

AI Safety

Rifky Firmansyah

 <http://orcid.org/0009-0000-9177-966X>

Universitas Esa Unggul, Indonesia

ABSTRACT

This study proposes a comprehensive framework for AI safety as the foundation for developing secure, controllable, and human-aligned artificial intelligence systems. Through a systematic review of literature from 2021 to 2025, it identifies three core pillars: technical robustness against disruptions, ethical value alignment, and socio-regulatory governance. The analysis covers real-world cases of medical diagnostic vulnerabilities, autonomous vehicle failures, and manipulation risks in large language models (LLMs). Mitigation approaches such as safe reinforcement learning, formal verification, and human-in-the-loop mechanisms are explored. The study also highlights the importance of global standards, including ISO/IEC 23894 and the NIST AI Risk Management Framework, in ensuring transparency and accountability. Findings suggest that cross-disciplinary collaboration among scientists, regulators, and civil society is crucial to building a trustworthy and ethically responsible AI ecosystem.

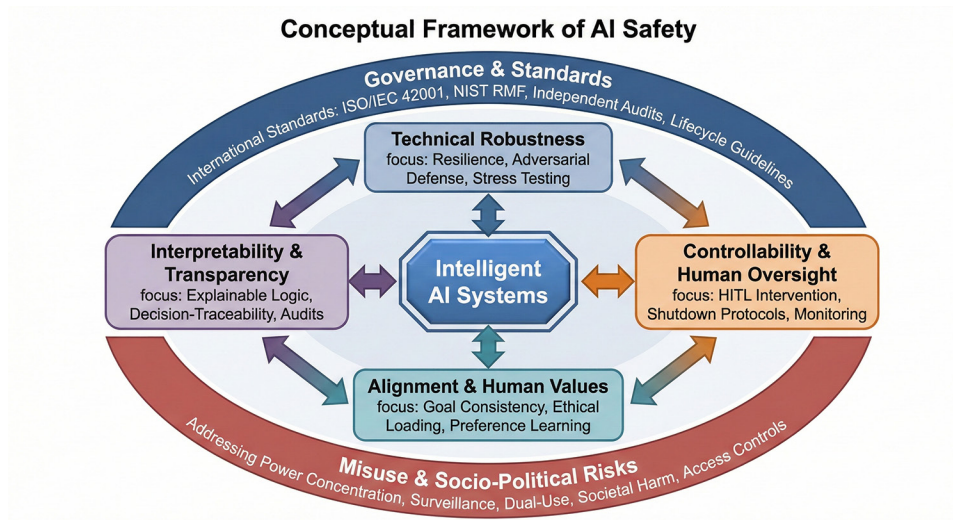
1. INTRODUCTION: AI SAFETY IN THE ERA OF COMPLEX AUTONOMOUS SYSTEMS

Artificial Intelligence (AI) systems are increasingly embedded in critical infrastructures such as healthcare, autonomous transportation, defense, and digital communication. As these systems scale in complexity and autonomy, they introduce systemic risks that extend beyond conventional software failures, including adversarial vulnerabilities, value misalignment, reduced transparency, and governance gaps. AI Safety therefore emerges as a socio-technical discipline concerned with ensuring that intelligent systems operate reliably, ethically, and accountably in real-world contexts. Addressing these challenges requires the integration of technical robustness, alignment with human values, interpretability, human oversight, and adaptive governance mechanisms across the AI system lifecycle.

DOI: 10.4018/406036

Copyright ©2027, IGI Global Scientific Publishing. Copying or distributing in print or electronic forms without written permission of IGI Global Scientific Publishing is prohibited. Use of this chapter to train generative artificial intelligence (AI) technologies is expressly prohibited. The publisher reserves all rights to license its use for generative AI training and machine learning model development.

Figure 1. Conceptual Framework of AI Safety



This figure presents AI Safety as a socio-technical framework centered on intelligent AI systems and organized around six interdependent pillars: technical robustness, alignment with human values, interpretability and transparency, controllability and human oversight, misuse and socio-political risks, and governance and standards. The framework highlights the bidirectional relationships between technical safeguards and institutional mechanisms across the AI system lifecycle.

1.1 Technical Robustness and Vulnerability of Modern AI Systems

Modern AI systems, particularly those based on deep learning architectures, have demonstrated remarkable performance across domains such as image recognition, language understanding, and autonomous control. However, this performance is often accompanied by an alarming vulnerability to small, targeted perturbations that can severely disrupt system behavior. Adversarial attacks exemplify this fragility; small, often imperceptible modifications to input data can cause high-confidence misclassifications, even in clinically validated systems. In medical imaging, adversarial noise has been shown to trigger dangerous diagnostic errors in radiology models, highlighting how seemingly minor input perturbations can undermine the reliability of safety-critical AI systems. These findings are not confined to healthcare; similar vulnerabilities exist across a range of AI applications, from autonomous vehicle perception to real-time surveillance and multimodal models.

These observations point to a broader concern regarding robustness in modern AI systems: increases in architectural complexity do not automatically translate into improved resilience against adversarial manipulation. Empirical studies in safety-critical domains indicate that highly parameterized models may exhibit brittle behavior when exposed to distributional shifts or adversarial inputs, particularly when robustness is not an explicit optimization objective. In response, emerging strategies such as model pruning and the integration of attention mechanisms have demonstrated potential for improving robustness while maintaining interpretability and performance (Chen et al., 2021). These technical solutions,

25 more pages are available in the full version of this document, which may be purchased using the "Add to Cart" button on the publisher's webpage:

www.igi-global.com/chapter/ai-safety/406036

Related Content

Color Image Segmentation of Endoscopic and Microscopic Images for Abnormality Detection in Esophagus

P. S. Hiremathand Iranna Y. Humnabad (2012). *Cross-Disciplinary Applications of Artificial Intelligence and Pattern Recognition: Advancing Technologies* (pp. 165-193).

www.irma-international.org/chapter/color-image-segmentation-endoscopic-microscopic/62690

A Hybrid Learning Framework for Imbalanced Classification

Eric P. Jiang (2022). *International Journal of Intelligent Information Technologies* (pp. 1-15).

www.irma-international.org/article/a-hybrid-learning-framework-for-imbalanced-classification/306967

A Multimodal Sentiment Analysis Model for Graphic Texts Based on Deep Feature Interaction Networks

Wanjun Changand Dongfang Zhang (2024). *International Journal of Ambient Computing and Intelligence* (pp. 1-19).

www.irma-international.org/article/a-multimodal-sentiment-analysis-model-for-graphic-texts-based-on-deep-feature-interaction-networks/355192

A Blockchain-Based Security Model for Cloud Accounting Data

Congcong Gouand Xiaoqing Deng (2023). *International Journal of Ambient Computing and Intelligence* (pp. 1-16).

www.irma-international.org/article/a-blockchain-based-security-model-for-cloud-accounting-data/332860

Artificial Intelligence Applications in Human Resource Management: A Bibliometric Content Analysis and Future Research Agenda

Chandni Raniand Tanveer Kajla (2023). *AI and Emotional Intelligence for Modern Business Management* (pp. 178-193).

www.irma-international.org/chapter/artificial-intelligence-applications-in-human-resource-management/332636