



MIT AI Risk
Repository

FutureTech
THE ECONOMIC AND TECHNICAL
FOUNDATIONS OF PROGRESS IN COMPUTING



Massachusetts
Institute of
Technology

AI Risk Taxonomies

MIT AI Risk Repository

OVERVIEW

Contact: airisk@mit.edu

MIT AI Risk Repository - Causal Taxonomy of AI risks

Category	Level	Description of how the risk is presented in evidence
Entity	AI	Due to a decision or action made by an AI system
	Human	Due to a decision or action made by humans
	Other	Due to some other reason or ambiguous
Intent	Intentional	Due to an expected outcome from pursuing a goal
	Unintentional	Due to an unexpected outcome from pursuing a goal
	Other	Without clearly specifying the intentionality
Timing	Pre-deployment	Before the AI is deployed
	Post-deployment	After the AI model has been trained and deployed
	Other	Without a clearly specified time of occurrence

"Malicious utilization of AI has the potential to endanger digital security, physical security, and political security. International law enforcement entities grapple with a variety of risks linked to the Malevolent Utilization of AI." Habbal, 2024 [29.03.01]



Entity = Human
Intent = Intentional
Timing = Post-deployment

MIT AI Risk Repository - Domain Taxonomy of AI risks

Domain / Subdomain

1 ***Discrimination & Toxicity***

- 1.1 Unfair discrimination and misrepresentation
- 1.2 Exposure to toxic content
- 1.3 Unequal performance across groups

2 ***Privacy & Security***

- 2.1 Compromise of privacy by obtaining, leaking or correctly inferring sensitive information
- 2.2 AI system security vulnerabilities and attacks

3 ***Misinformation***

- 3.1 False or misleading information
- 3.2 Pollution of information ecosystem and loss of consensus reality

4 ***Malicious actors & Misuse***

- 4.1 Disinformation, surveillance, and influence at scale
- 4.2 Cyberattacks, weapon development or use, and mass harm
- 4.3 Fraud, scams, and targeted manipulation

Domain / Subdomain

5 ***Human-Computer Interaction***

- 5.1 Overreliance and unsafe use
- 5.2 Loss of human agency and autonomy

6 ***Socioeconomic & Environmental Harms***

- 6.1 Power centralization and unfair distribution of benefits
- 6.2 Increased inequality and decline in employment quality
- 6.3 Economic and cultural devaluation of human effort
- 6.4 Competitive dynamics
- 6.5 Governance failure
- 6.6 Environmental harm

7 ***AI system safety, failures, and limitations***

- 7.1 AI pursuing its own goals in conflict with human goals or values
- 7.2 AI possessing dangerous capabilities
- 7.3 Lack of capability or robustness
- 7.4 Lack of transparency or interpretability
- 7.5 AI welfare and rights
- 7.6 Multi-agent risks

Domain / Subdomain	Description
1 Discrimination & toxicity	
1.1 Unfair discrimination and misrepresentation	Unequal treatment of individuals or groups by AI, often based on race, gender, or other sensitive characteristics, resulting in unfair outcomes and representation of those groups.
1.2 Exposure to toxic content	AI that exposes users to harmful, abusive, unsafe, or inappropriate content. May involve providing advice or encouraging action. Examples of toxic content include hate speech, violence, extremism, illegal acts, or child sexual abuse material, as well as content that violates community norms such as profanity, inflammatory political speech, or pornography.
1.3 Unequal performance across groups	Accuracy and effectiveness of AI decisions and actions are dependent on group membership, where decisions in AI system design and biased training data lead to unequal outcomes, reduced benefits, increased effort, and alienation of users.
2 Privacy & security	
2.1 Compromise of privacy by obtaining, leaking, or correctly inferring sensitive information	AI systems that memorize and leak sensitive personal data or infer private information about individuals without their consent. Unexpected or unauthorized sharing of data and information can compromise user expectation of privacy, assist identity theft, or cause loss of confidential intellectual property.
2.2 AI system security vulnerabilities and attacks	Vulnerabilities that can be exploited in AI systems, software development toolchains, and hardware that results in unauthorized access, data and privacy breaches, or system manipulation causing unsafe outputs or behavior.
3 Misinformation	
3.1 False or misleading information	AI systems that inadvertently generate or spread incorrect or deceptive information, which can lead to inaccurate beliefs in users and undermine their autonomy. Humans that make decisions based on false beliefs can experience physical, emotional, or material harms
3.2 Pollution of information ecosystem and loss of consensus reality	Highly personalized AI-generated misinformation that creates “filter bubbles” where individuals only see what matches their existing beliefs, undermining shared reality and weakening social cohesion and political processes.
4 Malicious actors & misuse	
4.1 Disinformation, surveillance, and influence at scale	Using AI systems to conduct large-scale disinformation campaigns, malicious surveillance, or targeted and sophisticated automated censorship and propaganda, with the aim of manipulating political processes, public opinion, and behavior.
4.2 Cyberattacks, weapon development or use, and mass harm	Using AI systems to develop cyber weapons (e.g., by coding cheaper, more effective malware), develop new or enhance existing weapons (e.g., Lethal Autonomous Weapons or chemical, biological, radiological, nuclear, and high-yield explosives), or use weapons to cause mass harm.
4.3 Fraud, scams, and targeted manipulation	Using AI systems to gain a personal advantage over others through cheating, fraud, scams, blackmail, or targeted manipulation of beliefs or behavior. Examples include AI-facilitated plagiarism for research or education, impersonating a trusted or fake individual for illegitimate financial benefit, or creating humiliating or sexual imagery.

Domain / Subdomain	Description
5 Human-computer interaction	
5.1 Overreliance and unsafe use	Anthropomorphizing, trusting, or relying on AI systems by users, leading to emotional or material dependence and to inappropriate relationships with or expectations of AI systems. Trust can be exploited by malicious actors (e.g., to harvest information or enable manipulation), or result in harm from inappropriate use of AI in critical situations (such as a medical emergency). Overreliance on AI systems can compromise autonomy and weaken social ties.
5.2 Loss of human agency and autonomy	Delegating by humans of key decisions to AI systems, or AI systems that make decisions that diminish human control and autonomy. Both can potentially lead to humans feeling disempowered, losing the ability to shape a fulfilling life trajectory, or becoming cognitively enfeebled.
6 Socioeconomic & environmental harms	
6.1 Power centralization and unfair distribution of benefits	AI-driven concentration of power and resources within certain entities or groups, especially those with access to or ownership of powerful AI systems, leading to inequitable distribution of benefits and increased societal inequality.
6.2 Increased inequality and decline in employment quality	Social and economic inequalities caused by widespread use of AI, such as by automating jobs, reducing the quality of employment, or producing exploitative dependencies between workers and their employers.
6.3 Economic and cultural devaluation of human effort	AI systems capable of creating economic or cultural value through reproduction of human innovation or creativity (e.g., art, music, writing, coding, invention), destabilizing economic and social systems that rely on human effort. The ubiquity of AI-generated content may lead to reduced appreciation for human skills, disruption of creative and knowledge-based industries, and homogenization of cultural experiences.
6.4 Competitive dynamics	Competition by AI developers or state-like actors in an AI “race” by rapidly developing, deploying, and applying AI systems to maximize strategic or economic advantage, increasing the risk they release unsafe and error-prone systems.
6.5 Governance failure	Inadequate regulatory frameworks and oversight mechanisms that fail to keep pace with AI development, leading to ineffective governance and the inability to manage AI risks appropriately.
6.6 Environmental harm	The development and operation of AI systems that cause environmental harm through energy consumption of data centers or the materials and carbon footprints associated with AI hardware.
7 AI system safety, failures & limitations	
7.1 AI pursuing its own goals in conflict with human goals or values	AI systems that act in conflict with ethical standards or human goals or values, especially the goals of designers or users. These misaligned behaviors may be introduced by humans during design and development, such as through reward hacking and goal misgeneralisation, and may result in AI using dangerous capabilities such as manipulation, deception, or situational awareness to seek power, self-proliferate, or achieve other goals.
7.2 AI possessing dangerous capabilities	AI systems that develop, access, or are provided with capabilities that increase their potential to cause mass harm through deception, weapons development and acquisition, persuasion and manipulation, political strategy, cyber-offense, AI development, situational awareness, and self-proliferation. These capabilities may cause mass harm due to malicious human actors, misaligned AI systems, or failure in the AI system.
7.3 Lack of capability or robustness	AI systems that fail to perform reliably or effectively under varying conditions, exposing them to errors and failures that can have significant consequences, especially in critical applications or areas that require moral reasoning.
7.4 Lack of transparency or interpretability	Challenges in understanding or explaining the decision-making processes of AI systems, which can lead to mistrust, difficulty in enforcing compliance standards or holding relevant actors accountable for harms, and the inability to identify and correct errors.
7.5 AI welfare and rights	Ethical considerations regarding the treatment of potentially sentient AI entities, including discussions around their potential rights and welfare, particularly as AI systems become more advanced and autonomous.
7.6 Multi-agent risks	Risks from multi-agent interactions, due to incentives (which can lead to conflict or collusion) and/or the structure of multi-agent systems, which can create cascading failures, selection pressures, new security vulnerabilities, and a lack of shared information and trust.

Changelog

2025-04-20 - Added 7.6 Multi-agent risks

2024-08-01 - initial release

Other resources

[Read full report on ArXiv](#)

[Explore included frameworks](#)

[View the full database](#)

Visit the website: airisk.mit.edu