

1

Introduction and Background

On the global market, several nations are racing to achieve a global innovation advantage in AI as it is understood that AI is a foundational technology that can boost competitiveness, increase productivity, protect national security, and help solve societal challenges. Comparing China, the European Union, and the United States in terms of their relative standing in the AI economy by examining six categories of metrics: talent, research, development, adoption, data, and hardware, the United States leads in absolute terms, with China coming in second, and the European Union lags further behind. This order could change in the coming years depending on a range of policy actions that can propel each nation or region to improve its AI capabilities [11, 14]. AI technology developments significantly impact electronic and component systems, semiconductor design, and production, as the amount of data processed and stored by AI applications continues to increase. Semiconductor architectural improvements are needed to address data use in AI-integrated circuits, and improvements in semiconductor design for AI are requested to enhance overall performance, speed, memory capacity, with increased energy efficiency. In this context, major initiatives have started globally to address the development of the semiconductor industry, such as the European Chips Act, which aims to bolster Europe's competitiveness, resilience, and help achieve both the digital and green transition [12]. In this context, edge AI represents a paradigm shift in deploying and utilising AI technologies, marking a transformative evolution from centralised data processing systems to decentralised, edge-oriented solutions. The edge AI deployments are characterised by massive scale and heterogeneity. A single system may comprise many devices with diverse hardware and software stacks, creating significant challenges for deployment, interoperability, and management. These devices are often deployed in uncontrolled environments, making them vulnerable

2 *Introduction and Background*

to physical tampering and environmental hazards, which introduces a class of security threats not typically considered in secure data centres. This tightly coupled system of trade-offs, forced by a resource-scarce environment, makes a holistic systems engineering approach a must.

This transition underscores the capability of executing AI algorithms directly on intelligent edge devices and embedded systems across the edge continuum, including micro-, deep- and meta-edge, including avoiding the need for constant connectivity to cloud-based processing centres. As edge AI is developed and applied across various sectors, understanding the edge AI system's functional and non-functional requirements becomes imperative to harness the full potential of technology.

Edge AI refers to the deployment of AI algorithms, machine learning (ML), deep learning (DL) and generative AI models directly on local, interconnected edge devices, such as Internet of Things (IoT) devices, smartphones, security cameras, gateways, embedded systems and on-premises servers using edge computing processing. By processing data locally, edge AI systems can deliver real-time analytics and decision-making with significantly reduced latency, which is critical for applications like autonomous vehicles and industrial automation. This approach enhances data privacy and security by minimising the transmission of potentially sensitive information to the cloud, reduces the demand for network bandwidth and associated costs, and enables operational autonomy in environments with intermittent or non-existent network connectivity.

The design of edge AI systems is fundamentally constraint-driven. These constraints are not minor implementation details but are the primary architectural drivers that shape every design decision. The most prominent is the severe limitation on resources. Edge devices are typically constrained by processing power, available memory (RAM and flash storage), and, most critically, energy consumption, which directly impacts battery life and thermal management.

These resource limitations impose direct constraints on the AI models themselves. Large, complex models are not appropriate at the edge, which necessitates the use of lightweight model architectures and aggressive model compression techniques, such as quantisation (reducing the precision of model weights), pruning (removing unnecessary connections), and knowledge distillation (training a smaller “student” model to mimic a larger “teacher” model). While on-device training offers benefits for privacy and adaptation to new data, it is exceptionally challenging due to these same resource constraints.

Edge AI systems must also contend with network constraints. Although they are designed to operate with less reliance on the cloud, they are not entirely disconnected. The need for model updates, data synchronisation, or federated learning means that systems must be robust to environments with limited, unreliable, or costly network connectivity.

The rapid integration of AI and edge AI into various sectors has moved the technology from a research area to a driving force of digital transformation.

The technical, market and social developments in automation and integration of AI in industrial environments advance the topic of ethics of AI, dependability combined with industrial AI trust, which focuses on achieving the desired outcome for AI-based technologies and applications in various industrial sectors while complying with legal rules and adhering to ethical norms. Addressing the complex interrelation between ethics and AI comes with notable dynamics, controversial issues, a lack of standards and no common agreement on principles about ethics. In addition, trust in AI and edge AI systems has multiple dimensions combining system dependability characteristics (e.g., privacy, security, safety, reliability, availability, resilience, connectability and maintainability) with human and machine behaviour, which require a greater understanding of how individuals interact with machines and how machines/things interact with other machines/things to extend trust [70].

In these conditions, technology companies are focused on building AI platforms that meet the stakeholders' needs for optimised performance, profitability and security. In doing so, they are partnering across the AI ecosystem of semiconductor and integrated circuits companies, hyperscalers, large language models, data and software companies, and engaging with global trade policy unknowns and resource constraints. The trends in new AI frontiers and the focus on companies in these environments include AI reasoning, custom silicon, edge and cloud balanced migrations, systems to measure AI efficacy and building an agentic AI future [15].

These transitions and trends necessitate a structured approach to governance and technical oversight, which is the primary role of standardisation and regulations. Standards provide a common language and a set of established principles for developing, deploying, and maintaining AI and edge AI systems. They are crucial for ensuring that AI technologies are not only innovative but also safe, secure, reliable, and aligned with societal values [7, 115].

Regulatory bodies are increasingly looking to standards as a means to ensure compliance and manage the complexities introduced by AI

4 *Introduction and Background*

[115, 145]. For instance, the European Union’s AI Act [6, 115, 129] references harmonised standards as a mechanism for demonstrating conformity with its legal requirements.

Similarly, initiatives like the U.S. National Institute of Standards and Technology (NIST) AI Risk Management Framework (AI RMF) provide guidance that is being adopted and referenced by organisations globally [8]. This synergy between regulation and standardisation is creating a framework for responsible AI innovation, pushing developers and organisations to adopt best practices in their AI engineering [131] disciplines.

A global ecosystem of standardisation bodies is actively working to address the multifaceted challenges of AI and edge AI. Each organisation brings a unique perspective and expertise, contributing to a holistic standards landscape. Their collaborative efforts are essential for avoiding fragmentation in standardisation activities and ensuring that standards are globally relevant and applicable across diverse industries and use cases. Fragmentation in standardisation activities refers to the situation where the lack of coordination or agreement can lead to multiple, competing standards, hindering the benefits of standardisation such as interoperability and market efficiency. As a result, this can display as diverging national or regional standards, or even within the same sector or industry, creating a fragmented landscape that can be costly and complex to navigate.

The proliferation of AI and edge AI has created an urgent need to move beyond research-oriented development towards a more disciplined, engineering-focused approach. The concepts of AI and edge AI engineering [131] have emerged to meet this need, advocating for the application of established principles from systems engineering, software engineering [130], requirements engineering and human-centred design to construct AI and edge AI systems that are functional, reliable, secure, and aligned with human values and mission objectives.

This evolution from experimental development to robust engineering is occurring in parallel with, and is significantly influenced by, the global push for regulation. Governments and international bodies are grappling with the societal, ethical, and economic implications of AI, leading to landmark legislative efforts. The adoption of the Artificial Intelligence Act (AIA) by the European Union in June 2024 introduced the first horizontal rules addressing the risks to health, safety and fundamental rights posed by AI systems. The AI Act framework categorises AI systems by risk as illustrated in Figure 1.1 and imposes stringent requirements on those deemed “high-risk” [115, 6].

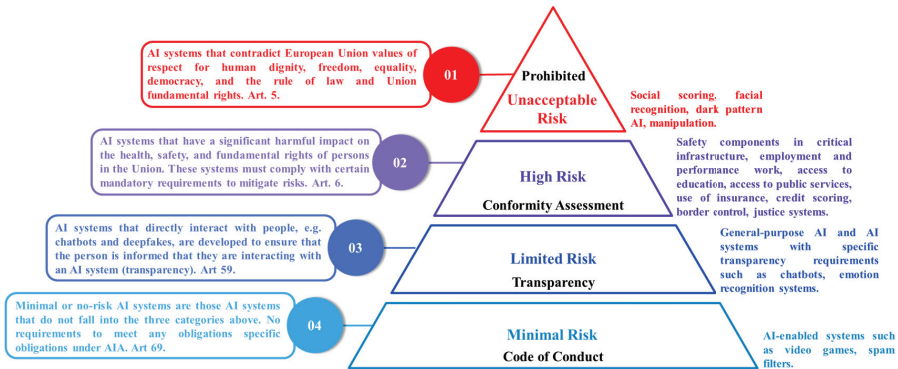


Figure 1.1 The EU AIA's pyramid of risks.

Most AI systems and applications pose minimal or no risks. Specific AI systems are subjected to transparency obligations, e.g., when they interact with natural persons or pose risks of impersonation or deception. High-risk AI systems are limited to those that could have a significant, harmful impact on the health, safety, or fundamental rights of individuals. For these, the AIA defines a clear set of requirements, as extensively discussed in this document. The AIA regulation prohibits certain AI practices that pose unacceptable risks [144].

The AIA provides for the development of harmonised European standards ('hENS') to address these risks. Adherence to these standards by AI providers will facilitate conformity assessments and provide a presumption of compliance with all or parts of the AI Act's requirements, to be assessed on a case-by-case basis.

The legislation empowers European Standards Organisations (ESOs) to develop these harmonised standards, which, when used by developers, provide a "presumption of conformity" with the law's technical requirement [9].

AI standardisation remains voluntary, also under the AIA. The market-driven nature of standardisation and differ depending on ESOs internal organisation (i.e., industry representation vs. national representation). The linkage between law and technical specification is transforming the AI and edge AI standardisation and can influence the market access.

The result is a dynamic and complex global standardisation landscape, where multiple Standards Development Organisations (SDOs) are working to create the necessary frameworks, guidelines, and technical specifications.

6 *Introduction and Background*

This report provides an exhaustive overview of these efforts, analysing the work of the International Organization for Standardization (ISO), the International Electrotechnical Commission (IEC), the Institute of Electrical and Electronics Engineers (IEEE), the International Telecommunication Union - Telecommunication Standardization Sector (ITU-T), the European Telecommunications Standards Institute (ETSI), and the European Committee for Standardization (CEN) and European Committee for Electrotechnical Standardization (CENELEC). Each organisation brings a unique perspective and focus, contributing distinct yet complementary pieces to the global puzzle of AI governance.

This book aims to structure and explain the comprehensive set of requirements essential for the successful realisation of edge AI systems.

Requirements in edge AI systems are classified into functional requirements (FRs), which define the functionality that the system must offer or what the system should do, and non-functional requirements (NFRs), known as quality requirements, that represent the desired qualities of the system, how well it should perform its functions, and include aspects like performance, reliability, and security. Constraints are pre-existing limitations on the design or development process, use of specific technology or standard, that cannot be influenced by the engineering design.

The book defines the functional and non-functional requirements and provides examples of their KPIs, measures (quantitative, qualitative), monitoring, ongoing assessment, transparency, and alignment with societal values and needs.

The FRs discussed address the core features, capabilities and tasks that edge AI systems perform to fulfil their intended purpose. These capabilities are critical considering the performance, latency, bandwidth, response times, energy efficiency, data processing, machine learning and deep learning model execution and real-time decision-making at the edge.

The NFRs focus on edge AI systems' quality attributes or characteristics that impact the performance. Topics such as functional suitability, performance efficiency, compatibility, interaction capability, reliability, security, maintainability, flexibility and safety are analysed to provide a holistic view of what is needed to ensure that edge AI systems can operate effectively in diverse and often constrained environments.

The book links the definition of functional and non-functional requirements of edge AI systems with the concepts of edge AI system dependability and trustworthiness.

Dependability is essential to the performance [20] of real-time edge AI systems and reflects edge AI systems' operational requirements while reflecting the degree of trustworthiness in the system. As a result, trustworthiness is the ability of an edge AI system to meet functional and non-functional requirements in a verifiable way, meaning that it can be checked for correctness by a person or tool through verification, validation, testing, and benchmarking and supported by explainability and interpretability methods [72].

By delineating and analysing the FRs and NFRs requirements, this book provides stakeholders, including researchers, designers, developers, system architects, and decision-makers, with a framework for designing and implementing edge AI solutions that meet operational objectives while delivering sustained value and efficiency. As we delve into the specifics of these requirements, the analysis highlights the challenges and opportunities associated with deploying AI at the edge yet guiding efforts to leverage this transformative technology in real-world applications.

