approaches. Reliance on qualitative risk matrices, still prevalent in practice, has been empirically shown to produce arbitrary and inconsistent results, failing to capture the complexity of AI trustworthiness dimensions [15].

## 1.2 The Methodological Gap

The challenge is further compounded by the interdisciplinary nature of AI risks. Threats to AI trustworthiness are not isolated; human oversight failures (e.g., bias, lack of explainability) directly impact cybersecurity (e.g., loss of data integrity), and vice-versa [24]. Existing cybersecurity frameworks are ill-equipped to handle these dependencies. For instance, a "Model Inversion" attack is technically a data privacy breach, but its financial impact is realized through regulatory fines and reputational damage, a linkage often missed by purely technical threat models [13]. Furthermore, the rapid pace of AI development, particularly in Generative AI, introduces operational risks that traditional software development lifecycles do not capture, such as "hallucinations" or "prompt injection" [6].

This paper argues that robust Quantitative Risk Assessment (QRA) for AI is impossible without first establishing a valid, comprehensive taxonomy of risk scenarios. To bridge the gap between theoretical probability modeling and practical decision-making, we present the AI System Threat Vector Taxonomy. This structured ontology classifies threats into nine critical domains, mapping technical vulnerabilities directly to business loss categories (Confidentiality, Integrity, Availability, Compliance, and Reputation).

## 1.3 Research Contribution

This study addresses the operational gap between AI governance principles and quantitative risk assessment practice. This paper presents three integrated contributions:

- **A Structured Threat Taxonomy:** Nine domains encompassing 53 operationally defined sub-threats spanning the AI system lifecycle. Unlike existing technical threat lists [16], each domain includes prevalence guidance and explicit mapping to business loss categories.

- **A Quantification Bridge:** This study demonstrates how to use this taxonomy as the "input layer" for probabilistic risk modeling enabling evidence-based reserve setting and control ROI analysis.

- **A Regulatory Integration Framework:** We provide an explicit mapping of threat domains to NIST AI RMF functions and ISO/IEC 42001 controls, providing auditable compliance documentation pathways.

By standardizing the identification of risk scenarios, this framework provides the necessary inputs for the "Map" and "Measure" functions of the NIST AI RMF [19], allowing organizations to transition from qualitative heat maps to quantitative financial exposure analysis.

## 2 Background and Related Work

To situate the proposed framework, this study analyzes three distinct but currently disconnected bodies of literature: (1) AI threat taxonomies, (2) quantitative risk assessment methodologies, and (3) emerging AI governance regulations. While significant progress has been made in each silo, a comprehensive mechanism to link technical threat vectors to financial quantification and regulatory compliance remains absent.

### 2.1 Evolution of AI Threat Taxonomies

The categorization of AI threats has evolved in parallel with the expanding attack surface of Machine Learning systems. Early research focused predominantly on Adversarial Machine Learning, characterizing evasion attacks and poisoning strategies directed at classifiers [3, 8]. This research culminated in the MITRE ATLAS framework, which maps ML vulnerabilities to the tactic/technique structure of MITRE ATT&CK [16]. While ATLAS provides a granular vocabulary for red teams simulating cyber-attacks, it focuses heavily on malicious intent, largely excluding non-adversarial failures such as data drift, fairness issues, or unintended model behaviors, categories that often carry higher operational risks [15].

Conversely, the rapid adoption of Generative AI prompted the release of the OWASP Top 10 for Large Language Model Applications [22]. This taxonomy successfully highlights application-layer vulnerabilities like "Prompt Injection" and "Insecure Output Handling." However, its scope is limited to LLMs, neglecting broader ML architectures. Furthermore, ENISA's AI Threat Landscape attempts to broaden the scope to include supply chain and data quality issues [5], but remains descriptive rather than operational.

These frameworks function as static catalogs. They are siloed by domain (AppSec vs. Adversarial ML) and lack the mapping to business loss categories required for financial risk modeling. An engineer can identify a "Membership Inference Attack" using ATLAS, but existing frameworks do not guide the organization in translating that technical event into a probabilistic financial loss [13].

### 2.2 Risk Quantification in Cybersecurity and AI

Quantitative risk assessment has increasingly replaced qualitative matrices in cybersecurity to reduce subjectivity. The Factor Analysis of Information Risk (FAIR) framework models risk as a function of Loss Event Frequency and Probable Loss Magnitude, while ISO/IEC 31000 requires the use of the best available information. Recent work has demonstrated the efficacy of Monte Carlo simulations in modeling these compound uncertainties for operational risks [10].

However, applying these methods to AI introduces novel challenges. Traditional vulnerability databases (e.g., CVEs) do not track the frequency of algorithmic failures like "Hallucinations" or "Bias," making frequency estimation difficult without a specialized taxonomy. Furthermore, the impact of AI risks is often non-linear and multi-dimensional—spanning regulatory fines, reputational damage, and algorithmic remediation costs [1]. Existing quantification models have not yet been adapted to ingest the unique probability distributions associated with AI threat vectors (e.g., the continuous degradation of Model Drift versus the discrete event of a Prompt Injection).

## 2.3 Regulatory Frameworks and the "Methodology Void"

The urgency for a unified risk taxonomy is driven by a rapidly solidifying regulatory landscape.

- **The EU AI Act:** Enacted to categorize AI systems by risk severity, this regulation mandates that high-risk system providers establish a "risk management system" capable of identifying known and foreseeable risks [4]. The regulatory assessment of risk severity is based on societal harms, distinct from the internal exposure risks organizations face when onboarding AI assets.

- **NIST AI Risk Management Framework (AI RMF 1.0):** This framework defines four core functions: Govern, Map, Measure, and Manage. While the "Map" function explicitly calls for the contextualization of risks, NIST remains methodology-agnostic [19].

- **ISO/IEC 42001:2023:** The international standard for AI Management Systems requires organizations to "assess AI system risks" under Clause 6.1. However, ISO 42001 currently lacks a mature, unified control library mapped to specific AI threat vectors [12].

These frameworks are command-based; they mandate that risk be assessed but do not prescribe how. They create a compliance requirement for quantification without providing the operational taxonomy necessary to execute it.

# 3 Methodology

We employed a four-phase mixed-methods approach to develop, integrate, and validate the AI threat taxonomy. This research design ensures the framework is (1) comprehensive in coverage, (2) operationally quantifiable, (3) compliance-aligned, and (4) empirically grounded.

- Phase 1: Taxonomy Development (Systematic Literature Review, Domain Synthesis, Sub-Threat Identification)

- Phase 2: Quantification Integration (Loss Category Mapping, Distribution Selection, Convolved Adaptation)

- Phase 3: Regulatory Alignment (Mapping to NIST AI RMF, ISO 42001, EU AI Act)

- Phase 4: Empirical Validation (Incident Database Analysis n=133, AI Risk Frameworks n=4)

## 3.1 Phase 1: Taxonomy Development

### 3.1.1 Systematic Literature Review

This study conducted a systematic review following PRISMA guidelines [23] to identify existing AI threat taxonomies. Search queries were executed across ACM Digital Library, IEEE Xplore, ArXiv, and Google Scholar for the period January 2018 to November 2025. Inclusion criteria required: (1) peer-reviewed publications or authoritative technical reports, (2) explicit threat categorization for AI/ML systems, and (3) publication in English. This yielded 47 primary sources. Additionally, we reviewed regulatory and standards documents. Content analysis identified recurring threat categories and, crucially, gaps in coverage regarding operational risk factors.

### 3.1.2 Domain Structure Development

Analysis of the sources revealed overlapping but inconsistent categorization schemes. We synthesized these into nine threat domains based on three organizing principles:

- **Lifecycle Coverage:** Threats were mapped to the AI development pipeline (Data Collection, Model Training, Deployment, Operations).

- **Stakeholder Perspective:** Domains address specific organizational functions: Security, Privacy, Compliance, and MLOps.
- **Loss Category Alignment:** Each domain maps to distinct loss types per the CIA-L-R framework (Confidentiality, Integrity, Availability, Legal, Reputation).

### 3.1.3 Sub-Threat Identification

For each domain, we extracted specific attack vectors and failure modes. Sub-threats were included if they met two criteria: (1) Distinct Operational Manifestation and (2) Documented Evidence. This paper prioritized sub-threats within each domain by prevalence, drawing on OWASP frequency assessments [22] and the AI Incident Database [14]. This yielded 52 sub-threats across the nine domains.

## 3.2 Phase 2: Quantification Integration

### 3.2.1 Loss Category Mapping

To enable financial risk modeling, each threat domain was mapped to applicable loss categories using the CIA-L-R framework. This mapping translates technical threats into business impacts. For example, the Privacy domain maps primarily to Confidentiality loss and Legal loss (GDPR fines).

### 3.2.2 Probability Distribution Framework

We adapted the Convolved Monte Carlo framework [10] to AI threat characteristics by providing threat-specific distribution recommendations based on temporal patterns (discrete vs. continuous) and impact profiles (bounded vs. heavy-tailed).

### 3.2.3 Impact Modeling Approach

Impact quantification follows the mapped loss categories using forensics and notification cost models [11], regulatory penalty schedules, and customer churn models.

## 3.3 Phase 3: Regulatory Alignment

The taxonomy operationalizes the "Map" and "Measure" functions of the NIST AI RMF. Misuse threats map to GOVERN 1.5 and MANAGE 2.3; Privacy threats to MAP 1.2 and MEASURE 2.7; and Biases to MEASURE 2.11. Furthermore, it links to ISO/IEC 42001 Control 6.3.1 (Poisoning), Control 6.4.1 (Drift), and Control 6.2.2 (Biases). The taxonomy supports EU AI Act compliance by aligning threat domains with risk levels for documentation (Art. 9, Art. 10, Art. 72).

## 3.4 Phase 4: Empirical Validation

We validated taxonomy coverage through analysis of the AI Incident Database (AIID). Applying inclusion criteria (production systems, 2019-2025), we classified 133 incidents. Results indicated that Unreliable Outputs and Biases were the most prevalent failures, whereas academic literature often over-focuses on adversarial attacks.

# 4 AI System Taxonomy

This section presents the core contribution of this research: a structured ontology categorizing AI-specific threat vectors into nine critical domains.

## 4.1 Domain 1: Misuse

The Misuse domain encompasses scenarios where AI systems are utilized for unintended, unethical, or malicious purposes. Empirical evidence suggests Misuse is the most prevalent threat vector for Generative AI. The business impact is primarily reputational and legal.

- **Sub-Threats:** Prompt Injection, LLM Jailbreaking, Deepfake Generation, Disinformation Operations, Bot Abuse, Shadow AI Usage, Backdoor Attack (User-Side).

## 4.2 Domain 2: Poisoning

Poisoning refers to the injection of malicious data or components into training sets or models to corrupt behavior or logic. It represents a "high-impact, low-frequency" risk profile.

- **Sub-Threats:** Targeted Data Poisoning, Model Backdooring, Tainted Open-Source Models, Logic Corruption, Poisoned ML Libraries, Label Flipping Attacks, Gradient Manipulation, Poisoned Data Augmentation.

## 4.3 Domain 3: Privacy

The privacy domain covers the extraction or inference of sensitive information from trained models or user inputs. With the enforcement of the EU AI Act, privacy threats have shifted from theoretical concerns to immediate compliance liabilities.

- **Sub-Threats:** Model Inversion, Membership Inference, Personal Data Leakage, Sensitive Data Leakage, Inference Eavesdropping.

## 4.4 Domain 4: Adversarial

Adversarial threats involve designing harmful inputs (perturbations) to mislead or confuse AI models at run-time. While heavily cited in academia, its business relevance is sector-specific (e.g., high for autonomous vehicles, lower for standard business analytics).

- **Sub-Threats:** Evasion Attacks, Adversarial Patch/Image, Model Denial of Service (DoS), Query Flooding, Adversarial Reprogramming, Oracle/Extraction Attacks, Universal Perturbations, Adaptive Attacks.

## 4.5 Domain 5: Biases

This domain addresses models producing discriminatory, unfair, or biased outputs due to flawed data or design. Bias is a ubiquitous risk in systems making decisions about people.

- **Sub-Threats:** Representational Harm, Allocational Harm, Data Imbalance Bias, Proxy Discrimination, Algorithmic Amplification.

## 4.6 Domain 6: Unreliable Outputs

Scenarios where AI outputs are illogical, hallucinated, or non-factual without external manipulation. This is the single largest barrier to GenAI adoption in enterprise.

- **Sub-Threats:** Factual Hallucination, Source Fabrication, Logical Inconsistency, Incorrect Summarization, Unsafe Content Generation.

## 4.7 Domain 7: Drift

The deterioration of model accuracy or behavior as real-world data evolves over time. Drift is the "silent killer" of AI ROI, representing a chronic condition rather than an acute attack.

- **Sub-Threats:** Concept Drift, Data Distribution Drift, Upstream Data Changes, User Behavior Change, Feedback Loop Drift.

## 4.8 Domain 8: Supply Chain

Attacks propagated via third-party components, pre-trained models, data sources, or MLOps infrastructure.

- **Sub-Threats:** Compromised Pre-trained Model, Vulnerable ML Framework, Insecure Data Feeds/APIs, Container Image Poisoning, Compromised Annotation Tools.

## 4.9 Domain 9: IP Threat

The extraction of sensitive intellectual property, proprietary algorithms, or model weights from deployed systems.

- **Sub-Threats:** Model Extraction/Theft, Data Exfiltration, Proprietary Logic Theft, Hyperparameter Stealing, Watermark Removal.

# 5 From Taxonomy to Risk Quantification

The transition from qualitative threat identification to quantitative risk modeling is a fundamental requirement for the economic governance of AI systems.

## 5.1 Quantification Workflow

Application of the framework follows a rigorous six-step process designed to integrate into existing enterprise risk management workflows:

1. **Vulnerability Assessment:** Analyzing the AI system's architecture to determine the "exposure factor."

2. **Threat Identification:** Utilizing the taxonomy to identify relevant threat domains.

3. **Scenario Definition:** Translating abstract threats into concrete risk scenarios (e.g., "Misuse" -> "External actor uses prompt injection via API").

4. **Parameter Calibration:** Calibrating expected frequency and impact distributions using expert-derived parameters and historical data. Controls are modeled as reductions in frequency or loss magnitude [2].

5. **Reserve Setting and Reporting:** Establishing contingency reserves (e.g., VaR 95%) and documenting compliance.

## 5.2 The Economic Imperative for Quantification

Quantification allows organizations to assign a monetary value to model error rates, aligning technical metrics with business criticality [27]. It provides the ROI analysis necessary to prioritize security investments [26] and enables the structuring of warranties and insurance products [7, 21].

# 6    Results

The taxonomy underwent empirical validation through two complementary methods: analysis of documented AI incidents and comparative coverage analysis against existing frameworks.

## 6.1    Incident Database Coverage Analysis

We reviewed 133 distinct AI incidents reported in the AI Incident Database (AIID) between May 30, 2025, and November 17, 2025. The analysis confirmed that 100% of the reviewed cases could be successfully classified into the proposed taxonomy.

- **Misuse (n=81):** 61% of all cases, aligning with the proliferation of GenAI tools.

- **Unreliable Outputs (n=36):** 27% of cases, driven by hallucinations.

- **Supply Chain (n=7):** 5%.

The low counts for Biases and Drift likely reflect reporting bias, as these are often internal failures not publicly disclosed.

## 6.2    Framework Comparison Analysis

Comparative analysis against MITRE ATLAS, OWASP Top 10 for LLM, and ENISA Threat Landscape demonstrated coverage advantages. The taxonomy uniquely integrates threats across security, privacy, fairness, and reliability dimensions addressed separately in existing frameworks.

# 7    Discussion

The standardization of AI threat vectors provides a strategic governance tool bridging operational silos.

- **For AI Auditors:** The taxonomy serves as a definitive "Audit Checklist" for completeness.

- **For Red Teaming:** It provides a structured scope for penetration testing, moving beyond ad-hoc testing to scenario-based campaigns.

- **For Compliance:** It supports ISO 42001 and EU AI Act compliance by providing a defensible methodology for identifying "known and foreseeable risks."

# 8    Conclusion

As AI systems transition to critical infrastructure, governance must evolve from reactive "fire-fighting" to quantified risk management. This research addresses the foundational gap by providing the **AI System Threat Vector Taxonomy**. This structured ontology translates technical vulnerabilities into business loss categories, enabling Quantitative Risk Assessment (QRA). Empirical validation confirms that while academia prioritizes adversarial novelty, operational reality is dominated by Misuse and Unreliable Outputs. Organizations adopting this framework can move towards auditable, insurable, and resilient AI deployment.

# Data Availability

The full taxonomy, including updated scenario lists and mapping files, is available as an open-source repository at GitHub[1]. This dataset is licensed under CC-BY 4.0.

---

[1] https://github.com/hwyler/HernanHuwylerRiskManagement/blob/main/AIThreatTaxonomy

## Funding

## Conflicts of Interest

The author declares no conflicts of interest. The views expressed are those of the author and do not necessarily reflect the official policy or position of any affiliated agency.

## Author Biography

**Prof. Hernan Huwyler, MBA CPA** serves as the Academic Director for Compliance, Risk Management, and AI programs at IE Business School and IE Law School. Concurrently, he is a Senior Manager at Capgemini Invent's Applied AI Lab (Nordics), leading enterprise-wide AI governance, risk modeling, and control initiatives. A recognized expert in the intersection of technical AI security and regulatory compliance, Prof. Huwyler specializes in operationalizing frameworks such as the EU AI Act, ISO 42001, and NIST AI RMF.

## References

[1] A. Aluthwala and S. Wickramarathne, "Research of How Artificial Intelligence Impact on Project Risk Management," *Journal of Business and Management*, Oct. 2024.

[2] S. Bahadur and K. Tucker, "Identifying and addressing the risks of AI through regulations, compliance controls and technical design," *Journal of Financial Compliance*, vol. 8, no. 2, pp. 112–130, Dec. 2024.

[3] B. Biggio and F. Roli, "Wild patterns: Ten years after the rise of adversarial machine learning," *Pattern Recognition*, vol. 84, pp. 317–331, 2018.

[4] J. M. Camacho, A. C. Vieira, and D. Arroyo, "A Cybersecurity Risk Analysis Framework for Systems with Artificial Intelligence Components," *arXiv preprint arXiv:2401.01630*, 2024.

[5] ENISA, "AI Threat Landscape," European Union Agency for Cybersecurity, 2020.

[6] M. I. Faruk, F. W. Plabon, and U. S. Saha, "AI-Driven Project Risk Management: Leveraging Artificial Intelligence to Predict, Mitigate, and Manage Project Risks in Critical Infrastructure and National Security Projects," *Journal of Computer Science and Technology Studies*, vol. 7, no. 6, 2025.

[7] R. Gipiškis, A. S. Joaquin, and Z. S. Chin, "Risk Sources and Risk Management Measures in Support of Standards for General-Purpose AI Systems," *arXiv preprint arXiv:2410.23472*, 2024.

[8] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.

[9] Z. Guo and R. Li, "AI-Driven Risk Management for Sustainable Enterprise Development: A Review of Key Risks," *International Journal of Business and Management*, vol. 19, no. 6, p. 82, 2024.

[10] H. Huwyler, "Quantitative Risk Assessment in R: An Open-Source Convolutional Framework for Modeling Uncertainty and Reserves," Zenodo, 2025.

[11] IBM Security, "Cost of a Data Breach Report 2023," 2023.

[12] H. Li, M. Yazdi, and A. Nedjati, "Harnessing AI for Project Risk Management: A Paradigm Shift," in *Advances in Computational Intelligence*, Springer, 2024.

[13] M. Mahmoud, "The Risks and Vulnerabilities of Artificial Intelligence Usage in Information Security," *IEEE Conference on Computational Science and Computational Intelligence*, 2023.

[14] S. McGregor, "Preventing Repeated Real World AI Failures by Cataloging Incidents: The AI Incident Database," *AAAI Conference on Artificial Intelligence*, 2020.

[15] A. Metwally, S. A. Ali, and A. T. Mohamed, "Thinking Responsibly About Responsible AI in Risk Management: The Darkside of AI in RM," *IEEE Conference on Engineering and Technology*, 2024.

[16] MITRE ATLAS, "Adversarial Threat Landscape for Artificial-Intelligence Systems," 2024.

[17] J. C. Muria-Tarazón, J. V. Oltra-Gutiérrez, and R. Oltra-Badenes, "Uncovering Research Trends on Artificial Intelligence Risk Assessment in Businesses," *Applied Sciences*, vol. 15, no. 3, p. 1412, 2025.

[18] J. Newman, "A Taxonomy of Trustworthiness for Artificial Intelligence," *UC Berkeley Center for Long-Term Cybersecurity*, 2023.

[19] NIST, "Artificial Intelligence Risk Management Framework (AI RMF 1.0)," NIST.AI.100-1, 2023.

[20] NIST, "Adversarial Machine Learning—A Taxonomy and Terminology of Attacks and Mitigations," NIST.AI.100-2, 2024.

[21] C. Novelli, F. Casolari, and A. Rotolo, "AI Risk Assessment: A Scenario-Based, Proportional Methodology for the AI Act," *Digital Society*, 2024.

[22] OWASP, "OWASP Top 10 for Large Language Model Applications," 2023.

[23] M. J. Page *et al.*, "The PRISMA 2020 statement: an updated guideline for reporting systematic reviews," *BMJ*, vol. 372, 2021.

[24] N. Polemi, K. Kioskli, and T. P. Kearns, "A unified taxonomy of threats for AI trustworthiness," *Frontiers in Big Data*, 2024.

[25] N. I. Qureshi, A. Garg, and R. Singh, "AI and Corporate Risk Management: Identifying and Mitigating Technological and Ethical Risks," *IEEE International Conference on Knowledge Engineering and Communication Systems*, 2024.

[26] R. Schnitzer, A. Hapfelmeier, and S. Gaube, "AI Hazard Management: A Framework for the Systematic Management of Root Causes for AI Risks," *arXiv preprint arXiv:2310.16727*, 2024.

[27] R. Sharma, V. Harish, and G. Rana, "Navigating Risk In The Age Of Artificial Intelligence: Assessing And Identifying Risks With AI Strategies," *Kuey*, vol. 30, no. 4, 2024.

[28] P. Singh, "AI and Financial Risk Management Revolutionizing Risk Assessment and Mitigation," in *Artificial Intelligence for Financial Risk Management and Analysis*, IGI Global, 2025.