



# Trustworthiness of AI

WHITE PAPER

## Landscape of AI regulations, standards and publications

Tomislav Nad, Sebastian Scher, Florian Königstorfer

June 2024

Partners of SGS



SGS

## Abstract:

The development of artificial intelligence (AI) systems is progressing rapidly and so is their adoption in many industries, products and services. There is no question that AI is, and will even further, influence our societies and lives. Due to this influence and the progress in the development and adoption of AI systems, trust, ethics and social concerns need to be addressed. AI systems need to be reliable, fair, transparent – they need to be trustworthy.

This need is recognized by many organizations from governments, industry and academia. They have discussed and are still discussing how trust in AI systems can be established.

In this context, numerous white papers, proposals and standards have been published and are still in development. For someone who is just starting to look into this topic, the number of resources can be overwhelming.

This document aims to provide a summary and guidance through the jungle of documents about trust in AI. We look at existing standards, standards in development, reports, regulations, audit and test proposals, certification guidelines as well as any other informative white paper. On each of them, we provide a short summary and put them in context to the whole publication landscape. This document should assist the reader in becoming familiar with the topic.

## Contents:

<b>1. INTRODUCTION</b>	<b>3</b>
1.1 Trustworthiness of AI systems by the EU	4
1.2 Trustworthiness of AI systems by NIST	5
1.3 Trustworthiness of AI systems by OECD	5
1.4 Trustworthiness of AI systems by World Economic Forum (WEF)	6
1.5 Trustworthiness of AI systems by ChatGPT	6
1.6 Outline	7
<b>2. PUBLICATION LANDSCAPE</b>	<b>7</b>
2.1 Overview of overviews	7
2.2 Mapping publications to trust requirements	9
2.3 Standards and technical reports	18
2.4 Legal regulations	18
2.4.1 EU Artificial Intelligence Act	18
2.4.2 Artificial Intelligence Liability Directive	19
2.4.3 Blueprint for an AI Bill of Rights	19
2.4.4 UK Policy Paper	19
2.4.5 Japan AI Guidelines	19
2.5 White papers and reports	20
2.6 Audit catalogues	20
<b>3. CONCLUSIONS</b>	<b>21</b>
<b>4. ABOUT</b>	<b>22</b>
4.1 Authors	22
4.2 Organizations	22
<b>5. APPENDIX</b>	<b>24</b>
5.1 Table of published standards and reports	24
5.2 Table of standards and reports in development	29
5.3 Table of white papers and reports	33



## Introduction

Artificial intelligence (AI) promises to bring many benefits and hence is increasingly being adopted by a broad range of industries. With wide adoption comes the need to address trust, ethics and social concerns. AI systems need to ensure reliability, fairness and transparency – they need to be trustworthy.

In recent years, many organizations from governments, industry and academia have discussed the trustworthiness of AI systems. Numerous white papers, proposals and standards have been published and are still in development.

This demonstrates that the trustworthiness of AI systems is taken seriously across the globe and that interest in this topic is shared amongst all types of players in the field.

The most prominent activity with respect to trust in AI is undoubtedly the European Union (EU) AI Act<sup>1</sup>. It is the first comprehensive proposal for a regulation dealing with the risks related to the development and application of AI. However, it is by far not the only relevant document about trust in AI. Almost a countless number of publications have been published in recent years discussing this topic from different angles and with different interests in mind. For those who start to look into trustworthiness of AI themselves, it can be easily overwhelming, and one struggles with finding out where to start and identifying what is relevant. Considering that sooner or later everyone who develops or uses AI systems will be affected in some form by this discussion, it is important to understand what is going on, what the state of the art is and where future developments will head.

This document aims to provide a summary and guidance through the jungle of documents about trust in AI. We look at existing standards, standards in development, reports, regulations, audit and test proposals, certification guidelines as well as any other informative white paper. On each of them we provide a short summary and put them in context to the whole publication landscape.

Before we dive into the matter, we would like to set the stage by introducing the relevant concepts for trustworthy AI.

**AI systems**  
need to ensure  
reliability, fairness  
and transparency  
– they need to be  
**trustworthy.**

[1] [https://www.europarl.europa.eu/doceo/document/TA-9-2024-0138-FNL-COR01\\_EN.pdf](https://www.europarl.europa.eu/doceo/document/TA-9-2024-0138-FNL-COR01_EN.pdf)



## 1.1 Trustworthiness of AI systems by the EU

We start with the definition of what an AI system is. In its final version, the EU AI Act<sup>1</sup> provides a definition in Article 3 (1):

**‘AI system’ means a machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments.**

In this article we follow the broad definition of the EU.

Next, we need to define what a trustworthy AI system is. Therefore, we are citing the definitions developed by the High-Level Expert Group on Artificial Intelligence (AI HLEG). This independent expert group was set up by the European Commission (EC) in June 2018.

They published the Ethics Guidelines for Trustworthy Artificial Intelligence<sup>2</sup> which defines that trustworthy AI should have three components:

1. It should be lawful, complying with all applicable laws and regulations
2. It should be ethical, ensuring adherence to ethical principles and values, and
3. It should be robust, both from a technical and social perspective, since, even with good intentions, AI systems can cause unintentional harm

These led to the seven key requirements for trustworthy AI (shown in Figure 1):

1. Human agency and oversight including fundamental rights, human agency and human oversight
2. Technical robustness and safety including resilience to attack and security, fall back plan and general safety, accuracy, reliability and reproducibility
3. Privacy and data governance including respect for privacy, quality and integrity of data, and access to data
4. Transparency including traceability, explainability and communication
5. Diversity, non-discrimination and fairness including the avoidance of unfair bias, accessibility and universal design, and stakeholder participation
6. Societal and environmental wellbeing including sustainability and environmental friendliness, social impact, society and democracy
7. Accountability including auditability, minimization and reporting of negative impact, trade-offs and redress

---

[2] <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>

**AI** should be robust, both from a technical and social perspective, since, even with good intentions, AI systems can cause unintentional harm.

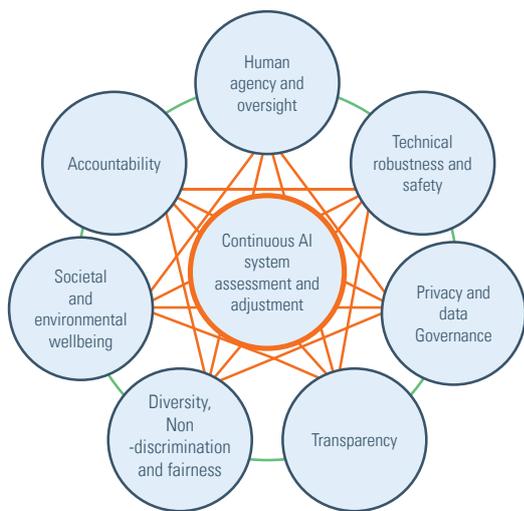


Figure 1: Requirements for trustworthy AI<sup>3</sup>

So, in other words, trust in AI can be established by sufficiently addressing these requirements and characteristics of AI systems. However, there is not always a clear and unambiguous definition of the key requirements. For example, the requirement of “fairness” is heavily defined by the context and what is considered to be fair might be sometimes intuitively easy to determine, but as discussed in a white paper about certifying fairness<sup>4</sup>, most often it is not.

Note that all publications on trustworthy AI address at least one of those aspects, where some even address all of them. It is important to note that not all publications break it down to the exact same seven aspects, however, they usually can be mapped to the definitions above.

## 1.2 Trustworthiness of AI systems by NIST

The National Institute of Standards and Technology (NIST) recently published the Artificial Intelligence Risk Management Framework (AI RMF 1.0)<sup>5</sup>. They define several characteristics of trustworthy AI systems as shown in Figure 2. As one can see, they are quite similar to the key requirements provided above.

[3] <https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai>

[4] <https://trustyour.ai/en/whitepaper/certifying-fairness-of-ai-applications-an-impossible-task/>

[5] <https://doi.org/10.6028/NIST.AI.100-1>



Figure 2: Characteristics of trustworthy AI systems by NIST<sup>6</sup>

## 1.3 Trustworthiness of AI systems by OECD

As an additional comparison, we would like to cite the five OECD AI Principles<sup>7</sup>. They can be mapped to the above seven key requirements.



[6] <https://doi.org/10.6028/NIST.AI.100-1>

[7] <https://oecd.ai/en/ai-principles>

## 1.4 Trustworthiness of AI systems by World Economic Forum (WEF)

As a more general example for digital trust, we would like to cite World Economic Forum and their take on trustworthy technologies<sup>9</sup> shown in Figure 4. Here the authors define digital trust as follows:

Individuals' expectation that digital technologies and services – and the organizations providing them – will protect all stakeholders' interests and uphold societal expectations and values.

They define three goals, namely security and reliability, accountability and oversight, and inclusive, ethical and responsible use. They are supported by several methods as shown in Figure 4.



Figure 4: Digital Trust Framework<sup>9</sup>

As we see, there is a common understanding between different organizations and stakeholders about digital trust and hence the trustworthiness of AI systems. In the following sections, we will refer to these requirements when discussing various publications around them.

**AI should be lawful, complying with all applicable laws and regulations.**

[8] <https://oecd.ai/en/ai-principles>

[9] <https://www.weforum.org/reports/earning-digital-trust-decision-making-for-trustworthy-technologies/>

## 1.5 Trustworthiness of AI Systems by ChatGPT

Finally, we asked ChatGPT<sup>10</sup> what trustworthy AI means.

### What does trustworthy artificial intelligence mean?

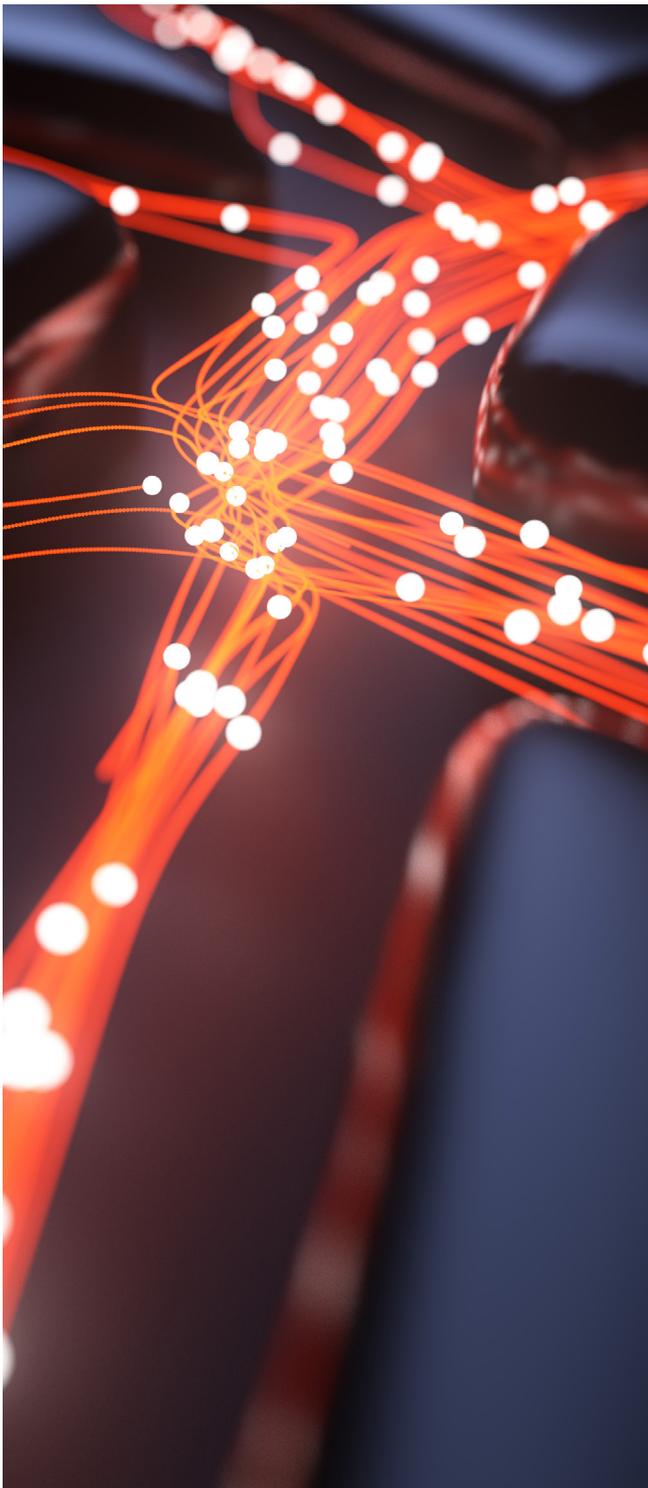
Trustworthy artificial intelligence (AI) refers to AI systems that are designed, developed and deployed in ways that ensure they are reliable, ethical and beneficial to society. Trustworthy AI is characterized by several key principles:

1. **Fairness:** AI systems should be fair and free from bias. They should not discriminate against individuals or groups based on race, gender, age or other protected characteristics. Ensuring fairness often involves rigorous testing and validation to identify and mitigate biases in data and algorithms
2. **Transparency:** AI systems should be transparent in their operations. This includes being clear about how decisions are made, providing understandable explanations for those decisions and allowing for scrutiny and auditability. Transparency helps users understand the AI's behavior and builds trust
3. **Accountability:** Developers, operators and users of AI systems should be accountable for their actions and the impacts of the AI systems. This involves establishing clear lines of responsibility and ensuring that there are mechanisms in place to address any harm or issues that arise from the use of AI
4. **Privacy:** AI systems must respect and protect the privacy of individuals. This includes implementing strong data protection measures, minimizing the amount of personal data collected and ensuring that data is used responsibly and in compliance with relevant laws and regulations
5. **Security:** AI systems must be secure from malicious attacks and unauthorized access. This involves implementing robust security measures to protect data and algorithms, as well as continuously monitoring and updating systems to address vulnerabilities
6. **Reliability:** AI systems should be reliable and perform consistently under various conditions. They should be tested extensively to ensure they function as intended and can handle unexpected inputs or situations without failure
7. **Ethical Use:** AI systems should be used ethically, in ways that align with societal values and norms. This includes considering the broader implications of AI deployment and ensuring that the technology is used to enhance human well-being and does not cause harm
8. **Human Oversight:** Trustworthy AI involves human oversight to ensure that AI systems are used appropriately and that their decisions can be challenged and corrected if necessary. Human oversight helps prevent misuse and allows for the integration of human judgment in critical areas

9. Sustainability: AI systems should be developed and used in ways that are environmentally sustainable. This involves considering the environmental impact of AI technologies, such as energy consumption and resource use, and striving to minimize their footprint

By adhering to these principles, AI systems can be more trustworthy, fostering greater confidence and acceptance among users and society at large.

[10] <https://chat.openai.com/>



## 1.6 Outline

The rest of this document is structured as follows. In section 2.1, we start with a list of overview documents similar to this one but with a different focus. In section 2.2, we provide a summary of the listed publications and put them in various context such that the reader may identify more quickly which document might be of interest. In section 2.3, we provide an overview of standards and technical reports published by standardization bodies. In section 2.4, we provide an overview of legal regulations. In section 2.5, we provide an overview of general white papers and reports published by various organizations. In section 2.6, we list audit catalogues relevant for AI. Finally, in section 3, we conclude and provide the complete list of included publications in Appendix 5.

## Publication landscape

### 2.1 Overview of overviews

This report is not the first of its kind. There have been several publications which provide an excellent overview of activities and publications related to the trustworthiness of AI. Each of them has a specific focus and are a great read in general. Hence, we would like to start to cite those that we discovered:

- The Stanford University runs a 100-year study on the development of AI and publishes a great report on a regular basis. The most recent one<sup>11</sup> has been published in September, 2021. It provides an overview about the progress in AI and related topics
- The StandICT.eu project published an overview document<sup>12</sup> about AI related standards. It is a static snapshot of a dynamically updated database compiled within StandICT.eu
- The European Union Agency for Cybersecurity (ENISA) recently published a report<sup>13</sup> which describes the standardization landscape covering AI with a focus on cybersecurity, by depicting the activities of the main standards organizations
- AI Watch, the EC knowledge service to monitor the development of AI, published a report<sup>14</sup> focusing on the mapping of the AI standards onto the requirements introduced by the EU AI Act
- RAND Europe published a report on labeling initiatives, codes of conduct and other self-regulatory mechanisms for artificial intelligence applications<sup>15</sup>. It shows how self-regulatory mechanisms related to trustworthiness of AI are being used in this context

[11] <https://ai100.stanford.edu/>

[12] [https://zenodo.org/record/5011179#.Y\\_f1HCbMJD8](https://zenodo.org/record/5011179#.Y_f1HCbMJD8)

[13] <https://www.enisa.europa.eu/publications/cybersecurity-of-ai-and-standardisation>

[14] <https://op.europa.eu/en/publication-detail/-/publication/36c46b8e-e518-11eb-a1a5-01aa75ed71a1/language-en/format-PDF>

[15] [https://www.rand.org/pubs/research\\_reports/RRA1773-1.html](https://www.rand.org/pubs/research_reports/RRA1773-1.html)



## 2.2 Mapping publications to trust requirements

In the following section, we provide a summary of all discussed publications in this report and map them in Table 1 to the seven key requirements as listed in section 1. Furthermore, we provide another view on these publications, mapping each of them to the potential target audience in Table 2. Finally, we provide a quick overview of the required AI expertise level for each of these publications, in Table 3.

This should allow the reader to quickly identify relevant documents for their own situation.

Note that some documents may not directly address a key aspect but is of general nature which can make it relevant for all key aspects.

**AI** should be ethical, ensuring adherence to ethical principles and values.

	Human agency and oversight	Technical robustness and safety	Privacy and data governance	Transparency	Diversity, non-discrimination and fairness	Societal and environmental wellbeing	Accountability
<b>Standards and technical reports</b>							
ISO/IEC 5338	x	x	x	x	x	x	
ISO/IEC 5339							
ISO/IEC 5392							
ISO/IEC 5469		x					
ISO/IEC 8183			x				
ISO/IEC 8200	x						
ISO/IEC 24027					x		
ISO/IEC 24028	x	x	x	x	x	x	x
ISO/IEC 24029		x					
ISO/IEC 24030							
ISO/IEC 24372							
ISO/IEC 38507	x	x	x	x	x	x	x
ISO 22100-5		x					
ISO/IEC 4213		x					
ISO/IEC 17903							
ISO/IEC 20546							
ISO/IEC 20547-1							
ISO/IEC 20547-2							
ISO/IEC 20547-3							
ISO/IEC 20547-4			x				
ISO/IEC 20547-5							
ISO/IEC 22989	x	x	x	x	x	x	x
ISO/IEC 23053	x	x	x	x	x	x	x
ISO/IEC 23894	x	x	x	x	x	x	x
ISO/IEC 24029-2		x					
ISO/IEC 24368					x	x	
ISO/IEC 24668			x				

	Human agency and oversight	Technical robustness and safety	Privacy and data governance	Transparency	Diversity, non-discrimination and fairness	Societal and environmental wellbeing	Accountability
ISO/IEC 25058	x	x	x	x	x	x	x
ISO/IEC 25059		x					
ISO/IEC 29119-11		x					
ISO/IEC 42001	x	x	x	x	x	x	x
DIN SPEC 92001-1	x	x	x	x	x	x	x
DIN SPEC 92001-2		x					
IEEE 2801			x				
IEEE 2802		x					
IEEE 2976				x			
IEEE 7000					x	x	
IEEE 7002			x				
IEEE 7001				x			
IEEE 7007						x	
IEEE 7010						x	
IEEE 7014						x	
IEEE 7015							
NIST AI Risk Management Framework	x	x	x	x	x	x	x
Standards and technical reports under development							
ISO/IEC 5259-1		x	x				
ISO/IEC 5259-2		x	x				
ISO/IEC 5259-3		x	x				
ISO/IEC 5259-4		x	x				
ISO/IEC 5259-5		x	x				
ISO/IEC 5259-6		x	x				
ISO/IEC 6254				x			
ISO/IEC 12791					x		
ISO/IEC 12792				x			
ISO/IEC 17847		x					
ISO/IEC 18988							
ISO/IEC 20226						x	
ISO/IEC 21221							
ISO/IEC 22440		x					
ISO/IEC 22443						x	
ISO/IEC 23281							
ISO/IEC 23282		x					
ISO/IEC 24029-3		x					
ISO/IEC 24970	x	x	x	x	x	x	x
ISO/IEC 25029							
ISO/IEC 42005	x	x	x	x	x	x	x

	Human agency and oversight	Technical robustness and safety	Privacy and data governance	Transparency	Diversity, non-discrimination and fairness	Societal and environmental wellbeing	Accountability
ISO/IEC 42006							
ISO/IEC 42102	x	x	x	x	x	x	x
ISO/IEC 42103		x	x				
ISO/IEC 42105	x						
ISO/IEC 42106							
ISO/IEC 27090			x				
IEEE P7003					x		
IEEE P7008						x	
<b>Regulations</b>							
EU AI Act	x	x	x	x	x	x	x
EU Liability Directive		x				x	x
Blueprint for an AI Bill of Rights	x	x	x	x	x		x
UK Policy Paper	x	x	x	x	x		x
Japan AI Guidelines	x	x	x	x	x		x
<b>White papers and reports</b>							
Towards Auditable AI Systems		x					
Artificial Intelligence Cybersecurity Challenges		x	x				
Trusted Artificial Intelligence Towards Certification of Machine Learning Applications	x	x	x	x	x	x	x
Trust and Artificial Intelligence	x	x	x	x	x	x	x
Artificial Intelligence and future directions for ETSI		x					
Securing Artificial Intelligence (SAI); Mitigation Strategy Report		x	x				
Securing Artificial Intelligence (SAI); The role of hardware in security of AI		x	x				
Securing Machine Learning Algorithms	x	x	x	x	x	x	x
Standardization Roadmap AI	x	x	x	x	x	x	x

Table 1: Mapping of publications to AI aspects

	Human agency and oversight	Technical robustness and safety	Privacy and data governance	Transparency	Diversity, non-discrimination and fairness	Societal and environmental wellbeing	Accountability
On Artificial Intelligence – A European Approach to Excellence and Trust	x	x	x	x	x	x	x
Policy and Investment Recommendations for Trustworthy Artificial Intelligence	x	x	x	x	x	x	x
Taxonomy of AI Risk	x	x	x	x	x	x	x
AI Risk Management Framework: Initial Draft	x	x	x	x	x	x	x
<b>Audit catalogues</b>							
Leitfaden zur Gestaltung ertrauenswürdiger Künstlicher Intelligenz	x	x	x	x	x	x	x
Auditing Machine Learning Algorithms- A White Paper for Public Auditors	x	x	x	x	x	x	x
Algorithmic Impact Assessment tool	x	x	x	x	x	x	x

Table 1: Mapping of publications to AI aspects

In the next table we map the publications to the assumed target audience of the document. This should further help the reader to identify the most relevant documents. “General” means, that the document is relevant for a general audience. “Developer” refers to organizations and people developing AI systems or AI components.

“Evaluator” is someone who assesses AI systems and processes. “User” is a user of an AI system.

Note that for some standards and technical reports which are still under development, not enough information is available to properly assume the target audience.

	General	Developer	Evaluator	User
<b>Standards and technical reports</b>				
ISO/IEC 5338	x	x	x	
ISO/IEC 5339		x		x
ISO/IEC 5392		x		
ISO/IEC 5469	x	x	x	x
ISO/IEC 8183	x	x	x	
ISO/IEC 8200	x	x	x	x
ISO/IEC 24027		x	x	
ISO/IEC 24028	x	x	x	x
ISO/IEC 24029-1		x	x	
ISO/IEC 24030	x	x	x	x
ISO/IEC 24372	x	x	x	
ISO/IEC 38507	x	x	x	x
ISO 22100-5		x	x	x
ISO/IEC 4213		x	x	
ISO/IEC 17903	x	x	x	
ISO/IEC 20546	x	x	x	
ISO/IEC 20547-1		x	x	
ISO/IEC 20547-2		x	x	
ISO/IEC 20547-3		x	x	
ISO/IEC 20547-4		x	x	
ISO/IEC 20547-5		x	x	
ISO/IEC 22989	x	x	x	x
ISO/IEC 23053	x	x	x	x
ISO/IEC 23894	x	x	x	x
ISO/IEC 24029-2		x	x	
ISO/IEC 24368	x	x	x	x
ISO/IEC 24668		x		
ISO/IEC 25058		x	x	
ISO/IEC 25059		x	x	
ISO/IEC 29119-11		x	x	
ISO/IEC 42001	x	x	x	x
DIN SPEC 92001-1	x	x	x	
DIN SPEC 92001-2		x	x	
IEEE 2801		x		

	General	Developer	Evaluator	User
IEEE 2802		x		
IEEE 7000		x		
IEEE 7002		x	x	
IEEE 7001		x	x	
IEEE 7007		x		
IEEE 7009		x		
IEEE 7010		x		
IEEE 7014		x		
IEEE 7015	x			
IEEE 2976		x	x	
NIST AI Risk Management Framework	x	x	x	x
Standards and technical reports under development				
ISO/IEC 5259-1		x		
ISO/IEC 5259-2		x		
ISO/IEC 5259-3		x		
ISO/IEC 5259-4		x		
ISO/IEC 5259-5		x		
ISO/IEC 5259-6		x		
ISO/IEC 6254	x	x	x	x
ISO/IEC 12791		x	x	
ISO/IEC 12792	x	x	x	x
ISO/IEC 17847		x	x	
ISO/IEC 18988				
ISO/IEC 20226	x	x	x	x
ISO/IEC 21221				
ISO/IEC 22440	x	x	x	x
ISO/IEC 22443	x	x	x	x
ISO/IEC 23281		x	x	
ISO/IEC 23282		x	x	
ISO/IEC 24029-3		x	x	
ISO/IEC 24970		x	x	
ISO/IEC 25029		x	x	
ISO/IEC 42005	x	x	x	x
ISO/IEC 42006			x	
ISO/IEC 42102	x	x	x	x
ISO/IEC 42103	x	x	x	x
ISO/IEC 42105	x	x	x	x
ISO/IEC 42106		x	x	
ISO/IEC 27090	x			x
IEEE 7003		x	x	x
IEEE 7008		x	x	

	General	Developer	Evaluator	User
<b>White papers and reports</b>				
Towards Auditable AI Systems		x	x	
Artificial Intelligence Cybersecurity Challenges		x		
Trusted Artificial Intelligence Towards Certification of Machine Learning Applications		x	x	
Trust and Artificial Intelligence	x	x		
Artificial Intelligence and Future Directions for ETSI		x	x	
Securing Artificial Intelligence (SAI); Mitigation Strategy Report		x		
Securing Artificial Intelligence (SAI); The Role of Hardware in Security of AI		x		
Securing Machine Learning Algorithms	x	x	x	
Standardization Roadmap AI	x	x	x	
On Artificial Intelligence – A European Approach to Excellence and Trust	x			
Policy and Investment Recommendations for Trustworthy Artificial Intelligence	x			
Taxonomy of AI Risk		x	x	
AI Risk Management Framework: Initial Draft		x	x	
<b>Audit catalogues</b>				
Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher Intelligenz		x	x	
Auditing Machine Learning Algorithms – A White Paper for Public Auditors		x	x	
Algorithmic Impact Assessment tool		x	x	x
<b>Regulations</b>				
EU AI Act	x	x	x	x
EU Liability Directive	x	x	x	x
Blueprint for an AI Bill of Rights	x	x	x	x
UK Policy Paper	x	x	x	x
Japan AI Guidelines	x	x	x	x

Table 2: Mapping of publications to target audience

In the next table we try to determine the minimum required level of AI expertise needed to follow the content documents. There is not enough information on the

standards and technical reports under development, hence they are not listed. Note that some documents require specific technical and legal expertise in addition.

	Novice	Competent	Proficient	User
<b>Standards and technical reports</b>				
ISO/IEC 5338	x			
ISO/IEC 5339	x			
ISO/IEC 5392	x			
ISO/IEC 5469	x			
ISO/IEC 8183	x			
ISO/IEC 8200	x			
ISO/IEC 24027			x	
ISO/IEC 24028	x			
ISO/IEC 24029-1			x	
ISO/IEC 24029-2				x
ISO/IEC 24030	x			
ISO/IEC 24372		x		
ISO/IEC 38507	x			
ISO 22100-5	x			
ISO/IEC 4213	x			
ISO/IEC 17903				x
ISO/IEC 20546		x		
ISO/IEC 20547-1		x		
ISO/IEC 20547-2		x		
ISO/IEC 20547-3		x		
ISO/IEC 20547-4		x		
ISO/IEC 20547-5		x		
ISO/IEC 22989	x			
ISO/IEC 23053	x			
ISO/IEC 23894	x			
ISO/IEC 24058			x	
ISO/IEC 24059			x	
ISO/IEC 24368	x			
ISO/IEC 24668	x			
ISO/IEC 42001	x			
ISO/IEC 29119-11			x	
DIN SPEC 92001-1	x			
DIN SPEC 92001-2			x	
IEEE 2801	x			
IEEE 2802	x			
IEEE 7000	x			
IEEE 7002	x			
IEEE 7001	x			
IEEE 7007			x	
IEEE 7009			x	
IEEE 7010	x			

	Novice	Competent	Proficient	User
IEEE 7014			x	
IEEE 7015			x	
IEEE 2976			x	
NIST AI Risk Management Framework	x			
<b>White papers and reports</b>				
Towards Auditable AI Systems			x	
Artificial Intelligence Cybersecurity Challenges				
Trusted Artificial Intelligence Towards Certification of Machine Learning Applications		x		
Trust and Artificial Intelligence		x		
Artificial Intelligence and Future directions for ETSI	x			
Securing Artificial Intelligence (SAI); Mitigation Strategy Report			x	
Securing Artificial Intelligence (SAI); The Role of Hardware in Security of AI			x	
Securing Machine Learning Algorithms	x			
Standardization Roadmap AI	x			
On Artificial Intelligence – A European Approach to Excellence and Trust	x			
Policy and Investment Recommendations for Trustworthy Artificial Intelligence	x			
Taxonomy of AI Risk		x		
AI Risk Management Framework: Initial Draft		x		
<b>Audit catalogues</b>				
Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher Intelligenz			x	
Auditing Machine Learning Algorithms - A White Paper for Public Auditors		x		
Algorithmic Impact Assessment tool		x		
<b>Regulations</b>				
EU AI Act	x			
EU Liability Directive	x			
Blueprint for an AI Bill of Rights	x			
UK Policy Paper	x			
Japan AI Guidelines	x			

Table 3: Mapping of publications requiring AI expertise level



## 2.3 Standards and technical reports

Many standards bodies are working on AI related topics and have published numerous documents. Most of them are technical reports which provide guidance for different stakeholders and input for the standardization process. Only a few of them are standards for example related to governance and the AI life cycle. However, many are currently in development by various bodies.

In this section we provide an overview of the various publications by the following standards organizations: ETSI, NIST, ISO, IEEE and DIN. We include a short summary to be able to put the content in context. We also mention the target audience for the document, i.e. if the content is targeting developers, users or evaluators. Lastly, we map the key aspect(s) of AI the document discusses. The full list of published standards and reports can be found in Appendix 5.1 and for those still in development in Appendix 5.2. All of the document titles and document abstracts/summaries included here are copied or directly derived from publicly available materials, however copyright of those original materials is retained by the respective owners.

Please note that there are more standards, reports and activities related to various aspects of AI systems. ETSI published a report<sup>16</sup>, which summarizes ETSI's ongoing and planned efforts in the field of AI. ISO/IEC runs a committee (JTC 1/SC 42) with focus on the standardization of AI. IEEE published several standards on technical aspects of AI and runs the initiative Autonomous and Intelligent Systems, which bundles their AI related activities. NIST runs several AI related programs with the goal to cultivate trust in AI technologies. DIN has recently published the second version of the Standardization Roadmap AI<sup>17</sup>, which comprehensively discusses the required actions for standards and certifications.

For more details of related activities, we refer to the corresponding web pages (ETSI, IEEE<sup>18,19</sup>, ISO<sup>20</sup>, DIN<sup>21</sup> and NIST<sup>22</sup>). In this document we focus on publications which directly address the trustworthiness of AI.

## 2.4 Legal regulations

Due to the novelty of AI, only a few AI specific regulations exist. However, policy makers have recognized the obvious need and started to propose AI regulations. ETSI provides<sup>16</sup> an overview of world-wide regulation actions. In this section we briefly summarize these policy activities.

### 2.4.1 EU Artificial Intelligence Act

The most prominent regulation is the EU AI Act<sup>23</sup>.

The proposed regulatory framework on AI has the following objectives:

- Ensure that AI systems placed on the Union market and used are safe and respect existing law on fundamental rights and Union values
- Ensure legal certainty to facilitate investment and innovation in AI
- Enhance governance and effective enforcement of existing law on fundamental rights and safety requirements applicable to AI systems
- Facilitate the development of a single market for lawful, safe and trustworthy AI applications and prevent market fragmentation

[16] <https://www.etsi.org/images/files/ETSIWhitePapers/ETSI-WP52-ETSI-activities-in-the-field-of-AI.pdf>

[17] <https://www.din.de/de/forschung-und-innovation/themen/kuenstliche-intelligenz/fahrplan-festlegen>

[18] <https://standards.ieee.org/initiatives/artificial-intelligence-systems/standards/>

[19] <https://standards.ieee.org/initiatives/autonomous-intelligence-systems/>

[20] <https://www.iso.org/committee/6794475.html>

[21] <https://www.din.de/en/innovation-and-research/artificial-intelligence>

[22] <https://www.nist.gov/artificial-intelligence>

[23] [https://www.europarl.europa.eu/RegData/commissions/imco/inag/2024/02-02/CJ40\\_AG\(2024\)758862\\_EN.pdf](https://www.europarl.europa.eu/RegData/commissions/imco/inag/2024/02-02/CJ40_AG(2024)758862_EN.pdf)

The EU AI Act defines different levels of risk and imposes different restrictions and requirements on AI applications with different levels of risk. A visualization of these categories can be seen in Figure 5.

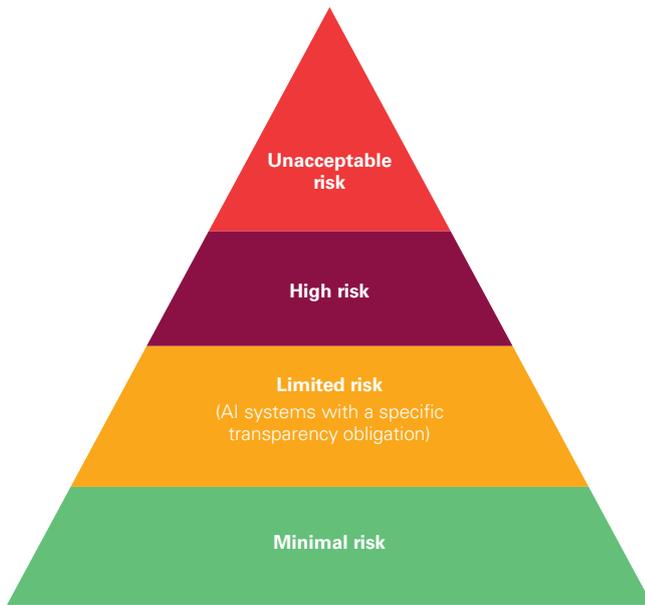


Figure 5: Levels of AI risks by EU AI Act

The EU AI Act focuses on the high-risk area where there is a high-risk to the health and safety or fundamental rights of natural persons. Annex III of the EU AI Act contains examples of such systems, e.g.

- **Biometrics, insofar as their use is permitted under relevant Union or national law:**
  - (a) Remote biometric identification systems. This shall not include AI systems intended to be used for biometric verification whose sole purpose is to confirm that a specific natural person is the person he or she claims to be;
    - (aa) AI systems intended to be used for biometric categorisation, according to sensitive or protected attributes or characteristics based on the inference of those attributes or characteristics;
    - (ab) AI systems intended to be used for emotion recognition.
- **Critical infrastructure:**
  - (a) AI systems intended to be used as safety components in the management and operation of critical digital infrastructure, road traffic and the supply of water, gas, heating and electricity.

## 2.4.2 Artificial Intelligence Liability Directive

In addition to the EU AI Act, the EC has published a proposal on a directive for liability rules for AI<sup>24</sup>. This will require the EU member states to harmonize their legislation for liability cases involving AI. The goal of the directive is (1) to make it easier for individuals who received damage from an AI system to claim their rights and (2) to provide a more harmonized legal framework across all member states of the EU, thus making it easier for AI providers to comply with the rules across different member states.

Even though the EC has not yet decided on the exact policy that it will propose, it is clear that the proposal will contain

Ensure that **AI systems** placed on the Union market and used are safe and respect existing law on fundamental rights and Union values.

measures that reduce a victim's burden to prove that the AI application has caused them harm. Additionally, the EC's preferred policy option also includes a review mechanism that enables the EC to reassess the need for harmonizing strict liability for AI use cases with a particular risk profile.

[24] [https://www.europarl.europa.eu/RegData/etudes/BRIE/2023/739342/EPRS\\_BRI\(2023\)739342\\_EN.pdf](https://www.europarl.europa.eu/RegData/etudes/BRIE/2023/739342/EPRS_BRI(2023)739342_EN.pdf)

## 2.4.3 Blueprint for an AI Bill of Rights

In the US, the White House Office of Science and Technology Policy (OSTP) published a "Blueprint for an AI Bill of Rights"<sup>25</sup>. This is a set of comprehensive guidelines for designing, using and deploying AI systems. It is comprised of five principles (safe and effective systems; algorithmic discrimination protections; data privacy; notice and explanation; human alternatives, consideration and fallback). Contrary to its name, it has no legal character.

## 2.4.4 UK Policy Paper

The UK is proposing a pro-innovation framework for regulating AI in their policy paper, "Establishing a pro-innovation approach to regulating AI"<sup>26</sup>. It has no legal binding.

## 2.4.5 Japan AI Guidelines

Japan published governance guidelines for implementation of AI principles<sup>27</sup>. It presents recommendations for AI systems operators and developers for the implementation of AI principles. It has no legal binding.



## 2.5 White papers and reports

The concept of auditing, testing and certifying AI applications has gotten much attention in the literature. A wide range of institutions – ranging from organizations that are already at their core close to auditing topics, such as the German Federal Office for Information Security (BSI)<sup>29</sup> to organizations from very different areas, such as the Austrian Chamber of Labor<sup>30</sup>, have published reports and white papers on the topic. Some of these reports and white papers have a broad focus on all relevant dimensions, while some focus on specific issues (e.g. security). For brevity, we describe all publications as “paper”, independently of whether they are called a white paper, report, or similar, by their authors. It is impossible to provide a complete overview of all publications that deal in some way with the topic of trust in AI. Here we give an overview of the most important publications.

A high-level approach to all topics surrounding, requirements, (necessary) legislation for AI, including governance and liability questions is given by the white paper of the EC on AI. The same high level expert group made recommendations for necessary policy and investment regarding trustworthy AI<sup>32</sup>.

BSI has proposed the “Certification Readiness Matrix”. This gives an overview on the state-of-the art on which technical aspects at which phases of the lifecycle of AI applications can already be audited with current methods. The European Union Agency for Cybersecurity (ENISA) published a report “Artificial Intelligence Cybersecurity Challenges - Threat Landscape for Artificial Intelligence<sup>32</sup>” on cybersecurity related aspects of AI systems, with the goal of “setting the ground for defining the AI threat landscape”. It includes a comprehensive list of definitions of relevant terms, including AI assets, and maps to general threat taxonomy already used by ENISA. TÜV Austria and Johannes Kepler University Linz published a white paper on the topic of certifying machine learning applications<sup>33</sup>. While not providing detailed guidelines on how to audit AI applications, it outlines all relevant dimensions and aspects that would be relevant, and it roughly outlines a certification procedure. Finally, it gives a brief overview of a certification catalog developed by TÜV Austria. This catalog is, however, not publicly available. It deals with all relevant aspects, from technical over ethical to (cyber) security aspects. ETSI has published several reports on the topic of AI, including<sup>34</sup> an overview on standardization of AI and a report discussing the role that hardware plays in AI security<sup>35</sup>. An overview of all relevant white papers and reports is presented in Appendix 5.3.

[25] <https://www.whitehouse.gov/wp-content/uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf>

[26] <https://www.gov.uk/government/publications/establishing-a-pro-innovation-approach-to-regulating-ai/establishing-a-pro-innovation-approach-to-regulating-ai-policy-statement>

[27] [https://www.meti.go.jp/shingikai/mono\\_info\\_service/ai\\_shakai\\_jisso/pdf/20220128\\_2.pdf](https://www.meti.go.jp/shingikai/mono_info_service/ai_shakai_jisso/pdf/20220128_2.pdf)

[28] [https://www.bsi.bund.de/SharedDocs/Downloads/EN/BSI/KI/Towards\\_Auditable\\_AI\\_Systems\\_2022.html](https://www.bsi.bund.de/SharedDocs/Downloads/EN/BSI/KI/Towards_Auditable_AI_Systems_2022.html)

[29] <https://vera.arbeiterkammer.at/#/>

[30] [https://ec.europa.eu/info/sites/default/files/commission-white-paper-artificial-intelligence-feb2020\\_en.pdf](https://ec.europa.eu/info/sites/default/files/commission-white-paper-artificial-intelligence-feb2020_en.pdf)

[31] <https://digital-strategy.ec.europa.eu/en/library/policy-and-investment-recommendations-trustworthy-artificial-intelligence>

## 2.6 Audit catalogues

An audit catalogue is a list of criteria that an application or product must fulfill, as well as a guideline on how to assess these criteria. AI audit catalogues thus contain criteria and requirements that an AI system must fulfill in order to be trustworthy with respect to that specific catalogue. Currently the most comprehensive available audit catalogue is the AI Assessment Catalog (“Guideline for Designing Trustworthy Artificial Intelligence”) published by Fraunhofer IAIS<sup>36</sup>. It considers all relevant dimensions and goes into high detail. It is mainly based on auditing of documentation, not on technical tests actually carried out by the auditor. The catalogue employs a risk-based approach. For each dimension, a risk analysis must be made, and then the AI developer must show evidence and reasoning how the risk is reduced to an acceptable level. TÜV Austria in collaboration with Johannes Kepler University published on certifying ML applications<sup>37</sup>. In a collaboration of several European public auditing institutions, an audit catalogue for audits of ML-algorithms by supreme audit institutions has been published<sup>38</sup>. It is, however, more a set of guidelines than a reference catalogue.

While only a few catalogues for AI-auditing have been published, more are in the making, or are already finished but not publicly available (e.g. <sup>39,40,41,42</sup>). Since their full text is not available, we do not include them in our overview, as we cannot judge which aspects they cover and what their target audience is.

[32] <https://www.enisa.europa.eu/publications/artificial-intelligence-cybersecurity-challenges>

[33] [https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper\\_Trusted-AI\\_TUeV-AUSTRIA\\_JKU.pdf](https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper_Trusted-AI_TUeV-AUSTRIA_JKU.pdf)

[34] [https://www.etsi.org/images/files/ETSIWhitePapers/etsi\\_wp34\\_Artificial\\_Intelligence\\_and\\_future\\_directions\\_for\\_ETSI.pdf](https://www.etsi.org/images/files/ETSIWhitePapers/etsi_wp34_Artificial_Intelligence_and_future_directions_for_ETSI.pdf)

[35] [https://www.etsi.org/deliver/etsi\\_gr/SAI/001\\_099/006/01.01\\_01\\_60/gr\\_SAI006v010101p.pdf](https://www.etsi.org/deliver/etsi_gr/SAI/001_099/006/01.01_01_60/gr_SAI006v010101p.pdf)

[36] [https://www.iais.fraunhofer.de/content/dam/iais/fb/Kuenstliche-intelligenz/ki-pruefkatalog/Fraunhofer\\_IAIS\\_AI\\_ASSESSMENT\\_Catalog\\_Web.pdf](https://www.iais.fraunhofer.de/content/dam/iais/fb/Kuenstliche-intelligenz/ki-pruefkatalog/Fraunhofer_IAIS_AI_ASSESSMENT_Catalog_Web.pdf)

[37] [https://www.jku.at/fileadmin/gruppen/219/LIT\\_AI\\_Lab/News-Seite/White\\_Paper\\_-\\_Trusted\\_Artificial\\_Intelligence\\_-\\_Towards\\_Certification\\_of\\_Machine\\_Learning\\_Applications\\_web\\_s.pdf](https://www.jku.at/fileadmin/gruppen/219/LIT_AI_Lab/News-Seite/White_Paper_-_Trusted_Artificial_Intelligence_-_Towards_Certification_of_Machine_Learning_Applications_web_s.pdf)

[38] <https://auditingalgorithms.net/auditing-ml.pdf>

[39] <https://www.aitest.ai/aitest-guide>

[40] [https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper\\_Trusted-AI\\_TUeV-AUSTRIA\\_JKU.pdf](https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper_Trusted-AI_TUeV-AUSTRIA_JKU.pdf)

[41] <https://www.brz.gv.at/was-wir-tun/Innovationen/Wie-Zukunftstechnologien-die-Verwaltung-modernisieren/Kuenstliche-Intelligenz-in-der-Verwaltung/Vertrauenswuerdige-KI.html>

[42] <https://forhumanity.center/independent-audit-of-ai-systems/>

Title	Type	Publisher	Summary	Audience	AI aspect
Guideline for Designing Trustworthy Artificial Intelligence	Audit catalogue	Fraunhofer IAIS	Complete audit catalogue for AI applications. Covers all dimensions. The audit is based on checking documentation requirements, including documentation of the results of tests that the catalogue specifies. Audit is based on a risk-based method, the goal is to check whether the remaining risk is acceptable. In German, English version announced.	Developer, Evaluator	All
Auditing Machine Learning Algorithms – A white paper for Public Auditors	Audit catalogue	Supreme Audit Institutions of Finland, Germany, the Netherlands, Norway and the UK	Despite its name, the document actually is closer to an audit catalogue than a white paper. Essentially, it is a set of guidelines.	Developer, Evaluator	All
Algorithmic Impact Assessment Tool	Questionnaire	Canadian Government	Online assessment tool for evaluating the risk of automated decision-making tools. It consists of 48 risk and 33 mitigation questions.	Developer, User, Evaluator	All

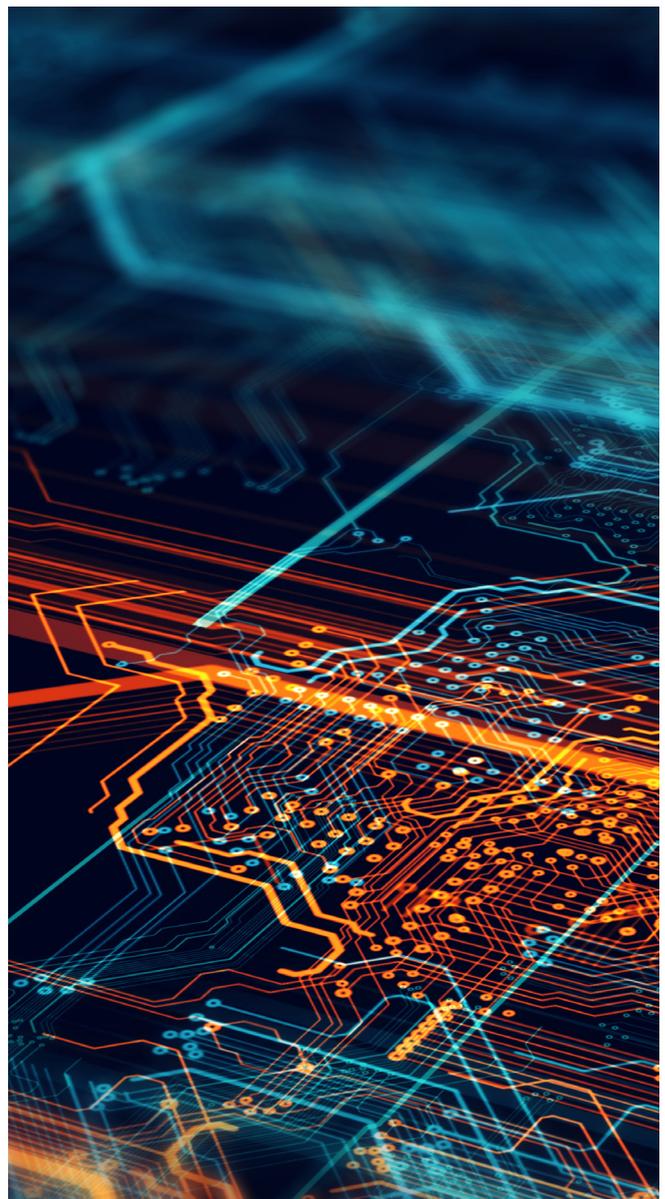
Table 4: Audit catalogues

## Conclusions

Trust in AI in general, and how to test or certify it, is an emerging topic of high relevance for a wide range of sectors and actors. In this report we provided an overview of relevant standards – both published standards and standards under development, legal regulations, white papers and reports, and audit catalogues. Most of the relevant standardization activities are still work in progress, as is research and legislation.

Therefore, new work is published regularly, and this overview can therefore only be a snapshot. Nevertheless, this overview will hopefully provide a starting point for all parties that are interested in trust in AI, be it from a developer, a user or an evaluator perspective.

While it is true that the relevant evaluation methods and schemes are still at their infancy, SGS as the world's largest testing, inspection and certification (TIC) company, will certainly be at the forefront providing professional services related to AI trustworthiness. We are constantly engaging with a large number of stakeholders to assess market needs and working with partners like Know-Center, a leading European innovation and research center for trustworthy AI and data science, to provide suitable solutions.



# About

## 4.1 Authors

### Tomislav Nad

Lead Innovation Technologist at SGS. He is exploring the potential and impact of emerging technologies and focusing on the security and trustworthiness of these technologies, amongst other aspects. He has worked in the field of cybersecurity for over 15 years as a researcher, engineer, consultant, evaluator, manager, speaker and lecturer, and helped to build various successful products and services.

Recently, he started to focus on the topic of trustworthiness of AI systems and since then delved deep into it. Tomislav holds a master's degree in Mathematics and a PhD degree in IT-Security and Cryptography from Graz University of Technology.

### Sebastian Scher

A physicist and holds a PhD degree in Atmospheric Sciences. He has worked on applying AI techniques to a wide range of research problems. Since 2022, he has been a senior researcher at Know-Center, a leading European research center for trustworthy AI, with strong ties both to the Technical University of Graz and to industry. Currently, Sebastian conducts research on how to make AI fair and robust.

### Florian Königstorfer

A PhD candidate at the Business Analytics and Data Science Center (BANDAS Center) at the Karl-Franzen University Graz. He is conducting research on suitable tools and methods to document AI applications. He holds a bachelor's degree in Econometrics and Operations Research and a master's degree in Business Intelligence and Smart Services from the University of Maastricht.

## 4.2 Organizations

### SGS Group of companies

The world's largest independent, international inspection, testing and certification organization. The Group comprises of more than 300 affiliated companies operating in over 140 countries with 99,600 employees in more than 2,600 offices and laboratories. Founded in 1878, today the SGS Group is unique in the international technical services sector, not only through its geographical coverage, but also through the comprehensive range of services offered. The SGS Group is a fully independent inspection and testing organization with no manufacturing, trading or financial interests which could compromise its independence. This, together with a reputation for quality, integrity and impartiality, is the basis of the Group's corporate development.

### The Know-Center

A leading European innovation and research center for trustworthy AI and data science. It makes data speak, recognizes patterns where no one looks anymore and creates unique solutions at the limits of what is conceivable today. For its customers in automation and logistics, energy and environmental management as well as medicine, these are decisive competitive advantages. For science and society, it is the database for solving the big questions of the future.

### The BANDAS Center

The digitization competence center of the University of Graz, School of Business, Economics and Social Sciences. The BANDAS Center deals with data-driven technologies in research and business. The focus of the center is how to apply data-driven technologies in the business world and how to evaluate the societal impact and consequences. One core of the BANDAS Center is trustworthy AI, with a particular focus on auditing and certification.





# Appendix

All the document titles and document abstracts/summaries included here are copied or directly derived from publicly available materials, however copyright of those original materials is retained by the respective owners.

## 5.1 Table of published standards and reports

Title	Type	Summary	Audience	AI aspect
<b>ISO/IEC</b>				
ISO/IEC 5338 Artificial intelligence – AI system life cycle processes	Standard	This document provides processes that support the definition, control, management, execution and improvement of the AI system in its life cycle stages. These processes can also be used within an organization or a project when developing or acquiring AI systems.	General, Developer, Evaluator	All
ISO/IEC 5339 Artificial Intelligence – Guidelines for AI applications	Technical report	Provides a macro-level view of an AI application to facilitate its understanding, development and use among all stakeholders. This includes an approach to identifying an AI application's stakeholders, context, functional characteristics and non-functional characteristics and guidelines for AI applications based on the make, use and impact perspectives.	Developer, User	N/A
ISO/IEC 5392 Information technology – Artificial intelligence – Reference architecture of knowledge engineering	Standard	This document defines a reference architecture of knowledge engineering (KE) in AI. The reference architecture describes KE roles, activities, constructional layers, components and their relationships amongst themselves and other systems from systemic user and functional views. This document also provides a common KE vocabulary by defining KE terms.	Developer	N/A
ISO/IEC 5469 Artificial intelligence – Functional safety and AI systems	Technical report	This document describes the properties, related risk factors, available methods and processes relating to: Use of AI inside a safety related function to realize the functionality; use of non-AI safety related functions to ensure safety for an AI controlled equipment; use of AI systems to design and develop safety related functions.	General, Developer, Evaluator, User	Safety
ISO/IEC 8183 Artificial intelligence – Data life cycle framework	Standard	This document provides an overarching data life cycle framework that is instantiable for any AI system from data ideation to decommissioning. This document is applicable to the data processing throughout the AI system life cycle including the acquisition, creation, development, deployment, maintenance and decommissioning.	General, Developer, Evaluator	Data governance
ISO/IEC 8200 Artificial intelligence – Controllability of automated artificial intelligence systems	Technical specification	This document defines a basic framework with principles, characteristics and approaches for the realization and enhancement for automated AI systems' controllability.	General, Developer, Evaluator, User	Human agency and oversight
ISO/IEC 17903 Artificial intelligence — Overview of machine learning computing devices	Technical report	An ML computing device can have a set of characteristics, including supported datatypes, ML operators, buffer settings, access and share mechanisms, memory access and addressing mechanisms, virtualization and sharing mechanisms, job scheduling mechanisms, topologies, data exchange mechanisms and memory interoperability mechanisms. This document surveys and provides information to AI stakeholders to assist them in understanding the representative characteristics of ML computing devices.	Developer, Evaluator	N/A

Title	Type	Summary	Audience	AI aspect
ISO/IEC 24027 Artificial intelligence (AI) – Bias in AI systems and AI aided decision making	Technical report	This document addresses bias in relation to AI systems, especially with regards to AI-aided decision-making.  Measurement techniques and methods for assessing bias are described, with the aim to address and treat bias-related vulnerabilities. All AI system lifecycle phases are in scope, including but not limited to data collection, training, continual learning, design, testing, evaluation and use.	Developer, Evaluator	Fairness
ISO/IEC 24028 Artificial intelligence – Overview of trustworthiness in artificial intelligence	Technical report	The goal of this document is to analyze the factors that can impact the trustworthiness of systems providing or using AI. The document briefly surveys the existing approaches that can support or improve trustworthiness in technical systems and discusses their potential application to AI systems. The document discusses possible approaches to mitigating AI system vulnerabilities that relate to trustworthiness. The document also discusses approaches to improving the trustworthiness of AI systems.	General, Developer, Evaluator, User	All
ISO/IEC TR 24029 Artificial Intelligence (AI) – Assessment of the robustness of neural networks – Part 1: Overview	Technical report	This document aims at providing an overview of the approaches available to assess risks related to robustness, with a particular focus on neural networks. Methods are categorized into three groups: statistical methods, formal methods and empirical methods. This document provides background on these methods to assess the robustness of neural networks.	Developer, Evaluator	Robustness and safety
ISO/IEC 24029 Artificial intelligence (AI) – Assessment of the robustness of neural networks – Part 2: Methodology for the use of formal methods	Standard	This document provides methodology for the use of formal methods to assess robustness properties of neural networks. The document focuses on how to select, apply and manage formal methods to prove robustness properties.	Developer, Evaluator	Robustness and safety
ISO/IEC TR 24030 Artificial intelligence (AI) – Use cases	Technical report	This document illustrates the applicability of the AI standardization work across a variety of application domains. It provides input to and reference for AI standardization work, etc. By investigating use cases, it is possible to find the new technical requirements (standardized demand) from the market, accelerating the transformation of science and technology achievements.	General, Developer, Evaluator, User	N/A
ISO/IEC TR 24372 Artificial intelligence (AI) – Overview of computational approaches for AI systems	Technical report	This document provides an overview of the state of the art of computational approaches for AI systems, by describing a) main computational characteristics of AI systems; b) main algorithms and approaches used in AI systems, referencing use cases contained in ISO/IEC TR 24030.	General, Developer, Evaluator	N/A
ISO/IEC 38507 Governance implications of the use of artificial intelligence by organizations	Technical report	The objective of this document is to provide guidance for the governing body of an organization that is using, or is considering the use of, AI. This document provides guidance on the role of a governing body with regard to the use of AI within their organization and encourages organizations to use appropriate standards to underpin their governance of the use of AI.	General, Developer, Evaluator, User	All
ISO/IEC TR 24372 Artificial intelligence (AI) – Overview of computational approaches for AI systems	Technical report	This document provides an overview of the state of the art of computational approaches for AI systems, by describing a) main computational characteristics of AI systems; b) main algorithms and approaches used in AI systems, referencing use cases contained in ISO/IEC TR 24030.	General, Developer, Evaluator	N/A
ISO/IEC 25059 Software engineering – Systems and software Quality Requirements and Evaluation (SQuaRE) – Quality model for AI systems	Standard	This document outlines a quality model for AI systems and is an application-specific extension to the SQuaRE series. The characteristics and sub-characteristics detailed in the model provide consistent terminology for specifying, measuring and evaluating AI system quality.	Developer, Evaluator	Robustness and safety

Title	Type	Summary	Audience	AI aspect
ISO/IEC 29119 11 Software and systems engineering – Software testing – Part 11: Testing of AI systems	Technical specification	This document describes testing techniques (including those described in ISO/IEC/IEEE 29119-4) applicable for AI systems in the context of the AI system life cycle model stages defined in ISO/IEC 22989. It describes how AI and ML assessment metrics can be used in the context of those testing techniques. It also maps testing processes, including those described in ISO/IEC/IEEE 29119-2, to the verification and validation stages in the AI system life cycle.	Developer, Evaluator	Robustness and safety
ISO/IEC 38507 Governance implications of the use of artificial intelligence by organizations	Technical report	The objective of this document is to provide guidance for the governing body of an organization that is using, or is considering the use of, AI. This document provides guidance on the role of a governing body with regard to the use of AI within their organization and encourages organizations to use appropriate standards to underpin their governance of the use of AI.	General, Developer, Evaluator, User	All
ISO/ TR 22100-5 Safety of machinery – Relationship with ISO 12100 – Part 5: Implications of embedded Artificial Intelligence machine learning	Technical report	This document addresses how artificial intelligence machine learning can impact the safety of machinery and machinery systems. This document describes how hazards being associated with AI applications machine learning in machinery or machinery systems, and designed to act within specific limits, can be considered in the risk assessment process.	Developer, Evaluator, User	Robustness and safety
ISO/IEC 4213 Artificial Intelligence – Assessment of machine learning classification performance	Technical specifications	This document specifies methodologies for measuring classification performance of machine learning models, systems and algorithms.	Developer, Evaluator	Robustness and safety
ISO/IEC 20546 Big data – Overview and vocabulary	Standard	This document provides a conceptual overview of the field of big data, its relationship to other technical areas and standards efforts, and the concepts ascribed to big data that are not new to big data.	General, Developer, Evaluator	N/A
ISO/IEC TR 20547-1 Big data reference architecture – Part 1: Framework and application process	Technical report	This document describes the framework of the big data reference architecture and the process for how a user of the document can apply it to their problem domain.	Developer, Evaluator	N/A
ISO/IEC TR 20547-2 Big data reference architecture – Part 2: Use cases and derived requirements	Technical report	This document provides examples of big data use cases with application domains and technical considerations derived from the contributed use cases.	Developer, Evaluator	N/A
ISO/IEC 20547-3 Big data reference architecture – Part 3: Reference architecture	Standard	This document specifies the big data reference architecture (BDRA). The reference architecture includes concepts and architectural views.	Developer, Evaluator	N/A
ISO/IEC 20547-4 Big data reference architecture – Part 4: Security and privacy	Standard	This document specifies the security and privacy aspects applicable to the big data reference architecture (BDRA) including the big data roles, activities and functional components and also provides guidance on security and privacy operations for big data.	Developer, Evaluator	Data governance, Privacy
ISO/IEC TR 20547-5 Big data reference architecture – Part 5: Standards roadmap	Technical report	This document describes big data relevant standards, both in existence and under development, along with priorities for future big data standards development based on gap analysis.	Developer, Evaluator	N/A
ISO/IEC 22989 Artificial intelligence – Artificial intelligence concepts and terminology	Standard	This document establishes terminology for AI and describes concepts in the field of AI. This document can be used in the development of other standards and in support of communications among diverse, interested parties or stakeholders.	General, Developer, Evaluator, User	All

Title	Type	Summary	Audience	AI aspect
ISO/IEC 23053 Framework for Artificial Intelligence (AI) Systems Using Machine Learning (ML)	Standard	This document establishes an AI and machine learning (ML) framework for describing a generic AI system using ML technology. The framework describes the system components and their functions in the AI ecosystem.	General, Developer, Evaluator, User	All
ISO/IEC 23894 Artificial intelligence – Guidance on risk management	Standard	This document provides guidance on how organizations that develop, produce, deploy or use products, systems and services that utilize AI can manage risk specifically related to AI. The guidance also aims to assist organizations to integrate risk management into their AI-related activities and functions. It moreover describes processes for the effective implementation and integration of AI risk management.	General, Developer, Evaluator, User	All
ISO/IEC TR 24368 Artificial intelligence – Overview of ethical and societal concerns	Technical report	This document provides a high-level overview of AI ethical and societal concerns. It includes an overview of international standards that address issues arising from AI ethical and societal concerns.	General, Developer, Evaluator, User	Fairness, societal and environmental wellbeing
ISO/IEC 24668 Artificial intelligence – Process management framework for big data analytics	Standard	This document provides a framework for developing processes to effectively leverage big data analytics across the organization irrespective of the industries or sectors.	Developer	Privacy and data governance
ISO/IEC 25058 Systems and software Quality Requirements and Evaluation (SQuaRE) — Guidance for quality evaluation of artificial intelligence (AI) systems	Standard	This document provides guidance for evaluation of artificial intelligence (AI) systems using an AI system quality model.	Developer, Evaluator	All
ISO/IEC 42001 Artificial intelligence – Management system	Standard	This document specifies the requirements and provides guidance for establishing, implementing, maintaining and continually improving an AI management system within the context of an organization. This document is intended for use by an organization providing or using products or services that utilize AI systems. This document helps the organization develop or use AI systems responsibly in pursuing its objectives and meet applicable regulatory requirements, obligations related to interested parties and expectations from them.	General, Developer, Evaluator, User	All
<b>DIN</b>				
DIN SPEC 92001-1 Artificial Intelligence – Life Cycle Processes and Quality Requirements – Part 1: Quality Meta Model	Standard	The purpose of this document is to establish a quality assuring and transparent life cycle of AI modules. This document presents a set of quality requirements that are structured in an AI specific quality metamodel. It proposes the differentiation between AI modules with regard to their safety, security, privacy and ethical relevance.	General, Developer, Evaluator	All
DIN SPEC 92001-2 Artificial Intelligence – Life Cycle Processes and Quality Requirements – Part 2: Robustness	Standard	This document provides specific requirements that ensure AI quality with respect to robustness are provided. Specifically, the AI quality pillar robustness is explained further and specific requirements for this pillar are listed. Furthermore, each requirement within the pillars is mapped to a set of lifecycle stages to facilitate the temporal classification of the requirements.	Developer, Evaluator	Robustness and safety
IEEE 2801-2022 Recommended Practice for the Quality Management of Datasets for Medical Artificial Intelligence	Standard	The document highlights quality objectives for organizations responsible for datasets. The document describes control of records during the lifecycle of datasets, including but not limited to data collection, annotation, transfer, utilization, storage, maintenance, updates, retirement and other activities.	Developer	Privacy and data governance

Title	Type	Summary	Audience	AI aspect
IEEE 2802-2022 Approved Draft Standard for the Performance and Safety Evaluation of Artificial Intelligence Based Medical Device: Terminology	Standard	This standard is aimed at establishing concepts and terminology for the performance and safety evaluation of artificial intelligence medical devices, which covers basic technology, dataset, quality characteristics, quality evaluation and application scenario.	Developer	Robustness and safety
IEEE 7000-2021 Standard Model Process for Addressing Ethical Concerns during System Design	Standard	The standard establishes a set of processes by which engineers and technologists can include consideration of ethical values throughout the stages of concept exploration and development, which encompass system initiation, analysis and design.	Developer	Fairness, societal and environmental wellbeing
IEEE 7001-2021 Standard for Transparency of Autonomous Systems	Standard	This standard provides a framework to help developers of autonomous systems both review and, if needed, design features into those systems to make them more transparent. The framework sets out requirements for those features, the transparency they bring to a system and how they would be demonstrated in order to determine conformance with this standard.	Developer, Evaluator	Transparency
IEEE 7002-2022 Standard for Data Privacy Process	Standard	The requirements for a systems/software engineering process for privacy-oriented considerations regarding products, services and systems utilizing employee, customer or other external user's personal data are defined by this standard. Specific procedures, diagrams and checklists are provided to perform conformity assessments on their specific privacy practices.	General, Developer, Evaluator	Privacy and data governance
IEEE 7007-2021 Ontological Standard for Ethically Driven Robotics and Automation Systems	Standard	This standard establishes a set of ontologies with different abstraction levels that contain concepts, definitions, axioms and use cases that are deemed relevant and appropriate to establish ethically driven methodologies for the design of robots and automation (R&A) systems.	Developer	Societal and environmental wellbeing
IEEE 7009 Standard for Fail-Safe Design of Autonomous and Semi-Autonomous Systems	Standard	This standard establishes a practical, technical baseline of specific methodologies and tools for the development, implementation and use of effective fail-safe mechanisms in autonomous and semi-autonomous systems.	Developer, Evaluator	Safety
IEEE 7010-2020 Recommended Practice for Assessing the Impact of Autonomous and Intelligent Systems on Human Well-Being	Standard	The impact of AI or autonomous and intelligent systems (A/IS) on humans is measured by this standard. The positive outcome of A/IS on human wellbeing is the overall intent of this standard. Scientifically valid wellbeing indices currently in use and based on a stakeholder engagement process ground this standard. Product development guidance, identification of areas for improvement, risk management, performance assessment and the identification of intended and unintended users, uses and impacts on human wellbeing of A/IS are the intents of this standard.	Developer	Societal and environmental wellbeing
IEEE 7014 Standard for Ethical Considerations in Emulated Empathy in Autonomous and Intelligent Systems	Standard	This standard defines a model for ethical considerations and practices in the design, creation and use of empathic technology, incorporating systems that have the capacity to identify, quantify, respond to or simulate affective states, such as emotions and cognitive states.	Developer	Societal and environmental wellbeing
IEEE 7015 Standard for Data and Artificial Intelligence (AI) Literacy, Skills, and Readiness	Standard	To coordinate global data and AI literacy building efforts, this standard establishes an operational framework and associated capabilities for designing policy interventions, tracking their progress and empirically evaluating their outcomes.	General	N/A

Title	Type	Summary	Audience	AI aspect
IEEE 2976 Standard for XAI – eXplainable Artificial Intelligence - for Achieving Clarity and Interoperability of AI Systems Design	Standard	This standard defines mandatory and optional requirements and constraints that need to be satisfied for an AI method, algorithm, application or system to be recognized as explainable.	Developer, Evaluator	Transparency
<b>NIST</b>				
AI Risk Management Framework (AI RMF)	Standard	The goal of the AI RMF is to offer a resource to the organizations designing, developing, deploying or using AI systems to help manage the many risks of AI and promote trustworthy and responsible development and use of AI systems.	General, Developer, Evaluator, User	All

## 5.2 Table of standards and reports in development

Note that, not all documents under development drafts are available and hence a summary cannot be provided. Several documents are in an early stage.

For ease of reading, we use the final name of the documents as of February 27, 2023. The state and content of the documents might change during development according to the standardization bodies processes.

Title	Type	Summary	Audience	AI aspect
<b>ISO/IEC</b>				
ISO/IEC 5259-1 Artificial intelligence – Data quality for analytics and machine learning (ML) Part 1: Overview, terminology, and examples	Standard	This document provides the means for understanding and associating the individual document of the standard series and is the foundation for conceptual understanding of data quality for analytics and machine learning.	Developer	Robustness and safety, data governance
ISO/IEC 5259-2 Artificial intelligence – Data quality for analytics and machine learning (ML) – Part 2: Data quality measures	Standard	Not available.	Developer	Robustness and safety, data governance
ISO/IEC 5259-3 Artificial intelligence – Data quality for analytics and machine learning (ML) – Part 3: Data quality management requirements and guidelines	Standard	This document specifies requirements and provides guidance for establishing, implementing, maintaining and continually improving the quality for data used in the areas of analytics and machine learning.	Developer	Robustness and safety, data governance
ISO/IEC 5259-4 Artificial intelligence – Data quality for analytics and machine learning (ML) – Part 4: Data quality process framework	Standard	This document provides general common organizational approaches, regardless of type, size, or nature of the applying organization, to ensure data quality for training and evaluation in analytics and machine learning.	Developer	Robustness and safety, data governance
ISO/IEC 5259-5 Artificial intelligence – Data quality for analytics and machine learning (ML) – Part 5: Data quality governance	Standard	This document provides a data quality governance framework for analytics and machine learning to enable governing bodies of organizations to direct and oversee the implementation and operation of data quality measures, management and related processes with adequate controls throughout the data life cycle.	Developer	Robustness and safety, data governance

Title	Type	Summary	Audience	AI aspect
ISO/IEC 5259-6 Artificial intelligence — Data quality for analytics and machine learning (ML) Part 6: Visualization framework for data quality	Technical report	Not available.	Developer	Robustness and safety, data governance
ISO/IEC 6254 Artificial intelligence – Objectives and approaches for explainability of ML models and AI systems	Technical specification	This document describes approaches and methods that can be used to achieve explainability objectives of stakeholders with regards to ML models and AI systems' behaviors, outputs and results.	General, Developer, Evaluator, User	Transparency
ISO/IEC 12791 Artificial intelligence – Treatment of unwanted bias in classification and regression machine learning tasks	Technical specification	This document provides mitigation techniques that can be applied throughout the AI system life cycle in order to treat unwanted bias. This document describes how to address unwanted bias in AI systems that use machine learning to conduct classification and regression tasks.	Developer, Evaluator	Fairness
ISO/IEC 12792 Artificial intelligence – Transparency taxonomy of AI systems	Standard	This document defines a taxonomy of information elements to assist AI stakeholders with identifying and addressing the needs for transparency of AI systems. The document describes the semantics of the information elements and their relevance to the various objectives of different AI stakeholders.	General, Developer, Evaluator, User	Transparency
ISO/IEC 17847 Artificial intelligence – Verification and validation analysis of AI systems	Technical specification	This document describes approaches and provides guidance on processes for the verification and validation analysis of AI systems (comprising AI system components and the interaction of non-AI components with the AI system components) including formal methods, simulation and evaluation.	Developer, Evaluator	Robustness and safety
ISO/IEC 18988 Artificial intelligence — Application of AI technologies in health informatics	Technical report	Not available.	-	-
ISO/IEC 20226 Artificial intelligence — Environmental sustainability aspects of AI systems	Technical report	Not available.	General, Developer, Evaluator, User	Societal and environmental wellbeing
ISO/IEC 21221 Artificial intelligence – Beneficial AI systems	Technical specification	Not available.	-	-
ISO/IEC 22440 Artificial intelligence — Functional safety and AI systems	Technical specification	Not available.	General, Developer, Evaluator, User	Robustness and safety
ISO/IEC 22443 Artificial intelligence — Guidance on addressing societal concerns and ethical considerations	Technical specification	This document provides guidance on how an organization can identify and address societal concerns and ethical considerations during the life cycle of AI systems that can potentially harm individuals and society. The document expands existing AI system governance, management system and impact assessment standards.	General, Developer, Evaluator, User	Societal and environmental wellbeing

Title	Type	Summary	Audience	AI aspect
ISO/IEC 23281 Artificial intelligence — Overview of AI tasks and functionalities related to natural language processing	Technical report	This document describes the concept of AI tasks as applied to natural language. It proposes a landscaping of the AI tasks related to the analysis or generation of natural language, as well as other language related functionalities that are associated to those AI systems. It identifies existing and competing terminologies, co-existing variants of the same tasks and functionalities, and how specific tasks can be affected by language diversity in terms of their role or their challenges.	Developer, Evaluator	N/A
ISO/IEC 23282 Artificial Intelligence — Evaluation methods for accurate natural language processing systems	Standard	This document specifies the evaluation of natural language processing systems, in the sense of measuring the quality of a system's results to assess its functional suitability. It provides a definition of evaluation methods for those systems, together with guidance on how to select, implement and interpret those evaluation methods. This document covers quantitative metrics as well as other evaluation methods. It includes requirements on the implementation of the described metrics, and further requirements on the technical resources involved in the evaluation process.	Developer, Evaluator	Robustness
ISO/IEC 24029-3 Artificial intelligence (AI) — Assessment of the robustness of neural networks  Part 3: Methodology for the use of statistical methods	Standard	This document provides methodology for the use of statistical methods to assess robustness properties of neural networks. The document focuses on how to select, apply and manage statistical methods to assess robustness properties.	Developer, Evaluator	Robustness
ISO/IEC 24970 Artificial intelligence — AI system logging	Standard	This document describes common capabilities, requirements and a supporting information model for logging of events in AI systems. This document is designed to be used with a risk management system.	Developer, Evaluator	All
ISO/IEC 25029 Artificial intelligence — AI-enhanced nudging	Standard	This standard applies to nudging mechanisms enhanced by AI systems. This document provides definitions, concepts and guidelines to address AI-enhanced nudging mechanisms by organizations. This standard aims to support organizations to deal with AI-enhanced nudging mechanisms in alignment with existing AI standards. "AI-enhanced nudging mechanisms" are a sub category of digital nudges and which are enhanced by AI systems. It provides use-cases to illustrate AI-enhanced nudging mechanisms. It provides guidelines and requirements for designing responsible AI-enhanced nudging mechanisms. This includes horizontal processes and key indicators using specific vertical examples.	Developer, Evaluator	-
ISO/IEC 25058 Software and systems engineering – Systems and software Quality Requirements and Evaluation (SQuaRE) – Guidance for quality evaluation of AI systems	Technical specification	Not available.	-	-
ISO/IEC 42005 Artificial intelligence – AI system impact assessment	Standard	This document provides guidance for organizations performing AI system impact assessments for individuals and societies that can be affected by an AI system and its intended and foreseeable applications. It includes considerations for how and when to perform such assessments and at what stages of the AI system lifecycle, as well as guidance for AI system impact assessment documentation.	General, Developer, Evaluator User	All

Title	Type	Summary	Audience	AI aspect
ISO/IEC 42006 Artificial intelligence — Requirements for bodies providing audit and certification of artificial intelligence management systems	Standard	This document specifies additional requirements for ISO/IEC 17021-1 in order to enable accredited and or peer assessed certification bodies to reliably audit the management system for organizations that develop or use AI systems or both according to ISO/IEC 42001 and to make an evaluation and decision for certification. The application of this document enables the certification bodies to meet the specific technical features and the particular risks in dealing with AI systems according to ISO/IEC 42001.	Evaluator	-
ISO/IEC 42102 Artificial intelligence — Taxonomy of AI system methods and capabilities	Standard	This document provides guidance on the classification of AI systems by describing a taxonomy of methods and capabilities. The taxonomy enables AI stakeholders to describe and have a common understanding of an AI system. This document applies to all types of organizations involved in any of the lifecycle stages of AI systems as well as to any AI stakeholder roles.	General, Developer, Evaluator User	All
ISO/IEC 42103 Artificial intelligence — Overview of synthetic data in the context of AI systems	Technical report	This document provides an overview of synthetic data concepts, methods, uses and considerations in the context of AI systems.	General, Developer, Evaluator User	Robustness, privacy and data governance
ISO/IEC 42105 Artificial intelligence — Guidance for human oversight of AI systems	Standard	This document provides guidance on human control and monitoring of AI systems, which is referred to as human oversight. This document extends ISO/IEC TS 8200. This document is applicable to all types of organizations. This document is applicable throughout the AI system life cycle.	General, Developer, Evaluator User	Human oversight
ISO/IEC 42106 Artificial intelligence — Overview of differentiated benchmarking of AI system quality characteristics	Technical report	Not available.	Developer, Evaluator	-
ISO/IEC 27090 Cybersecurity – Artificial Intelligence – Guidance for addressing security threats and failures in artificial intelligence systems	Standard	This document provides guidance for organizations to address security threats and failures in AI systems. The guidance in this document aims to provide information to organizations to help them better understand the consequences of security threats to AI systems, throughout their lifecycle, and descriptions of how to detect and mitigate such threats.	General, User	Security
ISO/IEC 27563 Security and privacy in artificial intelligence use cases – Best practices	Technical report	Not available.	-	-
<b>IEEE</b>				
IEEE 7003 Algorithmic Bias Considerations	Standard	This document provides the means for understanding and associating the individual document of the standard series and is the foundation for conceptual understanding of data quality for analytics and machine learning. This standard describes specific methodologies to help users certify how they worked to address and eliminate issues of negative bias in the creation of their algorithms.	Developer, Evaluator, Users	Fairness
IEEE 7008 Standard for Ethically Driven Nudging for Robotic, Intelligent and Autonomous Systems	Standard	“Nudges” as exhibited by robotic, intelligent or autonomous systems are defined as overt or hidden suggestions or manipulations designed to influence the behavior or emotions of a user. This standard establishes a delineation of typical nudges. It contains concepts, functions and benefits necessary to establish and ensure ethically driven methodologies for the design of the robotic, intelligent and autonomous systems that incorporate them.	Developer, Evaluator,	Societal and environmental wellbeing

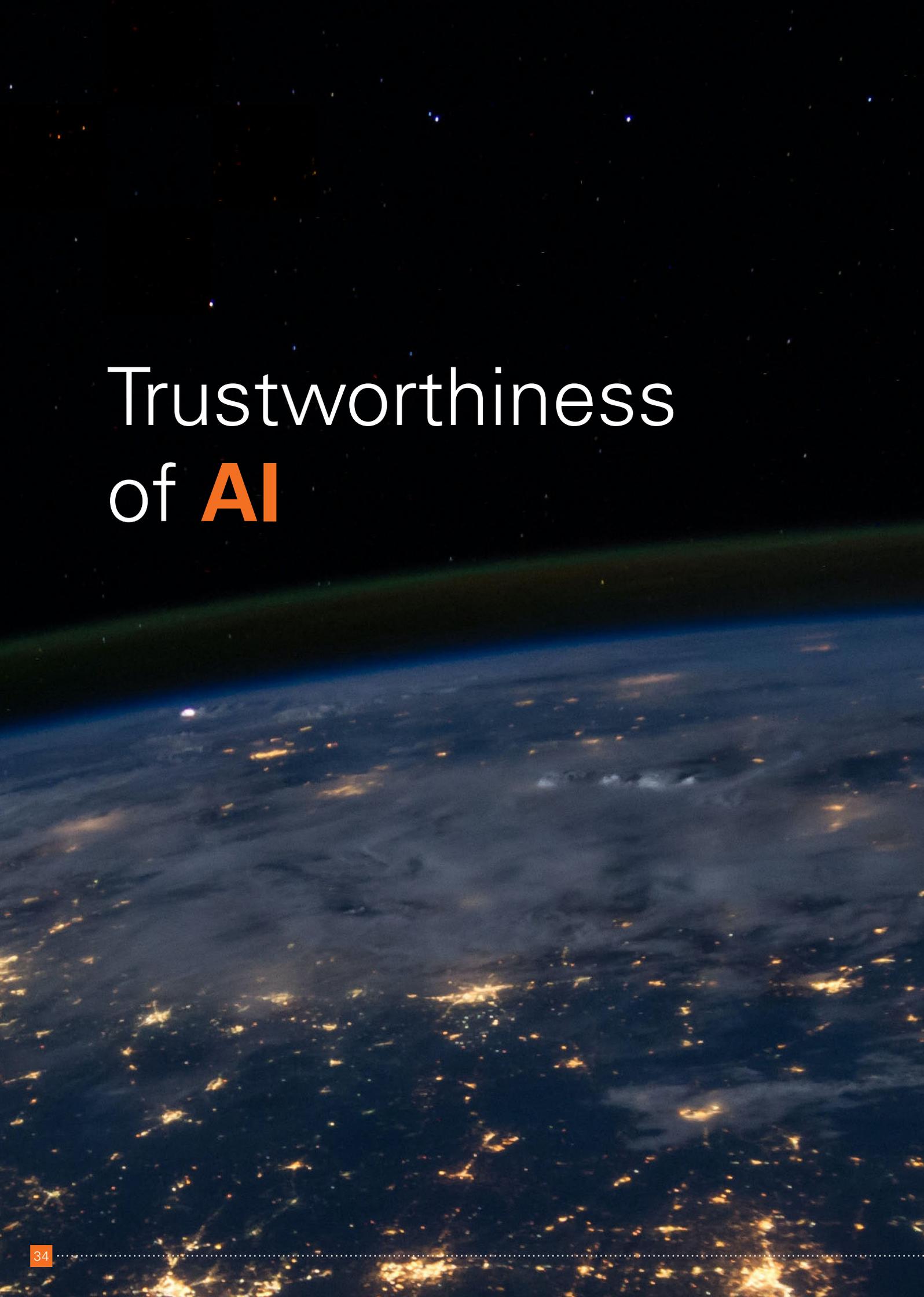
## 5.3 Table of white papers and reports

Title	Publisher	Summary	Audience	Context	AI aspect
Towards Auditable AI Systems <sup>43</sup>	TÜV-Verband, BSI, Fraunhofer HHI	<p>Overview of what can already be done in terms of certifying technical aspects of AI. Presents a “Certification Readiness Matrix” (CRM). The CRM has two dimensions: 1) stages in AI lifecycle, 2) technical aspects (e.g. security, safety, etc.). Each entry in the matrix gets a score for auditability scoring.</p> <p>The white paper does not give guidance on how to assign the scores, instead they are based on experience and intuition of the authors.</p> <p>The CRM is not really meant for a specific application, but as a general tool to show what is currently possible with regard to auditing AI applications. The results are based on a workshop held for that specific task.</p>	Developers, Auditors	Technical, Certification	Robustness and safety
Artificial Intelligence Cybersecurity Challenges <sup>44</sup>	ENISA	<p>Main purpose of the paper is “setting the ground for defining the AI threat landscape”, with focus on cybersecurity of AI. According to the paper, the relationship between AI and cybersecurity can be seen along three dimensions: cybersecurity for AI, AI to support cybersecurity, malicious use of AI. The paper provides a good overview of 1) AI lifecycle and 2) AI actors, and gives an overview of AI cybersecurity threat taxonomy.</p>	Developers	Technical	Robustness and safety, privacy and data governance
Trusted Artificial Intelligence Towards Certification of Machine Learning Applications <sup>45</sup>	TÜV Austria, Johannes Kepler Uni Linz, LIT AI Lab	<p>Overview on how to certify AI applications. The paper is detailed, but not detailed enough to be an actual audit catalogue.</p>	Developers, Evaluators	Certification	All

[43] [https://www.hhi.fraunhofer.de/fileadmin/News/2021/White\\_Paper/20210504\\_Whitepaper\\_\\_Towards\\_Auditable\\_AI\\_Systems\\_-\\_Current\\_status\\_and\\_future\\_directions\\_\\_final.pdf](https://www.hhi.fraunhofer.de/fileadmin/News/2021/White_Paper/20210504_Whitepaper__Towards_Auditable_AI_Systems_-_Current_status_and_future_directions__final.pdf)

[44] <https://www.enisa.europa.eu/publications/artificial-intelligence-cybersecurity-challenges/@@download/fullReport>

[45] [https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper\\_Trusted-AI\\_TUeV-AUSTRIA\\_JKU.pdf](https://www.tuv.at/wp-content/uploads/2022/03/Whitepaper_Trusted-AI_TUeV-AUSTRIA_JKU.pdf)

A photograph of Earth from space, showing the curvature of the planet and city lights at night. The background is a dark, starry sky. The text "Trustworthiness of AI" is overlaid on the image.

# Trustworthiness of **AI**

Title	Publisher	Summary	Audience	Context	AI aspect
Trust and Artificial Intelligence <sup>46</sup>	NIST	Study about why it is necessary that humans have trust in AI systems, not from a technical, but from a psychological viewpoint.	General, Developers	Certification, Standardization, General	All
Artificial Intelligence and future directions for ETSI <sup>47</sup>	ETSI	Overview over standardization activities, research activities, industry alliances and activities of SDOs regarding AI.	Developers, Evaluators	Standardization	Robustness and safety
Securing Artificial Intelligence (SAI); Mitigation Strategy Report <sup>48</sup>	ETSI	Detailed overview of AI-specific attacks and mitigation techniques against such attacks.	Developers	Standardization	Robustness and safety, privacy and data governance
Securing Artificial Intelligence (SAI); The role of hardware in security of AI <sup>49</sup>	ETSI	Paper on how hardware-assisted approaches can be used to ensure the integrity of possible untrusted ML platforms (such as cloud services). These hardware-assisted approaches include Trusted Execution Environments (TEEs), Roots of Trust (RoT) and specialized AI processing hardware.	Developers	Standardization	Robustness and safety, privacy and data governance
Securing Machine Learning Algorithms <sup>50</sup>	ENISA	The paper established a taxonomy of ML, gives an overview of threats for ML systems, existing security controls and recommendations on how the cybersecurity of systems that use ML can be enhanced.	General, Developers, Evaluators	Standardization	All
Standardization Roadmap AI <sup>51</sup>	DIN	Extensive overview of the needs and challenges surrounding standardization of AI, written from the perspective of Germany. Very extensive document.	General, Developers, Evaluators	Standardization, General	All
On Artificial Intelligence – A European approach to excellence and trust <sup>52</sup>	EC	High level white paper, presents policy options for how trustworthy and secure development of AI is possible in the EU, via aligning measures on European and national levels, and via a regulatory framework that builds an “ecosystem of trust”.	General	Policy making, Regulations	All
Policy and investment recommendations for trustworthy artificial intelligence <sup>53</sup>	HLEG/EC	Paper by the High-Level Expert Group on AI of the EC. High-level paper on trustworthy AI, featuring 33 concrete recommendations for policy and investments for European Institutions and member states, spanning private sector, public sector, research and academia, schools and regulation.	General	Policy making	All
Taxonomy of AI Risk <sup>54</sup>	NIST	Draft, overview of terms and definitions surrounding AI, and specifically AI risks.	Developers, Evaluators	Standardization	All

[46] <https://www.nist.gov/publications/trust-and-artificial-intelligence>

[52] <https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:52020DC0065>

[47] [https://www.etsi.org/images/files/ETSIWhitePapers/etsi\\_wp34\\_Artificial\\_Intelligence\\_and\\_future\\_directions\\_for\\_ETSI.pdf](https://www.etsi.org/images/files/ETSIWhitePapers/etsi_wp34_Artificial_Intelligence_and_future_directions_for_ETSI.pdf)

[53] [https://ec.europa.eu/newsroom/dae/document.cfm?doc\\_id=60343](https://ec.europa.eu/newsroom/dae/document.cfm?doc_id=60343)

[48] [https://www.etsi.org/deliver/etsi\\_gr/SAI/001\\_099/005/01.01.01\\_60/gr\\_SAI005v010101p.pdf](https://www.etsi.org/deliver/etsi_gr/SAI/001_099/005/01.01.01_60/gr_SAI005v010101p.pdf)

[54] [https://www.nist.gov/system/files/documents/2021/10/15/taxonomy\\_AI\\_risks.pdf](https://www.nist.gov/system/files/documents/2021/10/15/taxonomy_AI_risks.pdf)

[49] [https://www.etsi.org/deliver/etsi\\_gr/SAI/001\\_099/006/01.01.01\\_60/gr\\_SAI006v010101p.pdf](https://www.etsi.org/deliver/etsi_gr/SAI/001_099/006/01.01.01_60/gr_SAI006v010101p.pdf)

[50] <https://www.enisa.europa.eu/publications/securing-machine-learning-algorithms/@download/fullReport>

[51] <https://www.din.de/resource/blob/772610/e96c34dd6b12900ea75b460538805349/normungsroadmap-en-data.pdf>

**SGS.COM/DIGITAL**

## **CONTACT US**

### **SGS**

Emerging Technology

✉ [Enquiry.Emerging-Technology@sgs.com](mailto:Enquiry.Emerging-Technology@sgs.com)

### **KNOW CENTER**

Leading Research and Innovation Center for Trustworthy AI

🌐 <https://know-center.at/>

✉ [info@know-center.at](mailto:info@know-center.at)

### **BUSINESS ANALYTICS AND DATA SCIENCE-CENTER (BANDAS-CENTER)**

🌐 <https://business-analytics.uni-graz.at/de/>

✉ [bandas@uni-graz.at](mailto:bandas@uni-graz.at)

**SGS**

**When you need to be sure**