

# Emerging Threats in AI: A Detailed Review of Misuses and Risks Across Modern AI Technologies

Niyat Seghid<sup>1\*</sup>, Farkhund Iqbal<sup>1</sup>, Khalifa Al-Room<sup>2</sup>, Áine MacDermott<sup>3</sup>

<sup>1</sup>College of Technological Innovation, Zayed University, Abu Dhabi, UAE

<sup>2</sup>Dubai Police HQ, Dubai, UAE

<sup>3</sup>School of Computer Science and Mathematics, Liverpool John Moores University, Liverpool, UK

\* **Correspondence:** Niyat Seghid, College of Technological Innovation, Zayed University, UAE  
[niyat.seghid@zu.ac.ae](mailto:niyat.seghid@zu.ac.ae)

**Keywords:** Artificial Intelligence, AI Misuse, AI Risk, Deepfakes, Adversarial Attacks, Privacy Violations, Algorithmic Bias, AI Security

## Abstract

The swift evolution of artificial intelligence technologies (AI) has introduced unparalleled capabilities, alongside critical vulnerabilities that can be exploited maliciously or cause unintended harm. While numerous efforts have emerged to govern AI risks, there remains a lack of comprehensive analysis of how AI systems are actively being misused. This paper offers an in-depth review of AI misuses across modern technologies, analyzing attack mechanisms, documented incidents, and emerging threat vectors. We provide a brief review of AI risk repositories and existing taxonomic approaches to set the context, and then synthesize them into a comprehensive categorization of AI misuse across nine primary domains: (1) Adversarial Threats, (2) Privacy Violations, (3) Disinformation, Deception & Propaganda, (4) Bias & Discrimination, (5) System Safety & Reliability Failures, (6) Socioeconomic Exploitation & Inequality, (7) Environmental & Ecological Misuse, (8) Autonomy & Weaponization, and (9) Human Interaction & Psychological Harm. Across these domains, we identify and analyze distinct categories of AI misuses and risks, providing technical depth on exploitation mechanisms, documented cases with quantified impacts, and the latest developments including large language model vulnerabilities and multimodal attack vectors. We also assess the effectiveness of current mitigation strategies and countermeasures, evaluating technical security frameworks (e.g. MITRE ATLAS, OWASP Top 10 for Large Language Models (LLMs), MAESTRO), regulatory approaches (e.g. EU AI Act, NIST AI RMF), and compliance standards. Our analysis reveals significant gaps between AI capabilities and robustness of defensive measures, with adversaries holding persistent advantages across most attack categories. This work contributes to the field by: (1) systematically consolidating fragmented AI risk and misuse taxonomies and repositories, (2) developing a unified taxonomy of AI misuse patterns grounded in both theoretical models and empirical incident data, (3) critically evaluating the effectiveness of existing mitigation strategies, and (4) identifying priority research gaps to foster the development of more robust, ethical, and secure AI systems.

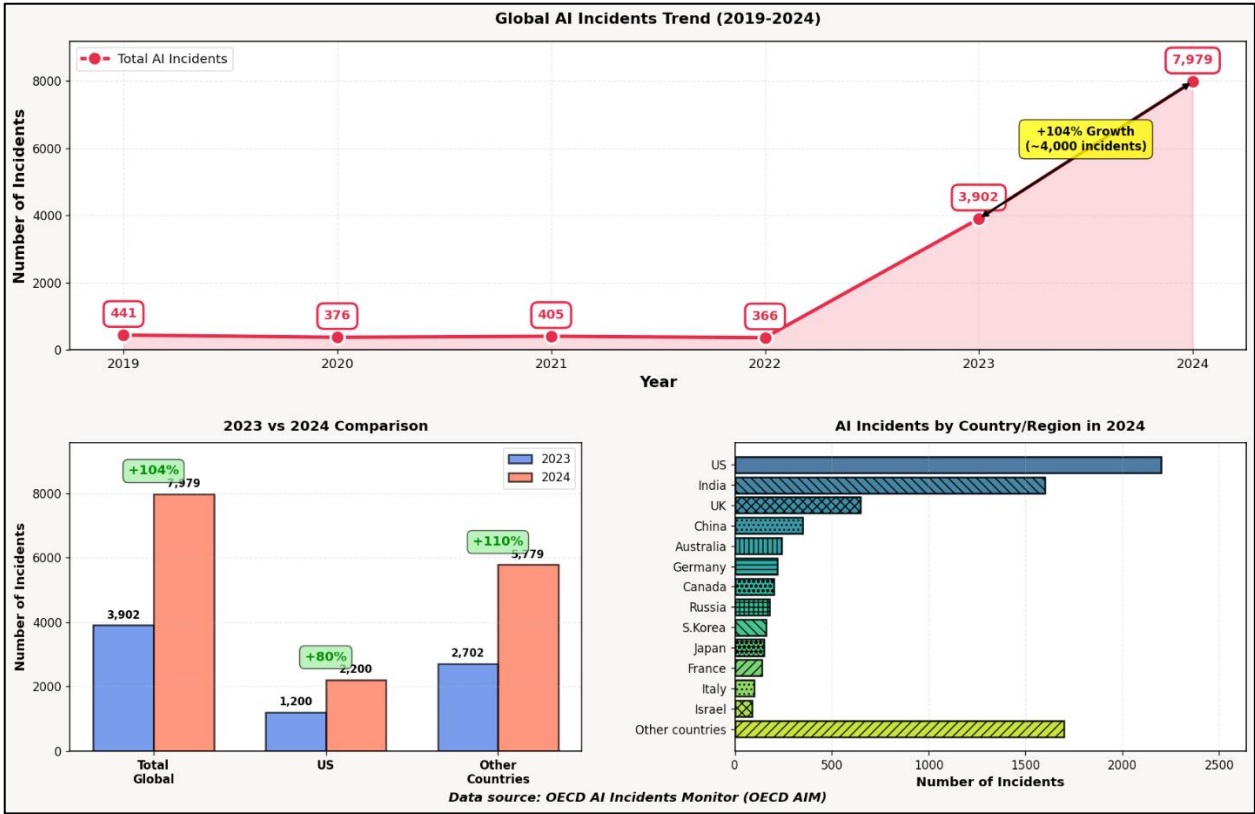
## 1 Introduction

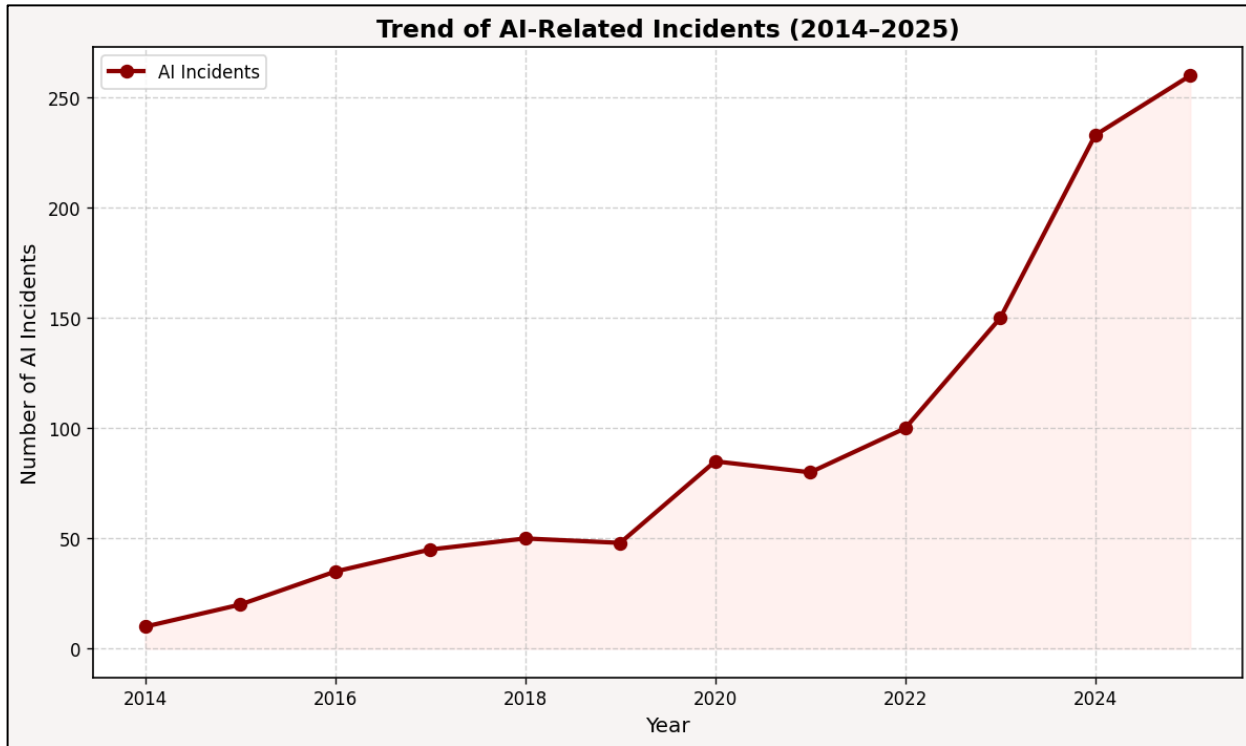
Artificial Intelligence (AI) has rapidly evolved into one of the most transformative technologies of the 21st century, reshaping industries, governance, and everyday life. Deep learning breakthroughs since 2012 (LeCun et al., 2015), proliferation of LLMs (Brown et al., 2020), advances in generative AI (OECD, 2019), and deployment of autonomous systems (Scharre, 2018) have created unprecedented

capabilities. However, these same technologies have also introduced critical vulnerabilities, offering new vectors for malicious exploitation.

The dual-use nature of AI lies at the core of this concern. Algorithms and models designed for beneficial purposes can be adapted for malicious or unethical use. Natural language models that enable intelligent assistants may be leveraged to produce convincing disinformation (Goldstein et al., 2023); generative models supporting creative industries can be used to fabricate realistic deepfakes (Chesney et al., 2019); computer vision systems designed for safety or accessibility may facilitate intrusive surveillance or unauthorized biometric profiling (Buolamwini et al., 2018); and recommendation algorithms designed to personalize user experiences can be exploited to manipulate behavior (Matz et al., 2017). As AI systems grow in capability, autonomy, and accessibility, their misuse potential increases in both scale and sophistication.

Recent incidents demonstrate that AI misuse is no longer theoretical but a pressing global issue with measurable consequences: deepfake-enabled fraud exceeding \$25 million (Stupp, 2019), AI-generated election disinformation affecting millions (DiResta et al., 2024), wrongful arrests from facial recognition errors (Garvie et al., 2016), algorithmic discrimination in healthcare affecting 200 million people annually (Obermeyer et al., 2019), and adversarial attacks on safety-critical systems (Eykholt et al., 2018). The proliferation of synthetic media, automated cyberattacks, and algorithmic discrimination reflects how AI can amplify deception, erode privacy, and reinforce social inequalities. Moreover, AI-driven automation and personalization have accelerated the scale and precision of harmful activities, from widespread disinformation campaigns to targeted phishing and identity manipulation. These developments highlight a growing mismatch between AI advancement and the capacity to detect, regulate, or mitigate its misuse, raising pressing ethical and security concerns. Recent statistics supporting these observations are shown in Fig.1.





(b) AI related incidents reported worldwide (2014 – 2024) from Artificial Intelligence Index Report 2025

Figure 1. Recent AI Misue Statistics

Despite the expanding body of literature on AI misuse, the research landscape remains highly fragmented. Existing studies often focus on specific domains, such as deepfakes, adversarial attacks, or data privacy, without integrating insights across technical, ethical, and societal dimensions (Slattery et al., 2024), (National Institute of Standards and Technology, 2023). Moreover, inconsistent terminology, varied categorization schemes, and rapidly evolving threat models further complicate efforts to develop a unified understanding of the full spectrum of misuse. This fragmentation creates challenges for those seeking to assess risks comprehensively or develop interdisciplinary strategies for prevention and response.

This review addresses that fragmentation by providing a systematic synthesis of AI misuse research across technical, ethical, and societal perspectives. Rather than proposing entirely new theoretical models, this paper organizes and consolidates existing knowledge to create a comprehensive and accessible overview of how AI technologies can be misused. To establish context, we briefly examine empirical AI risk and misuse taxonomies and repositories. Building upon insights from these sources, we propose a consolidated nine-domain categorization of AI misuse, each suitable for detailed technical and socio-ethical analysis. Across these domains, we provide technical depth on exploitation mechanisms, detailed real-world incidents, and discussion of countermeasure effectiveness.

Through extensive analysis of academic publications, industry reports, and documented misuse cases, this paper seeks to:

- Critically analyze existing AI risk and misuse taxonomies and repositories, examining how different research initiatives categorize AI threats and identifying convergences, gaps, and complementarities across classification schemes
- Synthesize a unified taxonomy of AI misuse that categorizes AI misuse and vectors across technical, social, and ethical dimensions,
- Analyze various forms of AI misuse, identifying common vulnerabilities and attack patterns, providing technical depth on exploitation mechanisms, and examining real-world case studies to understand the practical manifestations, impacts, and consequences.
- Evaluate existing mitigation strategies, assessing technical security frameworks, regulatory approaches, and compliance standards for their effectiveness, limitations, and applicability across different contexts.

The remainder of this paper is organized as follows: Section 2 reviews existing AI misuse and risk frameworks. Section 3 describes the review methodology. Section 4 categorizes AI misuse domains and presents synthesized findings. Section 5 discusses key incidents and implications across domains. Section 6 reviews current mitigation strategies and challenges and AI risk governance frameworks. Section 7 concludes with recommendations for future research and governance directions.

## 2 Background: Review of Existing AI Risk and Misuse Taxonomies

Several organizations and research groups have developed frameworks and/or taxonomies to classify and understand the risks and misuse of artificial intelligence, reflecting the growing need for systematic approaches to AI safety and governance. We briefly review major existing taxonomies to establish context before presenting our own categorization.

Among the most influential is the *MIT AI Risk Repository* developed by Slattery et al. (2024), which represents the most comprehensive effort to date, extracting and categorizing 1,612 risks from 65 existing taxonomies (Slattery et al., 2024). The framework organizes risks using a dual approach: a causal taxonomy classifying by entity, intentionality, and timing; and a domain taxonomy with seven domains and 24 subdomains covering discrimination and toxicity, privacy and security, misinformation, malicious actors and misuse, human-computer interaction, socioeconomic and environmental impacts, and AI system safety. While highly valuable for conceptual coverage, the repository largely abstracts away from detailed technical attack mechanisms and operational misuse pathways.

Complementing this work, incident-centered repositories provide empirical grounding. The *AI Incident Database* systematically catalogs real-world AI failures and misuse events, emphasizing recurrence patterns and socio-technical root causes (McGregor, 2021). Similarly, the *OECD AI Incident Monitor* aggregates reported AI-related incidents across jurisdictions, offering longitudinal insights into emerging misuse trends (OECD, 2023). These repositories shift the focus from hypothetical risks to documented harms, but do not provide fine-grained technical taxonomies. The OECD AI Incidents Monitor provides an international repository of documented AI incidents, collecting reports from multiple sources including news media, research papers, and direct submissions (OECD.AI, 2024). The database categorizes incidents by type (bias/discrimination, privacy violation, safety failure, etc.), sector (healthcare, finance, transportation, etc.), and AI technology involved (computer vision, NLP, recommendation systems, etc.). However, the limitation lies in its limited technical depth in incident descriptions.

Several security-oriented and adversarial taxonomies focus explicitly on malicious AI use. *MITRE Adversarial Threat Landscape for Artificial-Intelligence Systems (ATLAS)* provides a tactics–techniques–procedures (TTP) knowledge base documenting real-world attacks on AI systems, including data poisoning, model evasion, and supply-chain compromise (MITRE ATLAS, 2024).

The *ENISA AI Threat Landscape* similarly categorizes AI-related cybersecurity threats, emphasizing vulnerabilities, attacker capabilities, and systemic impacts (ENISA, 2020). The *OWASP Top 10 for LLMs* further refines this focus for generative and language models, identifying prompt injection, insecure output handling, training data poisoning, model denial-of-service, and supply-chain vulnerabilities as dominant misuse vectors (OWASP Foundation, 2024).

Beyond security, multiple domain-specific risk taxonomies have emerged. In healthcare, Golpayegani et al. (2022) propose a structured taxonomy covering clinical, ethical, and operational AI risks, highlighting patient harm and diagnostic bias (Golpayegani et al., 2022). In international security, UNIDIR synthesizes AI risks related to strategic stability, escalation dynamics, and confidence-building measures (UNIDIR, 2023). Mahmoud (2023) examines AI risks in information security, emphasizing automation-enabled attack amplification. The *IAA AITF AI Risks Taxonomy* (2024) introduces a three-level taxonomy tailored to actuarial and financial risk management, mapping AI-specific risk amplification onto traditional actuarial risk categories.

Additional academic contributions have expanded taxonomies to socio-technical and human-centered harms. Critch and Russell (2023) introduced their Taxonomy and Analysis of Societal-Scale Risks from AI (TASRA), examining macro-level dimensions including risk accountability and ethical alignment (Critch et al., 2023). TASRA focuses on long-term, systemic risks rather than near-term incidents, considering how AI could reshape power dynamics, decision-making authority, and social institutions. Weidinger et al. (2022) proposed taxonomies specifically targeting large language model risks, highlighting concerns such as discrimination, information hazards, and malicious uses (Weidinger et al., 2022). Marchal et al. (2024) focused on generative AI misuse, identifying threats including prompt injection, model leakage, and large-scale disinformation (Marchal et al., 2024). Moreover, Zhang et al. (2025) addressed the emerging domain of AI companionship applications, developing a taxonomy of harmful algorithmic behaviors that can occur in human-AI relationship, examining the psychological and relational harms that can emerge when AI systems are designed to form ongoing personal bonds with users, including emotional manipulation, unhealthy dependency, and intimate privacy violations (Zhang et al., 2025).

Collectively, these taxonomies and repositories provide complementary but fragmented views of AI misuse, varying in scope, granularity, and empirical grounding. Some emphasize technical attack vectors, others societal harms or domain-specific risks, and few attempt holistic integration. Building upon these efforts, this review consolidates and aligns them into a unified nine-domain taxonomy of AI misuse, collectively capturing the technical, ethical, and socio-technical dimensions of contemporary AI misuse, grounded in documented incidents and technical exploitation mechanisms, representing both underrepresented dimensions such as environmental sustainability as well as established concerns.

### 3 Methodology

This study employs a mixed-methods approach, combining a systematic literature review with in-depth case study analysis to develop a comprehensive understanding of AI misuse patterns and mitigation

strategies. By integrating these elements, the research bridges technical, social, and ethical perspectives, while grounding theoretical insights in real-world incidents.

### 3.1 Reporting Standards

This study presents a systematic review in accordance with the PRISMA 2020 guidelines as illustrated in Fig. 2 (Haddaway et al., 2022). Academic literature was retrieved from major scholarly databases including IEEE Xplore, Scopus, and the ACM Digital Library. In parallel, relevant case-based and policy-oriented materials were sourced from grey literature repositories and organizational databases such as the NIST repository, MIT AI Repository, and AI Incident Database. To capture developments coinciding with the rise of modern AI applications, the search covered the period from 2012 to 2025, aligning with the deep learning era and the acceleration of AI adoption across critical domains. The search queries combined key terms such as “artificial intelligence misuse”, “AI risks”, “AI security threats”, “adversarial attacks”, “AI safety”, “algorithmic bias”, “deepfakes”, and related variations.

A total of 128 records (104 from academic databases and 24 from other sources) were initially identified. After removing duplicate records, 125 unique studies were screened based on titles and abstracts. During this phase, 33 papers were excluded due to not meeting the inclusion criteria. Full-text retrieval was sought for 68 studies, of which seven were excluded after detailed assessment. The remaining 61 database-based studies were included in the final synthesis. From the additional 24 external sources, six reports were excluded, and 18 reports were retained. Altogether, 79 studies (61 database studies + 18 other sources) were included in the final review.

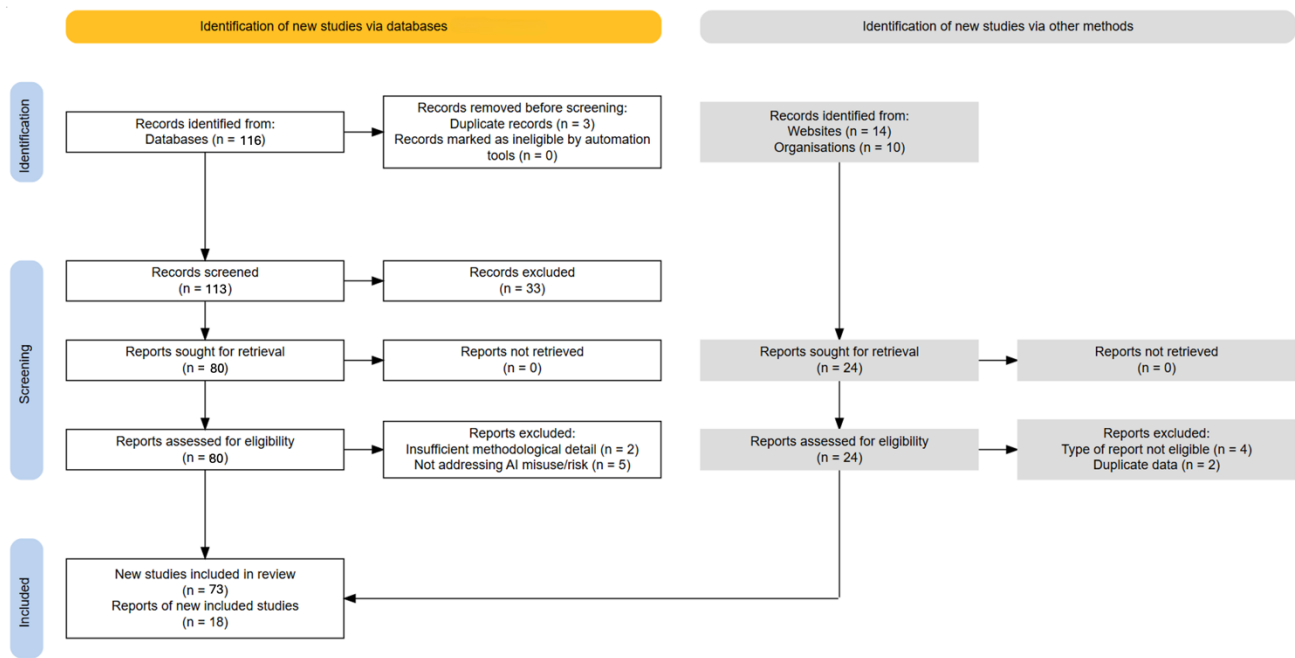


Figure 2. PRISMA flow diagram of study selection.

The diagram (Fig.2) outlines the identification, screening, eligibility assessment, and inclusion process for the reviewed studies.

### 3.2 Case Study Selection

Case studies were identified through systematic monitoring of multiple sources including the AI Incident Database, vulnerability disclosures, regulatory publications, industry transparency reports, and media coverage. The aim was to capture a diverse set of examples spanning different domains of misuse, levels of severity, and cultural or geographic contexts. Selected cases were required to have documented evidence verifying the incident. This ensured that the analysis addressed both technical and social dimensions of misuse and highlighted the ways AI vulnerabilities manifest in practice.

### 3.3 Taxonomy Development

The taxonomy of AI misuse was derived through a systematic synthesis of major AI risk frameworks, consolidating and extending them into a unified and comprehensive classification. Drawing on established taxonomies, we identified key conceptual overlaps and critical gaps, particularly in areas such as environmental sustainability and human-AI psychological manipulation. Accordingly, AI misuse was classified into nine domains encompassing both the technical and socio-technical dimensions of contemporary misuse. The detailed research workflow is given in Fig. 3, illustrating the methodological workflow adopted in the study, beginning with a systematic literature review and case study analysis, followed by framework synthesis and refinement, taxonomy construction, and the formulation of mitigation strategies.

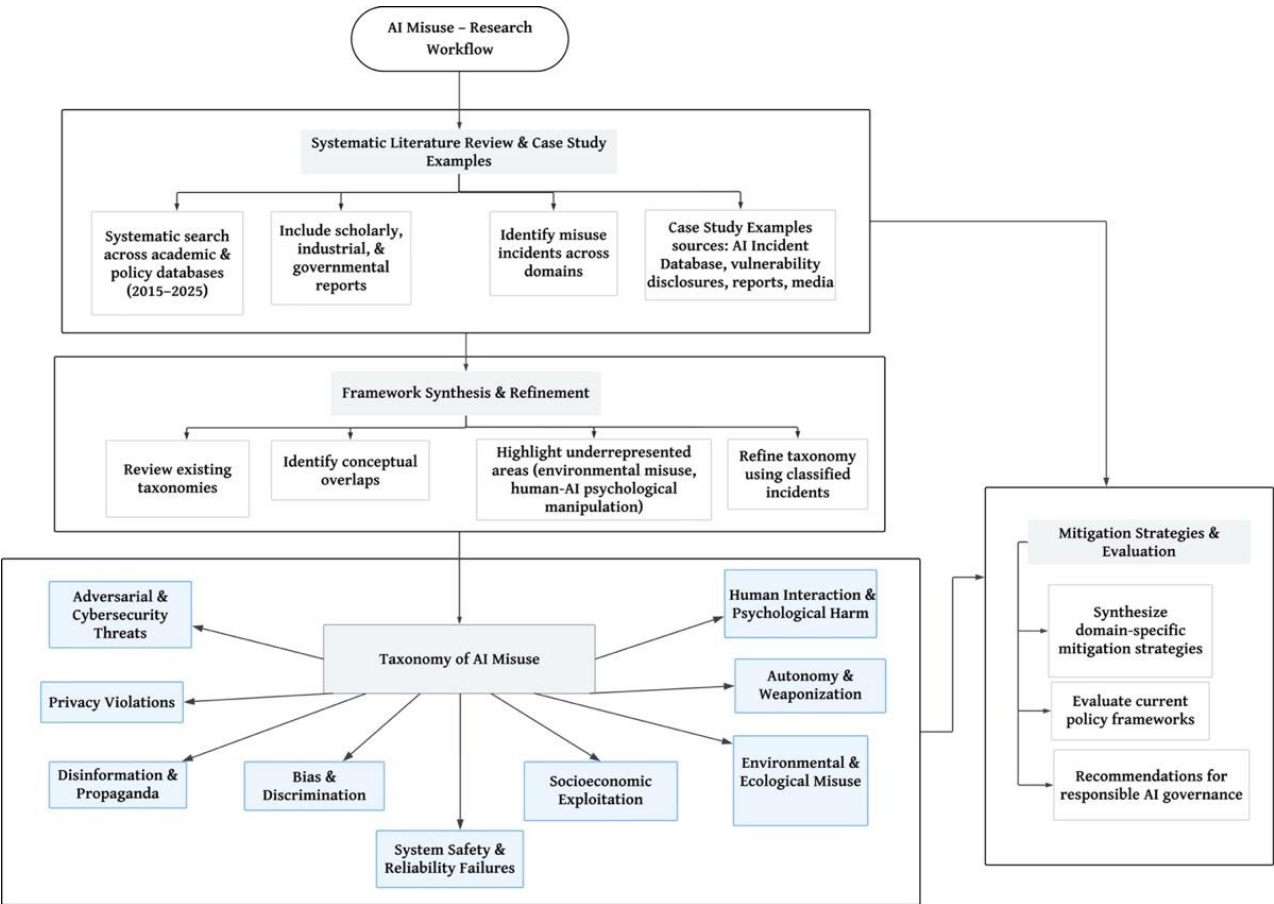


Figure 3. Methodological workflow of the research

## 4 Taxonomy of AI Misuse



To enable systematic analysis of AI misuse, this study develops a taxonomy that organizes threat vectors across nine primary domains, each further divided into subcategories. These domains encompass adversarial and cybersecurity threats, privacy violations, disinformation and synthetic media, bias and discrimination, system safety and reliability failures, socio-economic exploitation and inequality, environmental and ecological misuse, autonomy and weaponization, and human interaction and psychological harm. While these categories represent distinct manifestations of misuse, they are also deeply interconnected, with vulnerabilities in one domain frequently compounding risks in another. By organizing the landscape in this structured manner, the taxonomy provides a comprehensive framework for both researchers and practitioners to classify incidents, anticipate threats, and design targeted interventions (Slattery et al., 2024), (NIST, 2023).

**Table 1. Taxonomy of AI Misuse**

	Domain	Key Examples (Attack Types)	Mechanisms	Implications
1	<b><i>Adversarial Threats</i></b>	Evasion attacks, poisoning, backdoors, model extraction, membership inference, model inversion, supply chain attacks	Subtle perturbations to inputs, maliciously crafted training data, unauthorized model queries, compromised dependencies in AI pipelines	Compromise of AI integrity, intellectual property theft, inaccurate outputs in critical systems (e.g., autonomous vehicles, medical AI)
2	<b><i>Privacy Violations</i></b>	Sensitive attribute inference, re-identification, data leakage, unauthorized surveillance	Analysis of model outputs, correlational inference, generative model reconstruction	Breach of user confidentiality, regulatory violations, erosion of trust in digital services
3	<b><i>Disinformation, Deception, &amp; Propaganda</i></b>	Deepfakes, automated fake news, targeted propaganda, harmful/illegal content generation, prompt injection, erosion of trust	Generative models for text, image, video; automated amplification on social media	Misinformation at scale, manipulation of public opinion, destabilization of political and social systems
4	<b><i>Bias &amp; Discrimination</i></b>	Gender, racial, socioeconomic biases; opaque decision-making; stereotyping	Biased training data, reinforcement of historical inequities, algorithmic opacity	Unequal access to services, perpetuation of social inequities, reputational and legal risks for deploying organizations
5	<b><i>System Safety &amp; Reliability Failures</i></b>	Autonomous vehicle accidents, misdiagnoses in healthcare, industrial automation failures	Model misbehavior under unexpected conditions, inadequate validation and monitoring	Physical harm, operational disruption, loss of human life or safety incidents
6	<b><i>Socioeconomic Exploitation &amp; Inequality</i></b>	Job displacement, economic fraud, cheating, microtargeting, exploitation of vulnerable populations	Automation replacing human labor, AI-driven manipulation of financial and social systems	Increased economic disparities, reduced employment opportunities, ethical and legal challenges in AI governance
7	<b><i>Environmental &amp; Ecological Misuse</i></b>	High energy consumption of AI, carbon-intensive model training, automated harmful industrial practices	Resource-intensive model training, misuse of AI in environmental systems	Increased carbon footprint, ecological damage, sustainability concerns
8	<b><i>Autonomy &amp; Weaponization</i></b>	Autonomous drones, lethal AI weapons, cyber-physical attacks, Agentic AI systems	Decision-making without human oversight, AI-guided military systems	Escalation in conflict, ethical concerns over lethal AI, potential breaches of international law
9	<b><i>Human Interaction &amp; Psychological Harm</i></b>	Emotional manipulation via AI, addiction to AI interfaces, mental health impacts	Personalized content targeting, persuasive AI, immersive digital environments	Anxiety, depression, behavioral manipulation, loss of agency and autonomy

## 4.1 Adversarial Threats

Machine learning systems are vulnerable to a wide array of adversarial attacks that exploit both the data and model layers of the learning pipeline (as shown in Fig. 4). *Evasion attacks* exploit weaknesses in trained models by subtly perturbing inputs, causing misclassifications without visibly altering the underlying data (Biggio et al., 2013). These perturbations are often unnoticeable to humans but are designed to shift the input across the model’s decision boundary. Such attacks are particularly



concerning in real-time systems, where even minor perturbations to data can induce high-confidence yet incorrect predictions, compromising safety-critical applications such as healthcare decisions, autonomous navigation, or biometric authentication.

*Poisoning attacks*, in contrast, corrupt the training dataset itself, embedding malicious patterns that compromise the integrity of models even before deployment (Gu et al., 2017). Attackers may inject a small fraction of poisoned samples into the training sets to manipulate model behavior, either globally (causing widespread accuracy degradation) or specifically (triggering backdoor conditions under certain inputs). For instance, a backdoor poisoning attack might train a face recognition model to always classify images containing a specific pixel pattern as a trusted user, regardless of the actual identity. Such manipulation remains dormant during evaluation, evading detection, and activates only under attacker-controlled triggers. Because machine learning pipelines often rely on large, automatically scraped or user-contributed data, the injection of poisoned samples is both feasible and difficult to detect. Gu et al. (2017) introduced "BadNets," demonstrating how backdoor triggers embedded in training data enable attackers to maintain control over model behavior post-deployment. The poisoned model performs normally on clean inputs but exhibits attacker-specified behavior when triggered.

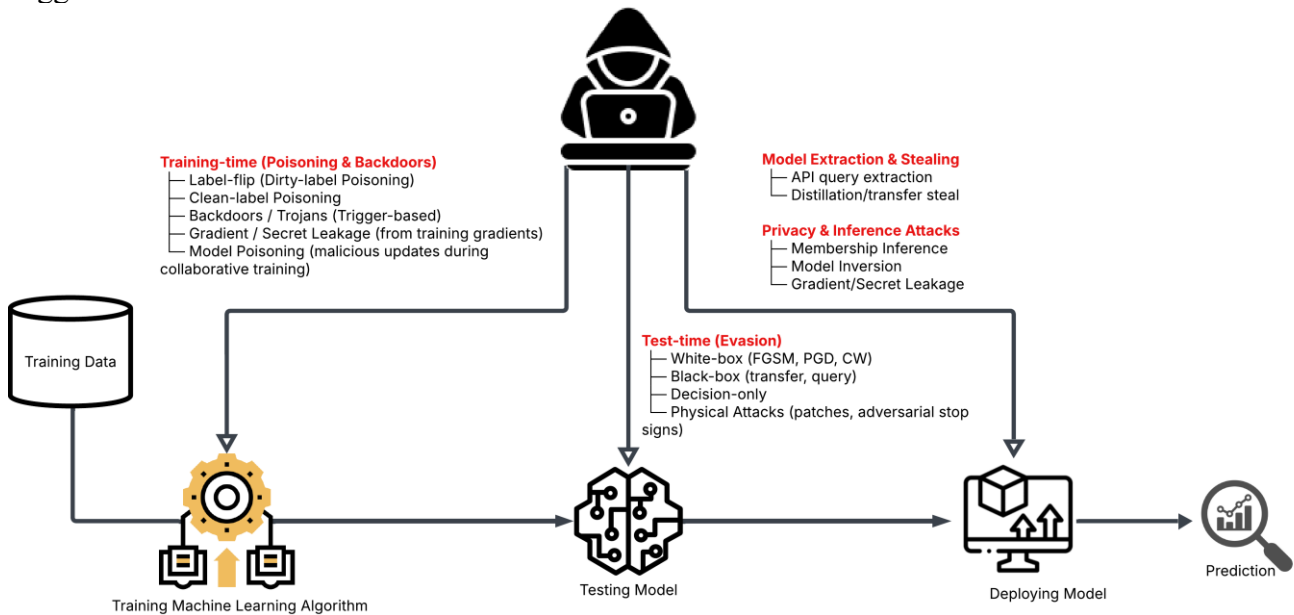


Figure 4. Examples of adversarial attacks across the machine learning lifecycle

*Model extraction attacks* further extend the adversarial threat surface by demonstrating how adversaries can reconstruct proprietary models through systematic querying and analyzing outputs (Papernot et al., 2017). By sending numerous inputs to a deployed model (often accessible via APIs) and recording the corresponding outputs, attackers can approximate the model's decision boundaries and replicate its functionality locally (see Fig. 5). This stolen surrogate model can then be exploited for additional purposes, such as launching more precise evasion attacks or performing model inversion to recover sensitive training data. In many cases, such extraction requires no privileged access, relying solely on adaptive query strategies and output probability vectors exposed by the API. Tramèr et al. demonstrated that prediction APIs expose sufficient information for attackers to build functionally equivalent models (Tramèr et al., 2016).



Figure 5. Model extraction attack

These attacks are not merely academic. In computer vision, adversarial perturbations have been shown to cause traffic sign recognition systems to misidentify stop signs as yield signs, with potentially catastrophic consequences for autonomous driving (Eykholt et al., 2018). In natural language processing, adversarial inputs can be crafted by substituting semantically similar words or introducing orthographic noise, allowing attackers to manipulate sentiment analysis models or bypass content moderation. The persistence of these vulnerabilities highlights the fragility of AI systems operating in adversarial environments (Madry et al., 2017).

## 4.2 Privacy Violations

AI systems frequently depend on vast amounts of personal and sensitive data, creating risks of privacy violations at multiple levels. At the individual level, models are vulnerable to attacks that expose whether particular data points were part of the training set, referred as membership inference attack (Hu et al., 2022). Model inversion attacks similarly enable the reconstruction of sensitive features from model outputs (Fredrikson et al., 2015). These vulnerabilities illustrate that AI systems, even when anonymized, can inadvertently leak private information.

At the systemic level, AI-driven surveillance technologies such as facial recognition amplify longstanding privacy concerns. It has been demonstrated by Sharif et al. (2016), that specially crafted eyeglass frames could enable individuals to impersonate others or evade facial recognition systems and access controls, making them particularly concerning for safety-critical applications (Sharif et al., 2016). Other studies have also revealed consistent accuracy disparities across demographic groups (Buolamwini et al., 2018), raising both technical and ethical questions about their use in law enforcement and public surveillance (Garvie et al., 2016). As AI systems become more deeply embedded in public and commercial infrastructures, the tension between utility and privacy continues to intensify. Without stronger safeguards, transparency, and privacy-preserving techniques, AI risks normalizing pervasive surveillance and eroding individual privacy.

## 4.3 Disinformation and Synthetic Media

The proliferation of generative models has dramatically transformed the landscape of disinformation. Deepfakes exemplify the capacity of AI systems to generate highly realistic yet fabricated content, including videos, audio, and images. These technologies have been used to create non-consensual intimate imagery, impersonate public officials, and manipulate political discourse (Chesney et al., 2019), (Vaccari et al., 2020). Large language models further extend this threat by enabling automated production of persuasive, coherent text at unprecedented scale (Goldstein et al., 2023). The

convergence of these technologies enables campaigns of influence that are more targeted, scalable, and difficult to attribute than traditional forms of propaganda.

The societal consequences of synthetic media are amplified by what Chesney and Citron (2019) describe as the “*liar’s dividend*”, whereby the mere existence of deepfakes undermines trust in authentic information (Chesney et al., 2019). Thus, AI-driven disinformation poses not only direct harm by spreading falsehoods but also indirect harm by eroding epistemic trust - the shared confidence in sources of knowledge - within societies.

#### 4.4 Bias and Discrimination

The embedding of bias into AI systems represents one of the most significant ethical challenges in contemporary deployment. Bias can arise at any stage of the machine learning pipeline, from the framing of research questions to data collection and algorithmic optimization (Barocas et al., 2016). Empirical evidence has repeatedly demonstrated how these biases translate into discriminatory outcomes. Buolamwini and Gebru (2018) showed that commercial facial recognition systems misclassified darker-skinned women at rates far higher than lighter-skinned subjects (Buolamwini et al., 2018). Obermeyer et al. (2019) also identified racial bias in healthcare algorithms that systematically underestimated the needs of Black patients (Obermeyer et al., 2019). Similarly, Angwin et al. (2016) documented how criminal justice risk assessment systems produced racially skewed predictions (Angwin et al., 2016). Such accuracy disparities have tangible real-world consequences. For example, in 2020, Robert Williams became the first documented case of wrongful arrest due to a facial recognition error, after Detroit Police Department’s system generated a false match (Evans, 2022).

Mitigating bias remains a profound challenge. Debiasing strategies, such as re-weighting datasets or modifying loss functions, have achieved partial success (Corbett-Davies et al., 2023), (Mehrabi et al., 2021). Yet scholars caution that fairness is a contested and multidimensional concept, with different definitions often mathematically incompatible (Green et al., 2020). Moreover, technical fixes alone cannot address the structural inequalities that biases both reflect and reinforce.

#### 4.5 System Safety & Reliability Failures

AI has become a dual-use technology in cybersecurity, serving both defensive and offensive roles. On the defensive side, machine learning enhances intrusion detection systems, anomaly detection, and malware classification. On the offensive side, adversaries have leveraged AI to automate phishing campaigns, discover software vulnerabilities, and craft adaptive malware (Brundage et al., 2018), (Apruzzese et al., 2018).

The emergence of large language models intensifies these threats by lowering the technical barriers to entry. Yao et al. demonstrated that such models can be prompted to generate functional malicious code, while Perez and Ribeiro showed how adversarial prompting can circumvent built-in safeguards (Yao et al., 2024), (Perez et al., 2022). These capabilities enable attackers with limited expertise to mount sophisticated operations, thereby expanding the threat landscape.

Apart from these, AI systems deployed in safety-critical applications present risks of catastrophic failures when models behave unexpectedly or incorrectly under operational conditions. For instance. Autonomous vehicles have been involved in multiple accidents resulting from perception failures, planning errors, and inadequate handling of edge cases.

## 4.6 Socioeconomic Exploitation and Inequality

AI technologies have significant implications for labor markets and economic structures. Automation driven by AI has displaced workers across various sectors, from manufacturing to customer service (Brynjolfsson et al., 2014), (Acemoglu et al., 2020). While some argue that new job categories will emerge, the transition period creates substantial economic disruption and exacerbates inequality (Autor et al., 2015).

AI also enables new forms of economic manipulation, including algorithmic pricing collusion, predatory microtargeting, and exploitation of vulnerable populations through personalized manipulation (Susser et al., 2019), (Calvano et al., 2020). These applications raise concerns about fairness, autonomy, and the concentration of economic power.

## 4.7 Environmental and Ecological Misuse

The environmental impact of AI training and deployment has gained increasing attention. Large-scale model training requires substantial computational resources, resulting in significant energy consumption and carbon emissions (Strubell et al., 2019). Additionally, AI can be misused to optimize environmentally harmful activities or bypass environmental regulations (Crawford et al., 2018). These risks are exacerbated by the growing scale and accessibility of AI technologies, which make it easier for actors with limited oversight to exploit systems in ways that harm ecological sustainability.

## 4.8 Autonomous Weaponization

Perhaps the most controversial domain of AI misuse concerns its application in military and defense systems. Lethal autonomous weapons systems (LAWS) have been identified as a critical area of concern, as they raise profound ethical, legal, and strategic dilemmas (Scharre, 2018). Scholars argue that delegating life-and-death decisions to machines undermines human accountability, risks lowering thresholds for armed conflict, and destabilizes international security (Russell, 2019). Despite calls for international regulation, progress toward binding agreements has been limited (Campaign to Stop Killer Robots, 2020).

Beyond lethal systems, AI has also been deployed for intelligence analysis, logistics optimization, and cyber operations, illustrating its broader role in military applications. The dual-use nature of these technologies complicates regulation, since advances intended for civilian purposes can be readily adapted for warfare (Horowitz et al., 2018).

Apart from these, agentic AI systems introduce novel attack vectors through their capacity for autonomous reasoning, tool use, and multi-step task execution. Unlike traditional AI systems that operate within narrowly defined boundaries, agentic systems can pursue goals through complex action sequences with minimal human oversight, creating opportunities for misuse. Agentic systems can autonomously chain together multiple attack steps, such as reconnaissance, exploitation, lateral movement, and data exfiltration, without requiring human intervention at each stage, challenging static security measures designed for simpler models (Ferrag et al., 2025), (Shrestha et al., 2025). Moreover, these systems can learn and adapt their strategies in real-time based on defensive responses, making static security measures less effective. With access to APIs, code execution environments, and system tools, agentic AI can misuse legitimate functionality to achieve unauthorized objectives, potentially escalating privileges through logical reasoning rather than traditional exploitation. Recent demonstrations have shown proof-of-concept scenarios where LLM-based agents autonomously

373 exploit vulnerabilities, conduct social engineering, or manipulate financial systems (Zhang et al.,  
374 2025). While large-scale malicious deployment remains limited, the rapid advancement of agentic  
375 capabilities warrants proactive security consideration.

## 376 **4.9 Human Interaction and Psychological Harm**

377 AI systems designed to engage users can have unintended psychological consequences. Persuasive AI,  
378 personalized content targeting, and immersive digital environments can lead to behavioral  
379 manipulation, addiction, and mental health impacts (Burr et al., 2020). The opacity of these systems  
380 makes it difficult for users to recognize when they are being manipulated, raising concerns about  
381 autonomy and wellbeing.

382 Although presented as distinct domains, these categories of misuse are deeply interconnected.  
383 Disinformation campaigns may be amplified by adversarially manipulated recommendation systems;  
384 bias in training data can exacerbate privacy violations; and cybersecurity threats can intersect with  
385 disinformation by spreading AI-generated propaganda through compromised platforms.  
386 Understanding these intersections is critical for developing holistic defensive strategies that address  
387 the complex ways in which AI misuse manifests across technological, social, and geopolitical contexts.

## 388 **5 Case Studies of AI Misuse**

389 This section presents detailed case studies illustrating real-world instances of AI misuse across multiple  
390 domains, highlighting technical mechanisms, impacts, responses, and lessons learned. These cases  
391 serve to contextualize the taxonomy of AI misuse described previously and underscore both the  
392 opportunities and risks inherent in AI technologies.

### 393 **5.1 Deepfake Pornography and Non-Consensual Intimate Imagery**

394 Since 2017, deepfake technology has been widely weaponized to generate non-consensual  
395 pornographic content, disproportionately targeting women, including celebrities, journalists,  
396 politicians, and private individuals (Ajder et al., 2019). Early deepfakes relied on GAN-based face-  
397 swapping models trained on publicly available images, but modern tools, such as DeepFaceLab<sup>1</sup> and  
398 commercial applications, have democratized creation, enabling realistic content creation with minimal  
399 technical expertise. Advances in model architecture and training methods have resulted in highly  
400 realistic outputs that are increasingly indistinguishable from authentic content.

401 The impact on victims is profound, including psychological distress, reputational damage, and  
402 sustained harassment. The rapid and widespread dissemination of such content online renders complete  
403 removal virtually impossible. Legal recourse remains limited in many jurisdictions, although some  
404 regions, such as Virginia, California, and the UK, have enacted laws criminalizing non-consensual  
405 deepfakes. While detection tools have emerged, they struggle to keep pace with increasingly  
406 sophisticated fakes, and platform enforcement remains inconsistent. This case highlights the  
407 inadequacy of purely reactive approaches, emphasizing the need for victim-centered strategies, robust  
408 legal frameworks, and platform accountability (MacDermott., 2025).

<sup>1</sup> DeepFaceLab is a leading software for creating deepfakes.

## 5.2 The 2024 Election Disinformation Campaigns

The 2024 US presidential election witnessed unprecedented deployment of AI-generated disinformation, including fabricated videos, AI-authored articles, and coordinated bot networks disseminating false narratives (DiResta et al., 2024). Large language models generated thousands of fake news articles and social media posts with human-level writing quality, while voice cloning enabled the creation of false audio of candidates making controversial statements. Moreover, automated accounts amplified content across platforms, and personalization algorithms targeted specific voter segments with tailored messaging.

Although direct electoral impact remains difficult to measure, the campaigns spread misinformation to millions of voters, complicated fact-checking efforts, and further eroded trust in information sources. This case underscores the scale and sophistication achievable with AI-powered disinformation and demonstrates that reactive detection approaches alone are insufficient without coordinated strategies involving platforms, governments, civil society, and technical researchers, etc. to defend users against manipulative content.

## 5.3 Clearview AI and Mass Surveillance

Clearview AI aggregated billions of facial images from social media and other publicly available sources without consent to build a facial recognition database marketed to law enforcement and private entities (Hill, 2020). The company collected approximately ten billion images with associated metadata, enabling searches for any individual across the internet from a single photograph. State-of-the-art deep learning models provided high recognition accuracy, raising concerns about pervasive surveillance and privacy violations.

The system facilitated monitoring of activists, protesters, and ordinary citizens, and disparities in accuracy generated discriminatory outcomes. Legal actions in multiple jurisdictions, including the EU, Canada, Australia, and several US states, resulted in fines and restrictions, while some law enforcement agencies ceased using the service. Nonetheless, the collected data cannot be retroactively "uncollected," and the company continues operations. This case illustrates the limitations of privacy frameworks designed for pre-AI contexts, demonstrating that proactive regulation to prevent data collection is essential, given the stark asymmetry between surveillance capability and individual privacy protection.

## 5.4 Algorithmic Bias in Healthcare Resource Allocation

In a study, Obermeyer et al. showed that a widely used algorithm in U.S. health systems systematically under-identified Black patients for enrollment into high-risk care management programs, relative to White patients with equivalent illness (Obermeyer et al., 2019). At the same risk score, Black patients were measurably sicker. The algorithm used health care costs as a proxy for medical needs, and because Black patients tend to incur lower costs for the same level of illness (due to unequal access and systemic barriers), the model underestimated their needs. In the studied sample, correcting for this bias would raise the share of Black patients flagged for extra care from 17.7 % to 46.5 %. In response, the algorithm developer committed to addressing the bias, prompting hospitals to audit other predictive tools. This case underscores how proxies correlated with sensitive attributes can encode bias, emphasizing the importance of understanding causal mechanisms rather than relying solely on correlations. It also highlights ethical considerations in defining optimization objectives and the necessity of comprehensive algorithmic auditing.

## 451 5.5 Voice-Cloning CEO Fraud

452 In March 2019, criminals exploited AI voice-cloning technology to impersonate a CEO's voice,  
453 successfully convincing a subordinate to transfer \$243,000 to fraudulent accounts (Stupp, 2019).  
454 Commercial voice synthesis tools trained on publicly available audio enabled the attackers to mimic  
455 speech patterns, tone, and accent convincingly. Beyond the immediate financial loss, the incident  
456 exposed vulnerabilities in voice-based authentication, previously considered secure, and demonstrated  
457 how AI can weaponize social engineering.

458 Organizations responded by implementing multi-factor authentication, out-of-band verification, and  
459 security training addressing voice-cloning risks. The case illustrates that AI capabilities can  
460 compromise traditional security assumptions, that low technical barriers facilitate broad exploitation,  
461 and that human factors often remain the weak link despite technical safeguards.

## 462 5.6 Adversarial Attacks on Autonomous Vehicle Systems

463 Research has demonstrated that autonomous vehicle vision systems can be misled by adversarial  
464 perturbations, such as strategically placed stickers on stop signs causing misclassification as speed  
465 limit signs (Eykholt et al., 2018). In these experiments, researchers used optimization algorithms to  
466 determine the smallest possible visual changes, that could consistently fool the vehicle's recognition  
467 model even under varying real-world conditions like different lighting, viewing angles, and distances.  
468 Although these attacks were conducted in controlled research environments rather than malicious  
469 settings, they expose fundamental weaknesses in safety-critical AI systems and highlight ongoing  
470 concerns about security, reliability, and potential misuse.

471 In 2018, an autonomous test vehicle in Tempe, Arizona, struck and killed a pedestrian, illustrating the  
472 real-world consequences of imperfect autonomous systems (Penmetsa et al., 2021). Tesla's Autopilot  
473 has also been involved in numerous crashes, some fatal, often occurring when the system fails to detect  
474 stationary obstacles, misinterprets road geometry, etc. The US National Transportation Safety Board  
475 has documented cases where drivers over-relied on automation and failed to maintain attention as  
476 required (Chu et al., 2023). Developers have begun incorporating adversarial training and robustness  
477 testing, yet comprehensive solutions remain elusive. These incidents emphasize that AI vulnerabilities  
478 extend from digital to physical domains, requiring security considerations from the design stage and  
479 defense-in-depth strategies rather than reliance solely on perceptual capabilities.

## 480 6 Mitigation Strategies and Evaluation

481 Mitigation of AI-related risks requires a multifaceted approach encompassing technical, regulatory,  
482 organizational, and social interventions. Each of these approaches is discussed in the following sections  
483 and a summary of the strategies is provided in Table 2.

### 484 6.1 Technical Countermeasures

485 *Adversarial robustness* techniques aim to improve the resilience of machine learning models against  
486 manipulative inputs. Adversarial training, which involves augmenting training datasets with  
487 adversarial examples (as shown in Fig.6), has demonstrated moderate effectiveness in enhancing  
488 robustness against known attacks; however, it struggles against adaptive adversaries and novel attack  
489 methods (Madry et al., 2017). This approach incurs significant computational costs that scale with the



complexity of the threat model and often involves trade-offs between accuracy and robustness. Consequently, it is most suitable for high-value targets where computational overheads are acceptable.

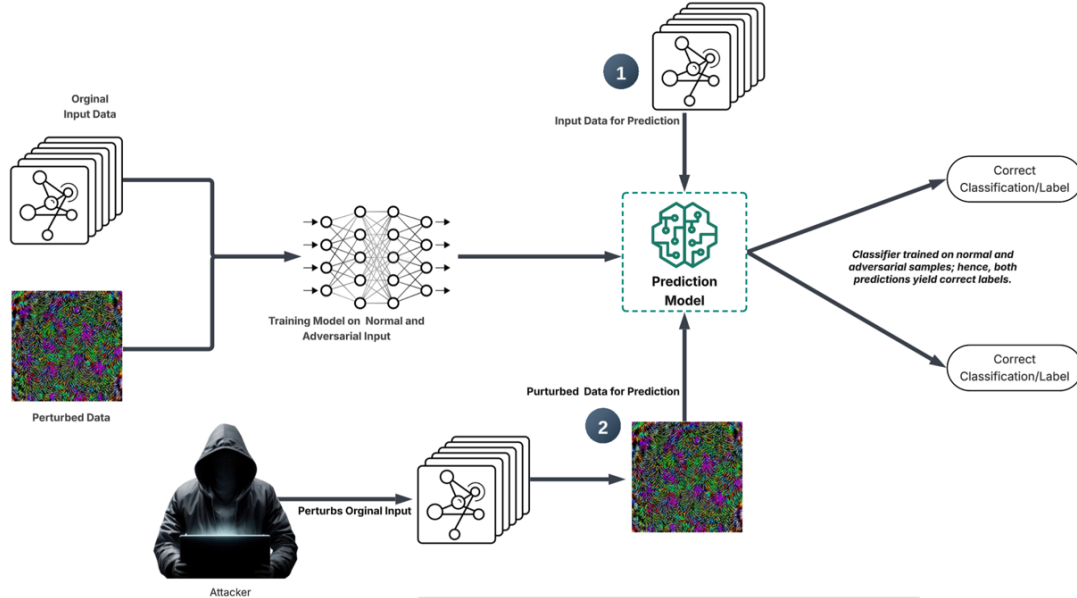


Figure 6. Adversarial training

Defensive distillation, which trains models with softened probability distributions to smooth decision boundaries, initially appeared promising (Papernot et al., 2016). While it may provide a layer of defense-in-depth, it is insufficient when deployed in isolation considering adaptive attacks. Input preprocessing methods, such as denoising, feature squeezing, or JPEG compression, can neutralize certain perturbations (Guo et al., 2017), yet these techniques degrade legitimate inputs and are easily circumvented by adaptive attackers.

Certified defenses also offer provable robustness guarantees within specified perturbation bounds, providing high theoretical value but with substantial practical limitations, including reduced accuracy and significant computational requirements (Cohen et al., 2019). Overall, no single technique currently provides comprehensive protection, and a defense-in-depth strategy combining multiple approaches represents the most viable option, albeit with persistent real-world limitations.

Apart from these, *deepfake detection technologies* have emerged to address the proliferation of synthetic media. Biological signal analysis, which detects irregularities in eye blinking, pulse, or breathing patterns, was moderately effective against early deepfakes (Wang et al., 2019) but is increasingly circumvented as generation techniques improve. GAN fingerprint detection can identify model-specific artifacts left by generative networks (Yu et al., 2019), proving useful for forensic attribution of known generators; however, it fails against unseen generators and adaptive attacks.

Temporal consistency analysis exploits frame-to-frame inconsistencies in video deepfakes, offering moderate effectiveness, particularly for video contents (Sabir et al., 2019). However, its utility diminishes as generation methods evolve. Multimodal inconsistency detection evaluates audio-visual synchronization and semantic coherence (Mittal et al., 2020), showing promise against poorly constructed deepfakes, though high-quality content often maintains consistency. Blockchain and cryptographic authentication can create verifiable chains of custody for authentic media (Hasan et al.,

2019), providing strong authenticity guarantees but requiring adoption at the point of capture, limiting applicability to existing media.

Collectively, detection approaches face an adversarial co-evolution, suggesting that proactive authentication mechanisms may prove more effective than reactive detection, albeit requiring substantial infrastructure development.

*Privacy-preserving machine learning* approaches, including differential privacy, federated learning, homomorphic encryption, and secure multi-party computation, aim to protect sensitive data while maintaining analytical capabilities. Differential privacy offers strong theoretical guarantees by introducing calibrated noise, though it necessitates careful parameter tuning to balance privacy and utility (Dwork et al., 2014). Federated learning allows decentralized training, reducing risks associated with centralized data storage (McMahan et al., 2017), but remains vulnerable to some inference attacks and incurs communication overhead.

Homomorphic encryption enables computation on encrypted data, providing theoretically strong privacy protection (Rahman et al., 2020). But this may be computationally prohibitive for complex operations. Secure multi-party computation facilitates joint computation without revealing individual inputs, offering robust privacy guarantees at the cost of significant communication and computational requirements. Overall, privacy-preserving techniques present effective protection but involve trade-offs in utility, performance, and implementation complexity.

In addition, *AI safety and alignment techniques* focus on guiding model behavior to reduce harmful outputs. Bai et al. (2022) came up with “*Constitutional AI*”, a method for training a harmless AI assistant through self-improvement, without human intervention to identify harmful outputs. It incorporates explicit principles to steer decisions, showing potential in mitigating undesired outputs but requiring careful selection of values. Reinforcement learning from human feedback (RLHF) also leverages human preferences to improve alignment (Ouyang et al., 2022) yet depends heavily on feedback quality and may inherit labeler biases.

Red teaming systematically probes system vulnerabilities (Perez et al., 2022), enabling targeted mitigation, but cannot exhaustively identify all risks and is expensive. Interpretability and explainability methods aid in understanding model decision-making (Molnar et al., 2020), which is valuable for building trust and identifying potential issues; however, explanation quality varies and post-hoc interpretations may be misleading. While these techniques advance safety, they remain incomplete, underscoring the necessity of complementary approaches for high-stakes applications.

## 6.2 Regulatory and Policy Interventions

Regulatory and policy interventions constitute a foundational layer in mitigating AI risks, particularly those related to privacy, accountability, and systemic harm. *Data protection and privacy regulations* establish essential frameworks for mitigating AI risks. The General Data Protection Regulation (GDPR) in the European Union exemplifies comprehensive privacy protection (European Parliament and Council, 2016), though enforcement challenges, jurisdictional limitations, and compliance burdens persist. Sector-specific regulations, such as HIPAA, GLBA, and COPPA, provide targeted protection for sensitive contexts but create fragmented coverage and may not fully address AI-specific risks.

In response to these limitations, AI-specific regulatory initiatives have emerged to address the unique challenges posed by AI systems. The EU AI Act represents a pioneering attempt at comprehensive,

risk-based AI regulation (European Commission, 2024), though its full effectiveness remains uncertain given ongoing implementation. Algorithmic accountability requirements, including audits, impact assessments, and transparency obligations, enhance visibility into AI systems but require technical expertise and standardization. While disclosure mandates (like informing users when AI-generated content is present), contribute to transparency, they fall short of preventing harm and often encounter challenges in enforcement and compliance. Overall, AI-specific regulatory frameworks remain fragmented and incomplete, necessitating global coordination that balances innovation with protective measures.

Beyond formal regulation, *content moderation and platform governance* constitute additional layers of policy intervention.. Platform self-regulation involves companies enforcing policies on AI-generated content, disinformation, and harmful material, with effectiveness varying across platforms. Proposals to reform Section 230 in the United States aim to adjust intermediary liability, though the potential impacts remain uncertain (Kosseff, 2019). Co-regulatory approaches, combining industry self-regulation with government oversight, such as the UK Online Safety Bill, may balance flexibility with accountability but require sustained political will and operational capacity. AI both amplifies the challenges of content moderation and offers potential solutions, indicating that multi-stakeholder governance is essential.

*International cooperation* is critical for addressing AI risks that transcend borders. Initiatives such as AI safety summits and agreements, exemplified by the Bletchley Declaration (i.e. a global agreement signed by 28 countries and the EU to foster a shared understanding of the risks and opportunities of advanced AI), facilitate shared understanding but remain non-binding and vulnerable to geopolitical tensions. Arms control frameworks propose restrictions on autonomous weapons and offensive cyber-AI, offering potential efficacy if adopted and enforced, though verification and enforcement challenges persist. International standards and best practices offer guidance on AI safety and security, promoting interoperability across systems. However, adherence is typically voluntary, and these standards often struggle to keep pace with rapid technological advancements. While global collaboration is essential, it remains inadequate in fully addressing the fast-evolving risks associated with AI.

## 6.2.1 AI Risk Governance Frameworks

Within this regulatory and policy landscape, AI risk governance frameworks play a critical complementary role by translating high-level regulatory goals into structured principles, processes, and operational guidance. Unlike legally binding regulations, these frameworks are designed to support organizations in identifying, assessing, and managing AI risks throughout the system lifecycle.

The *NIST AI Risk Management Framework* (2023) adopts a practical, implementation-oriented approach focused on organizational risk management in the US (National Institute of Standards and Technology, 2023). It structures AI risks around core trustworthy AI characteristics, including validity and reliability, safety, security and resilience, accountability and transparency, fairness with managed bias, and privacy enhancement. By emphasizing continuous risk assessment, governance integration, and lifecycle management, the NIST RMF provides actionable guidance well suited for organizational adoption across diverse sectors.

At a global level, the *OECD AI Principles* (2019) offer a high-level values-based framework adopted by 42 countries, covering inclusive growth, human-centered values, transparency, robustness and safety, and accountability (Organization for Economic Co-operation and Development, 2019). These principles provide important normative foundations and have achieved broad international consensus.

Multi-stakeholder governance initiatives further extend these efforts. The *Partnership on AI* (2021) developed a framework emphasizing responsible AI development across eight impact areas, including safety and robustness, fairness and non-discrimination, transparency and accountability, privacy and security, societal and environmental well-being, human control and autonomy, professional responsibility, and the promotion of human values (Partnership on AI, 2021). By integrating perspectives from academia, industry, civil society, and policymakers, such frameworks aim to bridge ethical principles with real-world deployment challenges.

For tackling risks associated with Agentic AI systems, technical threat models have been developed alongside these major frameworks. The *MAESTRO* (Multi-Agent Environment, Security, Threat, Risk, and Outcome), threat model provides a structured approach to identifying vulnerabilities in agentic AI systems across seven key dimensions such as model manipulation, adversarial inputs, privilege escalation, supply-chain compromise, training data poisoning, robustness failures, and output integrity issues (Huang, 2025). This technical threat modeling approach complements risk frameworks by focusing specifically on attack surfaces and defensive strategies for autonomous AI systems.

While the above comprehensive frameworks provide broad coverage, some domain-specific frameworks also address unique risks in specialized contexts. The *WHO Ethics and Governance of AI for Health framework* (WHO, 2021) identifies health-specific concerns including medical data privacy in AI-assisted diagnosis, algorithmic bias in health resource allocation, and AI-enabled health misinformation. Moreover, emergence of agentic AI systems has prompted development of specialized threat models. In biosecurity, frameworks address dual-use risks where AI capabilities for beneficial biological research can be misused for designing harmful biological agents or automating synthesis of dangerous compounds, effectively lowering technical barriers for bio-threat development (de Lima, 2024), (Trotsyuk, 2024). The UK's *AI Security Institute* (AISI) has developed safety case frameworks specifically for risk mitigation in biomedical research contexts, emphasizing structured argumentation for safety claims in high-stakes domains.

Together, these governance frameworks complement regulatory interventions by offering principled, operational, and technical approaches to managing AI risk. While none fully address the breadth of AI misuse in isolation, their combined application provides essential scaffolding for mitigating risks identified throughout this review.

### 6.3 Organizational and Social Interventions

*Organizational ethics programs and responsible AI frameworks* play a crucial role in internal governance. Ethics review boards can identify and address ethical concerns prior to deployment, but their effectiveness is contingent on institutional authority and resources. Responsible AI frameworks, such as Microsoft's RAI framework or Google's AI Principles, provide structured guidance for ethical AI development, though implementation quality varies. Bias auditing and testing help detect discriminatory system behavior, enabling targeted mitigation, yet defining fairness metrics remains contested and costly. Thus, genuine institutional commitment, supported by external accountability mechanisms, is essential for efficacy.

*Education and awareness initiatives* complement technical and regulatory measures. AI literacy programs educate the public on AI capabilities, risks, and critical evaluation of AI-generated content, fostering long-term societal resilience. Professional training for developers, policymakers, and domain experts enhances AI governance and responsible development, though rapid technological evolution challenges curriculum relevance. Media literacy and critical thinking programs further strengthen

644 resilience against disinformation. While essential, educational interventions cannot provide immediate  
645 protection and require sustained investment.

646 *Transparency and accountability mechanisms* are vital for monitoring AI deployment. Algorithmic  
647 impact assessments evaluate potential societal consequences before deployment (Reisman et al., 2018),  
648 while independent algorithmic auditing identifies issues post-deployment (Raji et al., 2020).  
649 Transparency reporting enables public scrutiny of system development and performance, though  
650 concerns regarding trade secrets, information overload, and technical complexity persist. Legal  
651 protections for whistleblowers facilitate internal accountability, provided they are genuinely enforced  
652 (Brown, 2017). Overall, transparency and accountability mechanisms remain underdeveloped relative  
653 to AI's societal impact and require urgent strengthening.

654 **Table 2. Mitigation Effectiveness Summary**

Approach Category	Representative Techniques	Effectiveness	Limitations	Deployment Status
Technical - Adversarial Robustness	Adversarial training, Certified defenses	Low-Medium	Trade-offs, Adaptive adversaries	Research/Limited deployment
Technical - Detection	Deepfake detection, Anomaly detection	Medium	Arms race dynamics	Active deployment but limited
Technical - Privacy	Differential privacy, Federated learning	Medium-High	Utility costs, Complexity	Growing deployment
Technical - Safety	Constitutional AI, RLHF	Medium	Incomplete, Research ongoing	Recent deployment
Regulatory - Privacy Laws	GDPR, CCPA	Medium-High	Enforcement challenges	Active in jurisdictions
Regulatory - AI-Specific	EU AI Act, Sector rules	Unknown	Early implementation	Emerging
Regulatory - Content Moderation	Platform policies, Co-regulation	Low-Medium	Inconsistent, Capture risk	Active but inadequate
Organizational - Ethics Programs	Review boards, Impact assessments	Low-Medium	Variable commitment	Mixed adoption
Organizational - Transparency	Audits, Reporting, Documentation	Medium	Access barriers, Standardization	Growing adoption
Social - Education	AI literacy, Media literacy	Medium (long-term)	Scale challenges, Time lag	Early stage
Ecosystem - Coordination	Standards, Information sharing	Medium	Cooperation barriers	Early stage

655 Despite growing mitigation efforts, significant gaps remain because many interventions are reactive,  
656 addressing known threats while adversaries continue to innovate. Offensive AI has access to resources  
657 comparable to defensive AI, enabling attackers to rapidly adopt the latest techniques and making it  
658 challenging for defenders to keep pace. Furthermore, policy verification and enforcement are often  
659 weak or inconsistent, and differences in regulations across jurisdictions create opportunities for  
660 regulatory arbitrage.

661 Compounding these challenges, AI's rapid evolution continues to outpace regulatory, educational, and  
662 societal adaptation. Persistent technical problems, such as adversarial robustness and deepfake  
663 detection, lack comprehensive solutions, and conflicting stakeholder priorities make it difficult to  
664 balance innovation, security, and privacy, while the widespread accessibility of AI tools amplifies the  
665 challenges of scaling effective defenses. Together, these factors underscore the persistent and growing  
666 difficulties in anticipating, managing, and mitigating AI misuse.

## 667 7 Conclusion

AI technologies hold immense transformative potential, yet they also introduce significant technical, social, and systemic risks. This paper has critically examined existing mitigation strategies, revealing that while technical defenses, regulatory frameworks, and organizational measures provide partial protection, they are often reactive, fragmented, and limited against adaptive threats. The emergence of advanced capabilities such as multimodal models and autonomous agents further amplify these risks, highlighting the need for proactive, integrated, and multi-stakeholder responses. To support this effort, we introduced a comprehensive taxonomy that organizes AI misuse into nine primary domains, providing a structured framework for understanding the full spectrum of risks - from technical vulnerabilities to socio-technical harms. The case studies presented demonstrate that AI misuse has tangible, measurable impacts, disproportionately affecting marginalized populations and eroding trust in digital systems and democratic institutions.

The trajectory of AI development presents society with critical choices about the values embedded in technological systems and the governance structures that shape their deployment. While AI capabilities continue to advance rapidly, our collective capacity to govern these technologies responsibly remains significantly underdeveloped. Addressing AI misuse requires moving beyond reactive, fragmented approaches toward proactive, integrated strategies that recognize the deeply socio-technical nature of these challenges. The stakes are high: unchecked misuse threatens privacy, security, democratic integrity, social equity, and human autonomy. Yet, with coordinated effort across technical, policy, and social domains, it remains possible to steer AI development toward beneficial outcomes that respect human rights, promote fairness, and enhance societal wellbeing.

## REFERENCES

- [1] Acemoglu, D., & Restrepo, P. (2020). Robots and jobs: Evidence from US labor markets. *Journal of Political Economy*, 128(6), 2188-2244.
- [2] Agile-index.ai. (2025). Available at: <https://agile-index.ai/publications/2025>.
- [3] Ajder, H., Patrini, G., Cavalli, F., & Cullen, L. (2019). The State of Deepfakes: Landscape, Threats, and Impact. Deeptrace. Available at: <https://scirp.org/reference/referencespapers?referenceid=3622764>
- [4] Angwin, J., Larson, J., Mattu, S., & Kirchner, L. (2016). Machine bias. ProPublica, May 23.
- [5] Reuel, A., Connolly, P., Meimandi, K.J., Tewari, S., Wiatrak, J., Venkatesh, D. and Kochenderfer, M. (2025). Responsible AI in the Global Context: Maturity Model and Survey. *Proceedings of the 2025 ACM Conference on Fairness, Accountability, and Transparency*, pp.2505–2541. doi:<https://doi.org/10.1145/3715275.3732165>.
- [6] Apruzzese, G., Colajanni, M., Ferretti, L., Guido, A., & Marchetti, M. (2018). On the effectiveness of machine and deep learning for cyber security. *10th International Conference on Cyber Conflict (CyCon)*. 371-390. 10.23919/CYCON.2018.8405026.
- [7] Autor, D. H. (2015). Why are there still so many jobs? The history and future of workplace automation. *Journal of Economic Perspectives*, 29(3), pp.3–30. doi.org/10.1257/jep.29.3.3.
- [8] Bai, Y., Kadavath, S., Kundu, S., Askill, A., Kernion, J., Jones, A., ... & Kaplan, J. (2022). Constitutional AI: Harmlessness from AI feedback. arXiv preprint arXiv:2212.08073.
- [9] Barocas, S., & Selbst, A. D. (2016). Big data's disparate impact. *California Law Review*, 104, 671-732.
- [10] Biggio, B., Corona, I., Maiorca, D., Nelson, B., Šrndić, N., Laskov, P., Giacinto, G. and Roli, F. (2013). Evasion Attacks against Machine Learning at Test Time. *Advanced Information Systems Engineering*, pp.387–402. doi:10.1007/978-3-642-40994-3\_25.

- [11] Brown, A. J. (2017). Whistleblowing in the Australian public sector: Enhancing the theory and practice of internal witness management in public sector organisations. ANU Press.
- [12] Brown, T. B., Mann, B., Ryder, N., Subbiah, M., Kaplan, J., Dhariwal, P., ... & Amodei, D. (2020). Language models are few-shot learners. *Advances in Neural Information Processing Systems*, 33, 1877-1901.
- [13] Brundage, M., Avin, S., Wang, J., Belfield, H., Krueger, G., Hadfield, G., ... & Anderljung, M. (2018). *The Malicious Use of Artificial Intelligence: Forecasting, Prevention, And Mitigation*. arXiv preprint arXiv:1802.07228.
- [14] Brynjolfsson, E., & McAfee, A. (2014). The second machine age: Work, progress, and prosperity in a time of brilliant technologies. WW Norton & Company.
- [15] Buolamwini, J., & Gebru, T. (2018). Gender shades: Intersectional accuracy disparities in commercial gender classification. *Conference on Fairness, Accountability and Transparency*, 77-91.
- [16] Burr, C., Taddeo, M., & Floridi, L. (2020). The ethics of digital well-being: A thematic review. *Science and Engineering Ethics*, 26(4), 2313-2343.
- [17] Calvano, E., Calzolari, G., Denicolò, V., & Pastorello, S. (2020). Artificial intelligence, algorithmic pricing, and collusion. *American Economic Review*, 110(10), 3267-3297.
- [18] Campaign to Stop Killer Robots. (2020). Banning killer robots.
- [19] Chesney, R., & Citron, D. (2019). Deep fakes: A looming challenge for privacy, democracy, and national security. *California Law Review*, 107, 1753-1820.
- [20] Chu, Y., Liu, P. (2023). Human Factor Risks in Driving Automation Crashes. In: Krömker, H. (eds) *HCI in Mobility, Transport, and Automotive Systems. HCII 2023. Lecture Notes in Computer Science*, vol 14048. Springer, Cham. [https://doi.org/10.1007/978-3-031-35678-0\\_1](https://doi.org/10.1007/978-3-031-35678-0_1)
- [21] Cohen, J., Rosenfeld, E., & Kolter, Z. (2019). Certified adversarial robustness via randomized smoothing. *International Conference on Machine Learning*, 1310-1320.
- [22] Corbett-Davies, S., Gaebler, J.D., Nilforoshan, H., Shroff, R. & Goel, S. (2023) ‘The measure and mismeasure of fairness’, *Journal of Machine Learning Research*, 24(1), pp. 312:1–312:117. Available at: <https://www.jmlr.org/papers/volume24/21-0341/21-0341.pdf>
- [23] Crawford, K., & Joler, V. (2018). Anatomy of an AI System: The Amazon Echo as an anatomical map of human labor, data and planetary resources. AI Now Institute and Share Lab.
- [24] Critch, A., & Russell, S. (2023). Taxonomy and Analysis of Societal-Scale Risks from AI (TASRA). arXiv:2306.06924.
- [25] de Lima, R. C., Sinclair, L., Megger, R., Maciel, M. A. G., Vasconcelos, P. F. D. C., & Quaresma, J. A. S. (2024). Artificial intelligence challenges in the face of biological threats: emerging catastrophic risks for public health. *Frontiers in artificial intelligence*, 7, 1382356. <https://doi.org/10.3389/frai.2024.1382356>
- [26] DiResta, R., Grossman, S., & Schafer, B. (2024). Generative AI and the 2024 US elections. Stanford Internet Observatory.
- [27] Dwork, C., & Roth, A. (2014). The algorithmic foundations of differential privacy. *Foundations and Trends in Theoretical Computer Science*, 9(3-4), 211-407.
- [28] European Commission. (2024). The EU Artificial Intelligence Act: A Risk-Based Framework for AI Governance. Brussels: European Union.
- [29] European Parliament and Council. (2016). Regulation (EU) 2016/679 (General Data Protection Regulation). *Official Journal of the European Union*, L119, 1-88.
- [30] Evans, L. (2022). 6 GEO. L. TECH. REV. FACIAL RECOGNITION AND A SYSTEMIC EFFECTS APPROACH TO FIRST AMENDMENT COVERAGE. [online] [https://georgetownlawtechreview.org/wp-content/uploads/2022/02/Evans\\_Facial-Recognition-and-A-Systemic-Effects-Approach-to-1A-Coverage\\_formatted.pdf](https://georgetownlawtechreview.org/wp-content/uploads/2022/02/Evans_Facial-Recognition-and-A-Systemic-Effects-Approach-to-1A-Coverage_formatted.pdf)



- [31] Eykholt, K., Evtimov, I., Fernandes, E., Li, B., Rahmati, A., Xiao, C., ... & Song, D. (2018). Robust physical-world attacks on deep learning visual classification. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition, 1625-1634.
- [32] Ferrag, M. A., Tihanyi, N., Hamouda, D., Maglaras, L., Lakas, A., & Debbah, M. (2025). From prompt injections to protocol exploits: Threats in LLM-powered AI Agents workflows. *ICT Express*. <https://doi.org/10.1016/j.ictexpress.2025.12.001>
- [33] Fredrikson, M., Jha, S., & Ristenpart, T. (2015). Model inversion attacks that exploit confidence information and basic countermeasures. Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security, 1322-1333.
- [34] Garvie, C., Bedoya, A., & Frankle, J. (2016). The perpetual line-up: Unregulated police face recognition in America. Georgetown Law Center on Privacy & Technology.
- [35] Goldstein, J. A., Sastry, G., Musser, M., DiResta, R., Gentzel, M., & Sedova, K. (2023). Generative language models and automated influence operations: Emerging threats and potential mitigations. arXiv preprint arXiv:2301.04246.
- [36] D. Golpayegani, J. Hovsha, L. W. S. Rossmaier, R. Saniei and J. Mišić, "Towards a Taxonomy of AI Risks in the Health Domain," 2022 Fourth International Conference on Transdisciplinary AI (TransAI), Laguna Hills, CA, USA, 2022, pp. 1-8, doi: 10.1109/TransAI54797.2022.00007.
- [37] Green, B., & Viljoen, S. (2020). Algorithmic realism: Expanding the boundaries of algorithmic thought. Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency, 19-31.
- [38] Gu, T., Dolan-Gavitt, B., & Garg, S. (2017). Badnets: Identifying vulnerabilities in the machine learning model supply chain. arXiv preprint arXiv:1708.06733.
- [39] Guo, C., Rana, M., Cisse, M., & Van Der Maaten, L. (2017). Countering adversarial images using input transformations. arXiv preprint arXiv:1711.00117.
- [40] Hasan, H. R., & Salah, K. (2019). Combating deepfake videos using blockchain and smart contracts. *IEEE Access*, 7, 41596-41606.
- [41] Hill, K. (2020). The secretive company that might end privacy as we know it. *The New York Times*, January 18.
- [42] Horowitz, M. C., Allen, G. C., Saravalle, E., Cho, A., Frederick, K., & Scharre, P. (2018). Artificial intelligence and international security. Center for a New American Security.
- [43] Hu, H., Salicic, Z., Sun, L., Dobbie, G., Yu, P.S. and Zhang, X. (2022). Membership Inference Attacks on Machine Learning: A Survey. *ACM Computing Surveys*. doi:<https://doi.org/10.1145/3523273>.
- [44] Huang, K. (2025, February 6). *Agentic AI Threat Modeling Framework: MAESTRO*. Cloud Security Alliance. <https://cloudsecurityalliance.org/blog/2025/02/06/agentic-ai-threat-modeling-framework-maestro>
- [45] Kosseff, J. (2019). The twenty-six words that created the internet. Cornell University Press.
- [46] LeCun, Y., Bengio, Y., & Hinton, G. (2015). Deep learning. *Nature*, 521(7553), 436-444.
- [47] Madry, A., Makelov, A., Schmidt, L., Tsipras, D., & Vladu, A. (2017). Towards deep learning models resistant to adversarial attacks. arXiv preprint arXiv:1706.06083.
- [48] Marchal, N., Xu, R., Elasmara, R., Gabriel, I., Goldberg, B., & Isaac, W. (2024). Generative AI Misuse: A Taxonomy of Tactics and Insights from Real-World Data. arXiv:2406.13843.
- [49] Matz, S. C., Kosinski, M., Nave, G., & Stillwell, D. J. (2017). Psychological targeting as an effective approach to digital mass persuasion. *Proceedings of the National Academy of Sciences*, 114(48), 12714-12719.

- [50] McGregor, S. (2021). Preventing Repeated Real World AI Failures by Cataloging Incidents: The AI Incident Database. *Proceedings of the AAAI Conference on Artificial Intelligence*, 35(17), 15458-15463. <https://doi.org/10.1609/aaai.v35i17.17817>
- [51] McMahan, B., Moore, E., Ramage, D., Hampson, S., & y Arcas, B. A. (2017). Communication-efficient learning of deep networks from decentralized data. *Artificial Intelligence and Statistics*, 1273-1282.
- [52] Mehrabi, N., Morstatter, F., Saxena, N., Lerman, K., & Galstyan, A. (2021). A survey on bias and fairness in machine learning. *ACM Computing Surveys*, 54(6), 1-35.
- [53] Mittal, T., Bhattacharya, U., Chandra, R., Bera, A., & Manocha, D. (2020). Emotions don't lie: An audio-visual deepfake detection method using affective cues. *Proceedings of the 28th ACM International Conference on Multimedia*, 2823-2832.
- [54] Molnar, C., Casalicchio, G., Bischl, B. (2020). Interpretable Machine Learning – A Brief History, State-of-the-Art and Challenges. In: Koprinska, I., *et al.* *ECML PKDD 2020 Workshops. ECML PKDD 2020. Communications in Computer and Information Science*, vol 1323. Springer, Cham. [https://doi.org/10.1007/978-3-030-65965-3\\_28](https://doi.org/10.1007/978-3-030-65965-3_28)
- [55] National Institute of Standards and Technology (NIST). (2023). AI Risk Management Framework (AI RMF 1.0). U.S. Department of Commerce.
- [56] Obermeyer, Z., Powers, B., Vogeli, C., & Mullainathan, S. (2019). Dissecting racial bias in an algorithm used to manage the health of populations. *Science*, 366(6464), 447-453.
- [57] Organisation for Economic Co-operation and Development (OECD). (2019). *OECD AI Principles*. Paris: OECD Publishing.
- [58] OECD.AI. (2024). OECD AI Incidents Monitor. <https://oecd.ai/en/incidents>
- [59] Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C., Mishkin, P., ... & Lowe, R. (2022). Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35, 27730-27744.
- [60] OWASP Foundation. (2024). OWASP Top 10 for Large Language Model Applications. <https://owasp.org/www-project-top-10-for-large-language-model-applications/>
- [61] Papernot, N., McDaniel, P., Goodfellow, I., Jha, S., Celik, Z. B., & Swami, A. (2017). Practical black-box attacks against machine learning. *Proceedings of the 2017 ACM Asia Conference on Computer and Communications Security*, 506-519.
- [62] Papernot, N., McDaniel, P., Wu, X., Jha, S., & Swami, A. (2016). Distillation as a defense to adversarial perturbations against deep neural networks. *2016 IEEE Symposium on Security and Privacy*, 582-597.
- [63] Partnership on AI. (2021). *Multi-Stakeholder Framework for Responsible AI Development*.
- [64] Penmetsa, P., Sheinidashtegol, P., Musaev, A., Adanu, E.K. and Hudnall, M. (2021). Effects of the Autonomous Vehicle Crashes on Public Perception of the Technology. *LATSS Research*, 45(4). doi:<https://doi.org/10.1016/j.iatssr.2021.04.003>.
- [65] Perez, E., Ringer, S., Lukošiūtė, K., Nguyen, K., Chen, E., Heiner, S., ... & Kaplan, J. (2022). Discovering language model behaviors with model-written evaluations. *arXiv preprint arXiv:2212.09251*.
- [66] Perez, F., & Ribeiro, I. (2022). Ignore previous prompt: Attack techniques for language models. *arXiv preprint arXiv:2211.09527*.
- [67] Rahman, M.S., Khalil, I., Atiquzzaman, M. and Yi, X. (2020). Towards privacy preserving AI based composition framework in edge networks using fully homomorphic encryption. *Engineering Applications of Artificial Intelligence*, 94, p.103737. doi:<https://doi.org/10.1016/j.engappai.2020.103737>.

- [68] Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T., Hutchinson, B., ... & Barnes, P. (2020). Closing the AI accountability gap: Defining an end-to-end framework for internal algorithmic auditing. *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*, 33-44.
- [69] Reisman, D., Schultz, J., Crawford, K., & Whittaker, M. (2018). Algorithmic impact assessments: A practical framework for public agency accountability. AI Now Institute.
- [70] Russell, S. (2019). *Human compatible: Artificial intelligence and the problem of control*. Viking.
- [71] Sabir, E., Cheng, J., Jaiswal, A., AbdAlmageed, W., Masi, I., & Natarajan, P. (2019). Recurrent convolutional strategies for face manipulation detection in videos. *Interfaces*, 3(1), 80-87.
- [72] Scharre, P. (2018). *Army of none: Autonomous weapons and the future of war*. WW Norton & Company.
- [73] Sharif, M., Bhagavatula, S., Bauer, L., & Reiter, M. K. (2016). Accessorize to a crime: Real and stealthy attacks on state-of-the-art face recognition. *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. Association for Computing Machinery, New York, NY, USA, 1528–1540. <https://doi.org/10.1145/2976749.2978392>
- [74] Shrestha, S., Banda, C., Mishra, A. K., Djebbar, F., & Puthal, D. (2025). Investigation of Cybersecurity Bottlenecks of AI Agents in Industrial Automation. *Computers*, 14(11), 456. <https://doi.org/10.3390/computers14110456>
- [75] Slattery, P., Saeri, A. K., Grundy, E. A. C., Graham, J., Noetel, M., Uuk, R., Dao, J., Pour, S., Casper, S., & Thompson, N. (2024). *The AI Risk Repository: A Comprehensive Meta-Review, Database, and Taxonomy of Risks from Artificial Intelligence*. <https://airisk.mit.edu/>
- [76] Stanford HAI. (2025). *The 2025 AI Index Report*. Stanford.edu. Available at: <https://hai.stanford.edu/ai-index/2025-ai-index-report>.
- [77] Strubell, E., Ganesh, A. and McCallum, A. (2019). Energy and Policy Considerations for Deep Learning in NLP. *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp.3645–3650. doi:<https://doi.org/10.18653/v1/p19-1355>.
- [78] Stupp, C. (2019). *Fraudsters Used AI to Mimic CEO's Voice in Unusual Cybercrime Case*. [online] Wall Street Journal. Available at: <https://www.wsj.com/articles/fraudsters-use-ai-to-mimic-ceos-voice-in-unusual-cybercrime-case-11567157402>.
- [79] Susser, D., Roessler, B. and Nissenbaum, H.F. (2019). Online Manipulation: Hidden Influences in a Digital World. *SSRN Electronic Journal*, 4(1). doi:<https://doi.org/10.2139/ssrn.3306006>.
- [80] Trotsyuk, A.A., Waeiss, Q., Bhatia, R.T. et al. Toward a framework for risk mitigation of potential misuse of artificial intelligence in biomedical research. *Nat Mach Intell* 6, 1435–1442 (2024). <https://doi.org/10.1038/s42256-024-00926-3>
- [81] Tramèr, F., Zhang, F., Juels, A., Reiter, M. K., & Ristenpart, T. (2016). Stealing machine learning models via prediction APIs. *25th USENIX Security Symposium*, 601-618.
- [82] Vaccari, C., & Chadwick, A. (2020). Deepfakes and disinformation: Exploring the impact of synthetic political video on deception, uncertainty, and trust in news. *Social Media + Society*, 6(1). <https://doi.org/10.1177/2056305120903408>
- [83] Wang, R., Juefei-Xu, F., Ma, L., Xie, X., Huang, Y., Wang, J., & Liu, Y. (2019). *FakeSpotter: A Simple yet Robust Baseline for Spotting AI-Synthesized Fake Faces*. *Proceedings of the 29th International Joint Conference on Artificial Intelligence (IJCAI 2020)*, pp. 3444-3451. doi:10.24963/ijcai.2020/476
- [84] Weidinger, L., Uesato, J., Rauh, M., Griffin, C., Huang, P.-S., Mellor, J., Glaese, A., Cheng, M., Balle, B., Kasirzadeh, A., Biles, C., Brown, S., Kenton, Z., Hawkins, W., Stepleton, T., Birhane, A., Hendricks, L.A., Rimell, L., Isaac, W. and Haas, J. (2022). Taxonomy of Risks

- posed by Language Models. *2022 ACM Conference on Fairness, Accountability, and Transparency*. doi:<https://doi.org/10.1145/3531146.3533088>.
- [85] WHO. (2021). Ethics and governance of artificial intelligence for health. World Health Organization.
- [86] Yao, Y., Duan, J., Xu, K., Cai, Y., Sun, Z. and Zhang, Y. (2024). A survey on Large Language Model (LLM) security and privacy: The Good, The Bad, and The Ugly. *High-Confidence Computing*, 4(2), p.100211. doi:<https://doi.org/10.1016/j.hcc.2024.100211>.
- [87] Yu, N., Davis, L. S., & Fritz, M. (2019). Attributing fake images to GANs: *Learning and analyzing GAN fingerprints*. Proceedings of the IEEE/CVF International Conference on Computer Vision, 7556-7566. doi:<https://doi.org/10.1109/ICCV.2019.00765>.
- [88] Zhang, J., Bu, H., Wen, H. *et al*. When LLMs meet cybersecurity: a systematic literature review. *Cybersecurity* **8**, 55 (2025). <https://doi.org/10.1186/s42400-025-00361-w>
- [89] Zhang, R., Li, H., Meng, H., Zhan, J., Gan, H. and Lee, Y.-C., 2025. *The dark side of AI companionship: A taxonomy of harmful algorithmic behaviors in human-AI relationships*. In: *Proceedings of the 2025 CHI Conference on Human Factors in Computing Systems (CHI '25)*. New York, NY, USA: Association for Computing Machinery, Article 13, pp.1–17. doi:10.1145/3706598.3713429
- [90] MacDermott, Á. Deepfake Forensics: Exploring the Impact and Implications of Fabricated Media in Digital Forensic Investigations, DFRWS EU 2025, Brno, Czech Republic, Apr. 2025. [Online]. Available: [https://dfrws.org/wp-content/uploads/2025/04/DFRWS\\_EU\\_2025\\_paperposter\\_115.pdf](https://dfrws.org/wp-content/uploads/2025/04/DFRWS_EU_2025_paperposter_115.pdf) [dfrws.org] [https://dfrws.org/wp-content/uploads/2025/04/DFRWS\\_EU\\_2025\\_paperposter\\_115.pdf](https://dfrws.org/wp-content/uploads/2025/04/DFRWS_EU_2025_paperposter_115.pdf)
- [91] Haddaway, N. R., Page, M. J., Pritchard, C. C., & McGuinness, L. A. (2022). PRISMA2020: An R package and Shiny app for producing PRISMA 2020-compliant flow diagrams, with interactivity for optimized digital transparency and Open Synthesis Campbell Systematic Reviews, 18, e1230. <https://doi.org/10.1002/cl2.1230>

## 8 Conflict of Interest

*The authors declare that the research was conducted in the absence of any commercial or financial relationships that could be construed as a potential conflict of interest.*

## 9 Funding

This study is supported by Research Incentive Funds (activity codes: R23064, R21096), Research Office, Zayed University, Dubai, United Arab Emirates.