

AI Assessment in Practice: Implementing a Certification Scheme for AI Trustworthiness

Carmen Frischknecht-Gruber   Philipp Denzel  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland, Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Monika Reif  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Yann Billeter  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Stefan Brunner 

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Oliver Forster 

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Frank-Peter Schilling  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Joanna Weng  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Ricardo Chavarriaga  

Zurich University of Applied Sciences ZHAW, Winterthur, Switzerland

Abstract

The trustworthiness of artificial intelligence systems is crucial for their widespread adoption and for avoiding negative impacts on society and the environment. This paper focuses on implementing a comprehensive certification scheme developed through a collaborative academic-industry project. The scheme provides practical guidelines for assessing and certifying the trustworthiness of AI-based systems. The implementation of the scheme leverages aspects from Machine Learning Operations and the requirements management tool Jira to ensure continuous compliance and efficient lifecycle management. The integration of various high-level frameworks, scientific methods, and metrics supports the systematic evaluation of key aspects of trustworthiness, such as reliability, transparency, safety and security, and human oversight. These methods and metrics were tested and assessed on real-world use cases to dependably verify means of compliance with regulatory requirements and evaluate criteria and detailed objectives for each of these key aspects. Thus, this certification framework bridges the gap between ethical guidelines and practical application, ensuring the safe and effective deployment of AI technologies.

2012 ACM Subject Classification Computing methodologies → Artificial intelligence; Social and professional topics → Computing / technology policy; Information systems → Information systems applications

Keywords and phrases AI Assessment, Certification Scheme, Artificial Intelligence, Trustworthiness of AI systems, AI Standards, AI Safety

Digital Object Identifier 10.4230/OASICS.SAIA.2024.1

Funding This work was co-financed by Innosuisse (101.650 IP-ICT). The contribution of R.C. was partially funded by the Wellcome Trust [Grant number: 226486/Z/22/Z].

Acknowledgements We would like to acknowledge the support and collaboration of CertX AG in the development of the certification scheme discussed in this paper.

1 Introduction

Global efforts are underway to implement frameworks for assessing and regulating artificial intelligence (AI) systems. The most imminent of these efforts is the EU Artificial Intelligence



© Carmen Frischknecht-Gruber et al.; licensed under Creative Commons License CC-BY 4.0

Symposium on Scaling AI Assessments (SAIA 2024).

Editor: Rebekka Görges, Elena Haedecke, Maximilian Poretschkin, and Anna Schmitz; Article No. 1; pp. 1:1–1:17

OpenAccess Series in Informatics



OASICS Schloss Dagstuhl – Leibniz-Zentrum für Informatik, Dagstuhl Publishing, Germany

36 Act [8]. The AI act gradually comes into force starting 1 August 2024, which means
37 organisations and certifiers are in dire need of building their capacity to prove and assess
38 compliance now. However, despite this and other forthcoming regulations around the
39 globe, there remains a significant lack of practical guidelines and methodologies for both
40 achieving and assessing the trustworthiness of AI-based systems (AIS). Although there
41 has been extensive development of ethical guidelines for AI, c.f., Jobin et al. [20], the
42 practical application of these principles remains vague. The lack of specificity in the
43 operationalisation of these guidelines presents a challenge to their effective implementation
44 across various AIS. The introduction and deployment of inadequately understood and
45 unreliable AI technologies can result in significant societal harm. These include the exclusion
46 or discrimination of minorities due to inherent biases and even physical injuries resulting from
47 erroneous decision-making by AIS, such as in human-robot interactions or misdiagnoses in the
48 healthcare sector. Furthermore, such technologies have the potential to exacerbate existing
49 educational disparities, lead to unfair legal outcomes, and increase inequality. There is also a
50 substantial risk of environmental damage, privacy breaches, and cybersecurity vulnerabilities.
51 It is, therefore, imperative to develop tools that allow for AIS to be thoroughly vetted
52 for responsibility and ethical considerations to mitigate these risks and protect societal
53 well-being.

54 To address this issue, the authors, in collaboration with a certification company, are
55 developing a certification scheme for AIS. This scheme is intended as a practical guide
56 and provides corresponding tools for developers and regulators to evaluate and certify the
57 trustworthiness of AIS throughout their lifecycle, including requirements, data acquisition,
58 model development, testing, deployment, and operation. It builds upon current standards and
59 guidelines of a number of bodies, including ISO/IEC, IEEE, EASA, as well as other guidance
60 documents [19, 18, 12, 32, 27], in addition to EU legislation. A total of 38 documents were
61 subjected to analysis, and the objectives and the various means of complying with them were
62 derived from these inputs.

63 This certification scheme effectively bridges the gap between regulatory requirements,
64 technical standards, and the specific scientific and technical methods needed to assess the
65 properties of machine learning (ML) models. Noteworthy, regulatory requirements and
66 technical standards do not provide clear instructions on which methods and metrics can be
67 used to assess the properties for trustworthy AIS. To fill this gap, we evaluated and identified
68 95 technical methods for assessing the transparency, explainability, reliability, robustness,
69 safety, and security of AI models. By doing so, the certification scheme complements existing
70 approaches in trustworthy AI certification by incorporating cutting-edge research from the
71 AI community on algorithmic techniques for determining and evaluating relevant model
72 properties. As a result, it provides a complete operational framework that links the regulatory
73 requirements to measurable objectives and methods to assess compliance with the EU AI
74 Act and supports regulations in other jurisdictions.

75 This paper outlines the implementation and application of the certification scheme, with
76 a particular focus on detailing the tools, workflows, and methodologies used to ensure both
77 comprehensive compliance and practical utility. Furthermore, it describes how these tools and
78 methodologies relate to objectives for means of compliance and demonstrates our approach
79 to assessing the given requirements.

80 Identifying, tracing, and documenting appropriate objectives, procedures, and technical
81 methods for assessing compliance requires adequate supporting tools. We address these needs
82 by implementing the certification scheme within the management platform Jira. This is
83 complemented by an automatised pipeline that implements algorithmic methods for assessing

84 the trustworthiness of AI models. This pipeline is implemented according to best practices
85 in AI engineering and Machine Learning Operations (MLOps) principles.

86 In the remainder of this paper, we give an overview of the current state of AI standard-
87 isation and regulatory efforts, highlighting key initiatives and guidelines. In Section 3, we
88 outline the certification scheme, detailing relevant regulatory requirements, criteria, and the
89 methodology for certification.

90 Then, we describe the implementation process, including tools and frameworks used to
91 verify compliance, and how these relate to the regulatory requirements (Section 4). Finally,
92 we summarise our findings and offer a discussion on the implications and future developments
93 in AI certification (Section 5).

94 **2** Background

95 The deployment and scalability of AI assessment frameworks face several key challenges,
96 particularly in balancing practical implementation with theoretical underpinnings. One of
97 the main obstacles lies in the aggregation of risks associated with AI systems, including bias,
98 transparency, security vulnerabilities and ethical considerations. Current frameworks often
99 address individual risks in isolation, but aggregating these risks in a way that provides a
100 holistic assessment is complex. Many frameworks still lack widely accepted methods for
101 this aggregation, leading to inconsistencies across industries and sectors. A significant need
102 for interdisciplinarity also poses a challenge in scaling AI assessment frameworks. Inputs
103 from law, ethics and computer science must be combined to form a coherent assessment
104 approach. Managing this complexity requires the integration of technical AI safety measures
105 with broader societal values, which is often challenging to operationalise at scale [37]. In
106 terms of approaches, the risk-based approach used in regulations such as the EU AI Act
107 offers a promising method for scaling up. This regulation categorises AI systems according
108 to the level of risk they pose, from low-risk applications such as spam filters to high-risk
109 systems such as healthcare AI. The EU AI law imposes strict regulatory requirements on
110 high-risk systems to ensure safety and accountability. Conversely, AI systems classified as
111 low risk are subject to a more flexible regulatory framework. Although these systems are not
112 subject to the same stringent requirements, they must still comply with transparency and
113 user information obligations. This risk-based classification ensures that regulatory oversight
114 is aligned with the potential impact of AI systems, thereby increasing overall regulatory
115 effectiveness while facilitating innovation in lower-risk areas. [9].

116 **2.1** Regulation and Standards

117 Currently, there are significant global efforts to establish regulatory frameworks for AI. The
118 EU has assumed a pioneering position with the AI Act, which is designed to establish a
119 comprehensive regulatory framework for AIS [8]. In a similar vein, the United States issued
120 an executive order in October 2023 with the objective of developing new standards for safe,
121 secure, and trustworthy AI [43].

122 Standards and guidelines play a pivotal role in supporting binding laws and regulations
123 by documenting best practices and providing a foundation for demonstrating compliance
124 and certification. A considerable number of national and international organisations are
125 engaged in a range of initiatives aimed at fostering trust in AI through the issuance of
126 standards and guidelines. The ISO/IEC standards [19] address a plethora of AI-related
127 aspects, including terminology, performance metrics, data quality, ethics, and human-AI
128 interaction. These standards are currently in place, with more anticipated in the future.

129 Similarly, the IEEE is developing a certification program with the objective of assessing the
130 transparency, accountability, bias, and privacy of AI-related processes [17]. The IEEE P7000
131 series [18] addresses the ethical implications of AI technologies. Other national entities, such
132 as the National Laboratory of Metrology and Testing's (LNE) AI certification program, have
133 established objective criteria for trustworthy AIS, emphasising ethics, safety, transparency,
134 and privacy [23]. The NIST framework [27] offers guidance on the management of risks, the
135 assurance of data quality, and the promotion of transparency and accountability in AIS, with
136 related principles also emphasised in the AI Risk Management Framework [28]. Moreover,
137 DIN/DKE offers comprehensive standardisation recommendations across all AI domains,
138 facilitating a unified language, principles for development and utilisation, and certification [11].
139 In the field of aviation, the European Union Aviation Safety Agency (EASA) has introduced
140 comprehensive guidelines for the safe utilisation of ML systems [42]. These guidelines provide
141 support to stakeholders in the aviation sector at each stage of the lifecycle of AIS, from
142 the initial stages of development through to operational use. The Fraunhofer Institute has
143 developed a guideline for the design of trustworthy AI systems [32]. The guideline employs a
144 six-dimensional evaluation framework to assess the trustworthiness of AIS, encompassing
145 fairness, autonomy and control, transparency, reliability, safety and security, and privacy. In
146 contrast to other contributions, the Fraunhofer guideline incorporates both process-related
147 measures and technical methods to enhance the evaluation of AIS.

148 **2.2 Frameworks**

149 Capturing, tracing, documenting, and systematically evaluating requirements throughout
150 the lifecycle of an AIS is an essential factor in trustworthy AI and its certification.

151 There are various methods and tools for the filtering and management of requirements,
152 from very basic text files or Excel sheets to dedicated frameworks such as Confluence, Jira,
153 Doorstop, Polarion, IBM Doors, Azure DevOps, and many more [3, 2, 5, 40, 33, 26]. In
154 practice, the simple solutions do not provide the necessary flexibility and overview of the
155 complicated relations between requirements. On the other hand, comprehensive requirement
156 management frameworks are flexible but often less intuitive in their use and relatively
157 expensive. After investigating several tools, we chose **Jira** (in its basic version, free) as a
158 requirement management tool for the certification of AIS. Jira is a project management and
159 issue-tracking software developed by Atlassian. It helps teams plan, track, and manage work
160 efficiently, offering features like customisable workflows, real-time reporting, and integration
161 with numerous other tools, making it a versatile solution for agile project management.

162 An important operational approach to scaling is the integration of Machine Learning
163 Operations (MLOps). The role of MLOps principles and best practices in AIS development
164 and operation, as well as its assessment, is twofold: First, Billeter et al.[4] and others [24] have
165 advanced the idea of MLOps as the enabler of trustworthy AI by design. This means that
166 following MLOps guidelines and principles during design, development and operation of an
167 AIS, will lead to increased trustworthiness of the AIS. These practices include version control,
168 continuous integration and deployment (CI/CD), automated testing, and monitoring. Second,
169 the assessment of the trustworthiness of AIS also requires comprehensive evaluations of many
170 objectives and means of compliance (MOC) derived from these requirements. Therefore,
171 concepts like following best practices in AI engineering and MLOps are indispensable not
172 just during AIS development but also during its assessment.

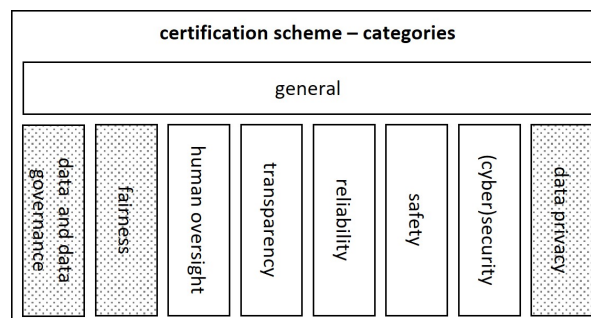
173 2.3 Algorithmic Tools for Trustworthy AI

174 While in some aspects of the verification of AI trustworthiness it is necessary to rely on
 175 qualitative results, in particular for model explainability or robustness, automated evaluation
 176 workflows mostly involve algorithmic methods with quantifiable output. Therefore, it is
 177 crucial to integrate assessment toolboxes which implement various algorithms and metrics,
 178 or rely on interfaces which allow for manual qualitative evaluation. There are a number
 179 of comprehensive toolboxes which implement appropriate technical methods paired with
 180 metrics, often isolated to assess specific aspects of AI trustworthiness such as transparency,
 181 reliability, or safety. For data and model explainability, industry-developed frameworks are
 182 Microsoft’s InterpretML [25], Seldon’s Alibi toolbox [38], IBM’s AIX360 toolbox [45], Sicara’s
 183 tf-explain [39], or PyTorch’s captum API [6]. Additionally, Quantus [16] is a relatively new
 184 and complementary explainability toolbox which implements a growing number of metrics
 185 and provides interfaces for other toolboxes such as captum or tf-explain. Toolboxes for
 186 testing the reliability, robustness, and safety of an AI model are, e.g., MIT’s Responsible
 187 AI Toolbox [41], Seldon’s Alibi-Detect [47], IBM’s ART [44] and UQ360 [46] toolboxes. In
 188 particular, there are many more toolboxes which implement specific tests for adversarial
 189 robustness, such as RobustnessGym [15], CleverHans [31], or Foolbox [34].

190 It is worth noticing that these toolboxes have been developed in parallel and, to a
 191 large extent, disconnected from the regulatory and certification frameworks. Hence, there
 192 suitability for compliance assessment is not entirely clear. Our analysis presents a significant
 193 step towards the integration of advances on both areas.

194 3 Overview of the Certification Scheme

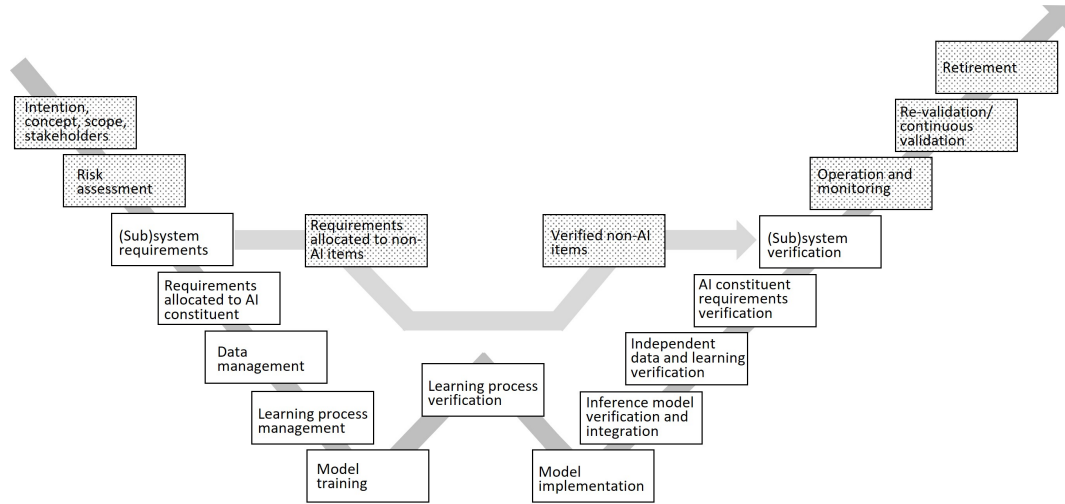
195 The developed certification scheme for AIS encompasses several principal key aspects of
 196 trustworthiness, such as human oversight, transparency, reliability (including robustness),
 197 safety and security [10] at the moment. So we cover with the actual version already some of
 198 the key aspects of the EU AI act, other aspects as described in 1. Each aspect is considered
 199 to ensure that AIS operate effectively, ethically, and safely across various applications.



■ **Figure 1** Extended key aspects of trustworthiness. The aspects of trustworthiness are as follows: data and data governance, fairness, human oversight, transparency, reliability, safety, (cyber)security, and data privacy. Within the certification scheme, the five non-shaded aspects are addressed, while the other three will be addressed at a later stage.

200 In addition, the certification scheme encompasses all relevant phases of the AIS lifecycle,
 201 as illustrated in Figure 2. The Certification Scheme employs a risk-based methodology in
 202 accordance with the EU AI Act. It commences with the concept of the system (including the

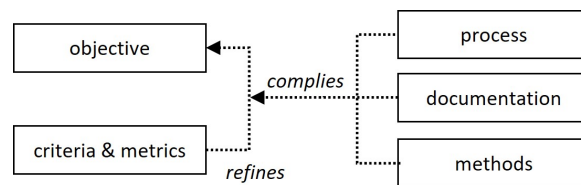
203 role of the AI part within the overall system) and the associated system risk, subsequently pro-
 204 gressing to the implementation of an AI model. The scheme culminates with the deployment,
 205 verification, validation, and operation of the system.



■ **Figure 2** Illustration of the lifecycles encompassed by the certification scheme, including risk assessment, (sub-)system requirements and design, data management, learning process management, model training, verification steps, and operation and monitoring.

206 **3.1 Key Aspects and Objectives**

207 For each phase of the lifecycle, the scheme identifies and addresses the critical key aspects
 208 through the establishment of corresponding objectives. These build the basis for proving
 209 compliance with regulations as the EU AI Act and are derived from the EU AI Act, existing
 210 ISO standards [19] and the EASA guidances [42]. The objectives are refined according to
 211 different qualitative criteria and quantitative metrics (see Figure 3). Different MOCs have
 212 been defined to achieve compliance with the aforementioned objectives. One group of MOCs
 213 describes the process means that must be in place for a thorough development, verification,
 214 or management process. Others describe the documentation means to cover, for example,
 215 auditability and other record-keeping aspects. The last group of MOCs define the technical
 216 methods that must be applied to achieve compliance with the objectives posed. These MOCs
 217 establish the link to the different technical methods of the second technical part of the
 218 certification scheme.



■ **Figure 3** The interrelationship between objectives, criteria and metrics, and compliance methods is illustrated in the diagram. The left side depicts the objectives and their refinement through the application of criteria and metrics, while the right side shows the processes, documentation and methods that ensure compliance with these objectives and criteria.

219 Initially, the certification scheme focused on transparency and reliability, encompassing 29
 220 and 44 objectives, respectively, with 100 and 156 MOCs. An updated version of the scheme
 221 additionally includes human oversight, safety and security alongside some general objectives
 222 relevant across multiple key aspects. Currently, the scheme covers:

- 223 ■ General Objectives: 5 objectives, 14 MOCs
- 224 ■ Human Oversight: 62 objectives, 65 MOCs
- 225 ■ Transparency: 29 objectives, 53 MOCs
- 226 ■ Reliability: 36 objectives, 105 MOCs
- 227 ■ Safety: 2 objectives, 6 MOCs
- 228 ■ (Cyber)Security: 5 objectives, 17 MOCs

229 The scheme includes a risk analyses and also addresses overlapping areas across key
 230 aspects, ensuring a comprehensive and integrated approach. Additional key aspects, such
 231 as data and data governance, will be implemented in the next step, and the key aspects of
 232 fairness and data privacy are planned for subsequent steps. In the following, we present two
 233 example objectives and their corresponding MOCs.

234 **Objective 1:** The applicant should define performance metrics to evaluate AIS performance
 235 and reliability.

- 236 ■ **MOC:** Define a suitable set of performance metrics for each high-level task to evaluate
 237 AIS performance and reliability.
- 238 ■ **MOC:** Define the expected performance with training, validation, and test data sets.
- 239 ■ **MOC:** Provide a comprehensive justification for the selection of metrics.

240 **Objective 2:** The applicant should identify and document the methods at AI/ML item
 241 and/or output level satisfying the specified AI explainability needs.

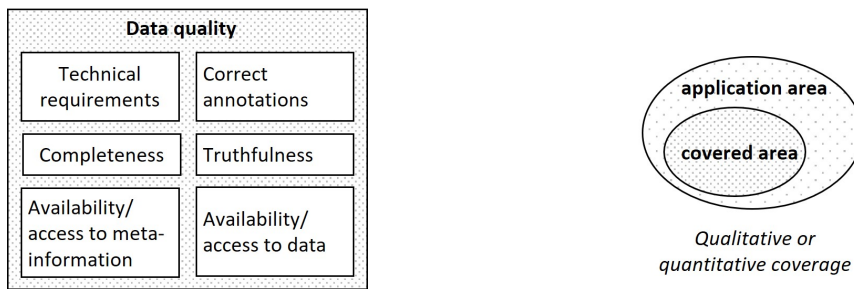
- 242 ■ **MOC:** Provide documentation of methods to provide explanations about the AI/ML
 243 item. The type and scope of the provided explanations should be chosen in terms of
 244 proportionality, considering the stakeholders.
- 245 ■ **MOC:** Specify the rules that apply to the current decision (e.g., for decision trees, list
 246 the selected branching next to the model output).
- 247 ■ **MOC:** Specify the most relevant attributes for a decision in linear regression models
 248 (e.g., for normalised inputs, the largest absolute coefficient value).
- 249 ■ **MOC:** For white-box models, use model-specific or model-agnostic methods for inter-
 250 pretability.

251 3.2 Key Aspects Overview

252 This section provides an overview of the key aspects covered by the scheme, including data
 253 governance, human oversight, transparency, reliability, and safety and (cyber)security.

254 3.2.1 Data and Data Governance

255 A dependable data set for a specific task requires careful attention to four key aspects:
 256 data quality, completeness, representativeness, and transparency. Data quality focuses on
 257 ensuring formal completeness and correctness and establishing reliability. The training,
 258 validation, and test data quality is assessed through qualitative and quantitative means
 259 (Figure 4). Correct annotations, task relevance, and data origin are crucial, alongside
 260 ensuring application coverage through metrics like class balance. Bias prevention requires
 261 unbiased training, validation, and test data, with fairness assessed via metrics like cosine
 262 similarity. Transparency ensures data is interpretable and preprocessing steps are clear,
 263 enabling verification by stakeholders.



(a) Data quality consists of six aspects. (b) Data coverage for the application.

■ **Figure 4** Data quality (formal data completeness and correctness) and data coverage.

264 **3.2.2 Human Oversight**

265 Human oversight of AIS, also referred to as autonomy and control, addresses potential risks
 266 that may arise when autonomous AI components limit the ability of users or experts to
 267 perceive or act. This aspect of AI safety ensures that system autonomy is appropriately
 268 constrained when it deviates from normal operation. To assess human oversight, AIS are
 269 categorised into four levels based on human involvement [29]. The first level, Human Control
 270 (HC), involves the AI acting solely as an assistive tool, where humans are responsible for every
 271 decision and subsequent action based on the AI’s output. At the Human-in-the-Loop (HIL)
 272 level, the AI operates partially autonomously but requires human intervention or confirmation,
 273 with humans monitoring and correcting its decisions as needed. The Human-on-the-Loop
 274 (HOL) level allows the AI to function almost autonomously, with limited human involvement
 275 for monitoring and occasional overrides. Finally, at the Human-out-of-the-Loop (HOOTL)
 276 level, the AI operates fully autonomously, handling tasks independently even in unexpected
 277 situations, with humans only involved in initial setup decisions like setting meta-commands
 278 in autonomous vehicles.

279 This key aspect includes objectives such as the implementation of human monitoring and
 280 control mechanisms, preservation of human decision-making capabilities, and ensuring the
 281 traceability of the AI component’s decision-making process.

282 **3.2.3 Transparency**

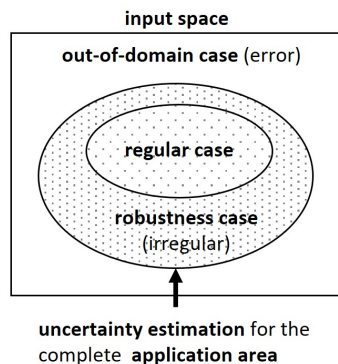
283 Transparency in AI is essential to prevent potential harm and ensure systems are under-
 284 standable to different stakeholders [36]. Transparency objectives are tailored to users, those
 285 affected (society), and experts (developers, providers, auditors and evaluators, authorities)
 286 (Figure 5). It involves setting criteria for interpretability and explainability, focusing on
 287 clarity, comprehensibility, and relevant metrics [7]. The interpretability of the ML model
 288 must be ensured through thorough documentation and visual aids like schematic diagrams.
 289 Explanation methods should be carefully chosen, justified, and documented, considering the
 290 audience’s qualifications. These methods should be evaluated statistically and by human
 291 reviewers, with a system in place for addressing user queries. For experts, transparency also
 292 involves validating decisions, ensuring technical traceability, and maintaining reproducibility.
 293 Key considerations include the scope, design, and stability of explanation methods relative
 294 to model outputs.

Society: Trust and understanding by clearly communicating the strengths and limitations.	Developers: Clear insights into the internal workings, and limitations.	Users: Transparent explanation of the decisions and results.	Authority: Assurance of regulatory compliance and operational transparency by thorough documentation.
	Providers: Monitoring capability by information on internal operations and performance	Auditors & Evaluators: Audit/evaluation capability	

■ **Figure 5** Transparency needs vary between stakeholders. The figure shows exemplary transparency requirements for some key stakeholders.

3.2.4 Reliability

Reliability in AIS is defined as the consistent execution of intended functions and also entails robustness, which pertains to maintaining performance under disturbances. An important concept is the Operational Design Domain (ODD), which delineates the specific conditions under which AIS can operate safely and effectively [35]. For developers, the ODD builds the basis for deriving detailed technical specifications that define the AIS input space, categorised into regular cases involving minor, expected disturbances; robustness cases where larger disturbances are encountered; and out-of-domain (OOD) cases, which involve data outside the application domain which may result in errors (Figure 6).



■ **Figure 6** Visualisation of the input space divided into regular, robustness, and out-of-domain cases.

Consequently, reliability is assessed in the three input spaces, in addition to the estimation of uncertainty. The regular case ensures reliable performance through data coverage, augmentation, and performance metrics evaluation (see Figure 7). Robustness tackles challenging conditions by addressing vulnerabilities and adversarial attacks. In out-of-domain (OOD) cases, the focus is on catching errors and improving generalisation, while uncertainty estimation involves setting appropriate metrics, assessing both intrinsic and extrinsic uncertainties, and developing mitigation measures.

Additional process steps include evaluating model architecture, implementing optimisation techniques such as pruning or quantisation, ensuring reproducibility, conducting regular assessments, and meticulously documenting all activities.

In the certification scheme, reliability assessment involves over 55 metrics and 95 methods, with a subset of 35 metrics and 50 methods selected for empirical testing. This selection was based on relevance, execution time, reliance on available information, and computational costs.

Performance metrics			
Regression	likelihood ratio	Confusion matrix	Completeness score
• (Mean) squared error	• True/False-negative rate	• Hinge loss	
• (Mean) absolute error	• True/False-positive rate		Ranking
Classification	• Precision-recall curve	Computer vision	• Mean reciprocal rank
• (Balanced) accuracy	• Receiver operation characteristics (ROC)	• Peak-signal-to-noise ratio (PSNR)	• Discounted cumulative gain
• Micro/macro average	• Lift	• Structural similarity	Natural Language Processing
• (Balanced) F1-score	• Matthew's correlation coefficient	• (Mean) intersection over union (mIOU)	• Perplexity score
• Prevalence	• Area under curve (AUC)	Clustering	• BLEU score
• Precision	• Cohen's Kappa	• Silhouette value	
• False discovery/omission rate		• Adjusted mutual information score	
• Positive/negative			

■ **Figure 7** List of performance metrics used for regression, classification, computer vision, clustering, ranking, and natural language processing.

317 Metrics vary across application domains and model objectives, so choosing the appropriate
 318 metric and method requires careful consideration of the model's goals, data characteristics,
 319 and desired outcomes. For example, formal verification employs logical and mathematical
 320 proofs to confirm system criteria, while model coverage analysis ensures comprehensive
 321 testing across various scenarios.

322 3.3 Safety and (Cyber)Security

323 The objective of safety is to minimise harm to people and the environment by designing
 324 AIS that incorporate corrective mechanisms for unexpected behaviours. This is of particular
 325 importance in contexts such as autonomous vehicles and healthcare, where errors can have
 326 significant and adverse consequences. (Cyber)security guarantees a system's integrity and
 327 availability by safeguarding it against unauthorised access, modification, or destruction. This
 328 encompasses the implementation of robust access controls, the assurance of data and model
 329 integrity, and the maintenance of system availability even in the event of an attack. Effective
 330 security measures are imperative for AIS in critical infrastructure, as breaches could result
 331 in significant damage. In order to enhance the security and resilience of AIS, adversarial
 332 training and verification are employed. Adversarial training is a method for enhancing
 333 the robustness of a model by exposing it to perturbations designed to deceive it, thereby
 334 identifying potential vulnerabilities.

335 4 Implementation of the Certification Scheme

336 The implementation and subsequent application of the AIS Certification Scheme to customers
 337 must meet established standards and regulations, requiring a carefully managed process. To
 338 achieve this, we evaluated several requirements management tools and ultimately selected
 339 Jira as the key tool for organising the objectives and associated means of compliance for our
 340 certification scheme. We then implemented an MLOps system based on state-of-the-art open-
 341 source tooling to perform the technical assessment of the AIS and evaluate the compliance
 342 with the defined objectives. In the following, we describe the requirement management
 343 system and the MLOps infrastructure.

344 4.1 Requirement Management Implementation

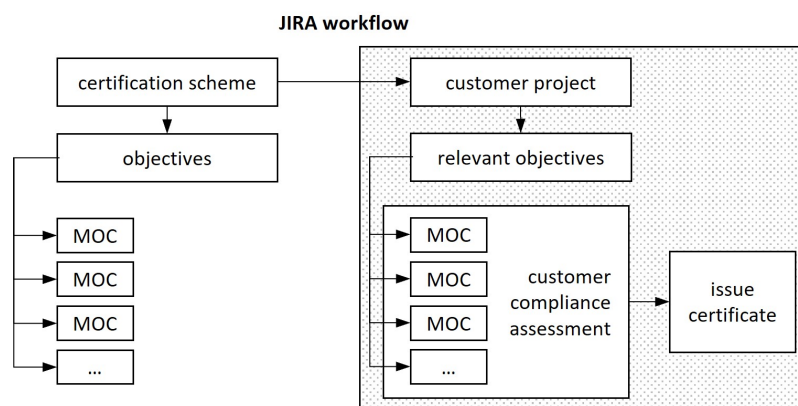
345 As written in section 2.2, Jira was chosen as requirement management tool to ensure traceability and effective management of the requirements. AI certification frameworks must adhere
346 to internationally recognised standards, including ISO 9001 (Quality Management Systems),
347 ISO/IEC 27001 (Information Security Management Systems), and the ISO/IEC 23894 (AI
348 - Guidance on Risk Management). Such standards necessitate meticulous documentation,
349 traceability, and periodic auditing to guarantee sustained compliance. The implementation of
350 such requirements in a manual or disparate system would increase the risk of inconsistencies
351 and errors, which would ultimately impact the efficiency and credibility of the certification
352 process. It is therefore imperative that robust requirements management tools are employed.
353

354 Using Jira for requirement management ensures each objective and MOC is meticulously
355 organised, facilitating clear communication and comprehensive oversight. Its ability to
356 maintain detailed records and provide real-time updates is crucial for this task. Real-time
357 collaboration and review capabilities are critical in aligning project tasks and reducing
358 errors. The platform supports multi-user editing, allowing teams to work simultaneously
359 from different locations. This live collaboration and features, such as decision tracking and
360 impact analysis, ensure that the development of the certification scheme remains agile and
361 responsive to changes. Additionally, the system's version control and history management
362 provide a complete audit trail, which is crucial for maintaining consistency and verifiability.

363 Centralised management of objectives and MOCs in a digital environment allows for
364 streamlined workflows and task alignment. We developed customised templates and dash-
365 boards for managing and tracing requirements. The possibility of sorting issues by attributes
366 such as the tag COMPLETE was proven to facilitate requirement tracking in the evaluation
367 we made of the platform. Each objective and MOC can be linked to others, showing rela-
368 tionships such as blocking issues and dependencies. The system's adaptability through the
369 reusability of issues across different projects and its capacity for baseline creation significantly
370 enhance the efficiency of the certification process. The platform facilitates organised and
371 efficient project management by enabling tasks such as editing, organising decision-making,
372 and managing tasks through a user-friendly interface. Integration with state-of-the-art tools,
373 such as Git integration platforms like GitHub or GitLab, as well as business communication
374 tools like Teams and Slack, along with the capability to create customisable pages, allows
375 the tool to be precisely tailored to specific project needs.

376 The certification scheme is structured with a parent-child relationship between objectives
377 and MOCs (Figure 8). Each issue type is defined by attributes, including description, main
378 category, additional categories, lifecycle phase, risk level, references, and approval status,
379 ensuring thorough documentation and easy information retrieval via specific filtering. This
380 structured approach facilitates organisational efficiency and enables the certification process
381 to be adapted as required.

382 For practical use, the certification scheme we developed has been implemented as a base
383 project; which can be readily exported, adapted, re-imported, or cloned to align with the
384 particular requirements of the customer or AI system to be assessed. For certification bodies
385 working with clients, the base project serves as the foundation from which the customer's
386 certification project is derived. The customer's AIS is then assessed against the MOCs from
387 the base scheme, supporting the issuance of the final certification.



■ **Figure 8** Visualisation of the certification workflow based on JIRA.

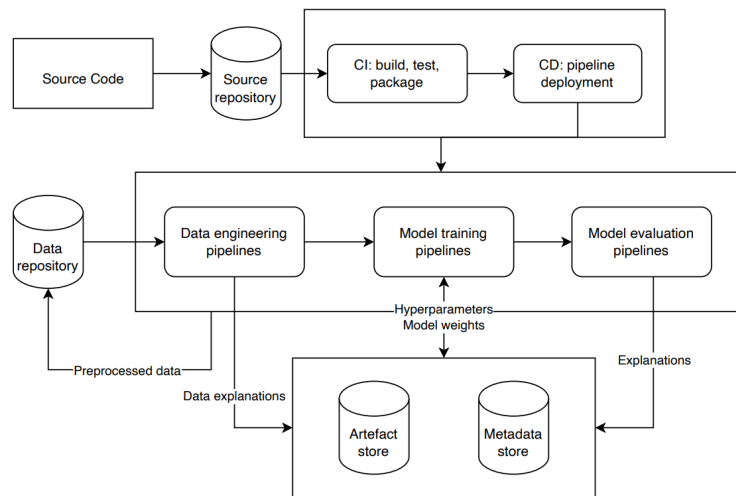
388 4.2 Machine Learning Operations Infrastructure

389 As argued in section 2.2. MLOps serves as enabler for trustworthy AI by design. It provides
 390 the necessary infrastructure and practices, ensuring that AIS are developed, deployed, and
 391 maintained reliably and efficiently. The adoption of MLOps thus facilitates the integration of
 392 trustworthy AI principles at every stage of the AI lifecycle, which are critical for regulatory
 393 compliance and societal acceptance [4]. MLOps extends DevOps practices to manage
 394 the complexities of bringing AIS into production, ensuring that they continuously meet
 395 trustworthiness standards [21].

396 The general architecture an AIS which adheres to MLOps principles supports the entire
 397 lifecycle of AIS and ensures that models are reproducible, reliable, and maintainable. This
 398 architecture includes project setup and requirements engineering, data engineering, model
 399 development, continuous integration/continuous deployment (CI/CD), and monitoring and
 400 maintenance. The requirement management system described in Section 4.1 will be part of
 401 the project setup and requirements engineering phase.

402 Complementing the requirements management, we implemented a full software pipeline
 403 for implementation, training and validation of AI models. This pipeline could be used (a) by
 404 AI developers for continuously tracking compliance with the certification requirements, or
 405 (b) by certifiers to perform systematic tests of their clients AI models. It thus demonstrates
 406 the benefits of MLOps best practices in terms of trustworthiness by design, implemented
 407 in the AIS development, as well as in terms of facilitating an efficient means of compliance
 408 tracking and verification as part of a certification.

409 In our pipeline, models are developed, trained and versioned using Git (through Git-
 410 Hub [13]) and MLflow [22], which document all changes to the models' code and parameters,
 411 respectively. MLflow also provides the tooling for tracking experiments, packaging code, and
 412 managing model deployment. The model development process involves experimentation with
 413 different algorithms and hyperparameters to optimise performance. GitHub Actions [14] and
 414 Apache Airflow [1] are used for workflow scheduling and monitoring and facilitate automated
 415 testing and deployment processes in CI/CD pipelines. Data is versioned using Oxen [30].
 416 A schematic of the system is shown in Figure 9. The system listens for modifications to
 417 the model source code or input dataset. Changes automatically trigger training and model
 418 evaluation pipelines, which execute tests based on the methods described in sections 3.2.3,
 419 3.2.4, and 3.3. For the certification scheme, we mainly relied on methods from captum, Alibi,
 420 AIX360, ART, and UQ360, as well as original implementations from academic papers. The



■ **Figure 9** Overview of the MLOps system architecture.

421 outputs, model parameters and similar artefacts, are stored and versioned. Additionally,
 422 data engineering pipelines are run, which prepare the data for training and evaluation, and
 423 perform data-related trustworthiness evaluations.

424 MLOps provides several benefits to both AIS development and certification. Traceability
 425 and documentation are maintained throughout the AI lifecycle, providing a clear audit trail
 426 and ensuring that all objectives and means of compliance are systematically recorded. Version
 427 control is critical to maintaining the integrity of AI models and datasets, allowing teams
 428 to revert to previous versions if necessary and ensuring that all changes are documented
 429 and traceable. Automation and testing are streamlined through CI/CD pipelines, ensuring
 430 that each change is rigorously tested for compliance with trustworthiness standards before
 431 deployment. Post-deployment, continuous monitoring of AIS ensures that they remain
 432 compliant and perform reliably in real-world conditions.

433 Our workflow and methods have been tested in two real-world computer vision use cases
 434 in medical applications and vehicle detection on construction sites [10]. These use cases
 435 correspond to distinct high-risk applications according to the EU AI act. These use cases
 436 provide a test bed for validating the tools for certification on different data types and sets of
 437 requirements.

438 5 Discussion

439 The proposed certification scheme introduces several significant innovations in the assessment
 440 and certification of AIS trustworthiness, addressing an important gap in current practices.
 441 Despite the existence of standards, ethical guidelines and regulations, there remains a significant
 442 gap in the availability of practical tools and methodologies to achieve and systematically
 443 assess compliance. Our certification scheme addresses this gap by providing structured tools
 444 that are crucial for the rigorous evaluation of AIS. The scheme is underpinned by an extensive
 445 review and integration of 38 key documents from various standards and regulatory bodies,
 446 such as ISO/IEC, IEEE, EASA, and the Fraunhofer Institute. This foundational research
 447 ensures that the certification objectives and their means of compliance are comprehensive
 448 and aligned with the best practices and requirements across industries.

449 An important aspect of the scheme's implementation was evaluating multiple requirements

450 management tools to support the certification workflow. Jira was selected for its robust
451 capabilities in managing the complex certification process. This choice was crucial for
452 maintaining systematic tracking of compliance objectives, ensuring that every requirement is
453 meticulously documented and traceable.

454 Moreover, the means of compliance entail the application of metrics and technical meth-
455 ods by the customer, which can also be employed in the technical assessment of the AIS.
456 Consequently, the scheme incorporates a technical assessment based on the implementation
457 of selected technical methods which are linked to the objectives. The selection is determined
458 through an evaluation of 95 well-established and cutting-edge methods, with the evaluation
459 criteria being their suitability in meeting the defined objectives, criteria, and metrics. These
460 methods were rigorously selected and empirically tested to ensure they provide effective com-
461 pliance across various key aspects of trustworthiness, such as human oversight, transparency,
462 safety, and (cyber)security. The workflow and methods developed within the certification
463 scheme were tested on two real-life use cases: skin lesion classification and vehicle detection
464 on construction sites. These practical applications demonstrate the scheme's effectiveness and
465 adaptability in diverse, real-world scenarios. In addition, an automated workflow was imple-
466 mented on a computing cluster following MLOps principles and best practices. This workflow
467 maps MLOps stages with Trustworthy AI principles and key aspects, ensuring continuous
468 compliance and efficient lifecycle management. By automating the certification process, the
469 scheme enhances reliability, reduces human error, and ensures that the certification remains
470 up-to-date with the latest developments in AI and ML technologies. Also, due to the dynamic
471 nature of AIS and their complex post-deployment environments, trust levels can fluctuate.
472 Continuous risk monitoring is essential to maintain trustworthiness, which is in line with the
473 iterative nature of MLOps and is driven by versioning, automation, testing, deployment, and
474 monitoring. Incorporating trustworthiness metrics alongside traditional performance metrics
475 enables continuous feedback loops that systematically address trustworthiness requirements
476 throughout the AI lifecycle [48].

477 The primary focus at the beginning of the development of the certification scheme
478 was on reliability and transparency, areas where technical implementations could be more
479 straightforwardly automated. As the scheme has developed, the scope has been expanded
480 to encompass additional key areas, such as human oversight, which present more intricate
481 challenges. These aspects are inherently linked to human interaction, which makes them
482 challenging to automate effectively. The absence of established technical methods and
483 metrics in these areas presents a significant challenge. As an illustration, the assessment of
484 fairness in AI systems is an evolving field with no universally accepted metrics. This makes
485 the certification process more challenging. The scheme provides a structured approach to
486 compliance, whether through design or iterative testing and improvement. However, the
487 absence of reliable metrics makes the implementation process less clear.

488 The tools and frameworks employed in the implementation of the certification scheme
489 are designed to be adaptable, allowing the scheme to evolve in response to advances in AI
490 techniques and changing requirements. One clear example is the increasing adoption of
491 foundational models (referred to as general-purpose AI models in the EU legislation), including
492 large language models (LLMs). These models, which are trained on vast and diverse datasets,
493 introduce significant complexity due to their context-dependent and sometimes unpredictable
494 behaviour. The subjective nature of their outputs and the difficulty of quantifying their
495 decision-making processes pose challenges for evaluating and validating their trustworthiness
496 within a standardised framework. As these models are increasingly deployed across many use
497 cases, the development of new requirements, MOCs, and methods tailored to these models

498 will be vital. Addressing these challenges will be essential for maintaining the relevance and
 499 applicability of the certification scheme as AI technologies continue to advance rapidly.

500 — References —

- 501 1 Apache Software Foundation. Airflow. URL: <https://airflow.apache.org/>.
- 502 2 Atlassian. Jira, 2002. URL: <https://www.atlassian.com/software/jira>.
- 503 3 Atlassian. Confluence, 2004. URL: <https://www.atlassian.com/software/confluence>.
- 504 4 Yann Billeter, Philipp Denzel, Ricardo Chavarriaga, Oliver Forster, Frank-Peter Schilling,
 505 Stefan Brunner, Carmen Frischknecht-Gruber, Monika Ulrike Reif, and Joanna Weng. MLOps
 506 as enabler of trustworthy AI. In *11th IEEE Swiss Conference on Data Science (SDS), Zurich,
 507 Switzerland, 30-31 May 2024*, 2024. doi:10.21256/zhaw-30443.
- 508 5 Jace Browning and Robert Adams. Doorstop: Text-based requirements management using
 509 version control, 2014. doi:10.4236/jsea.2014.73020.
- 510 6 Captum. Model interpretability for pytorch, 2023. URL: <https://captum.ai/>.
- 511 7 Chun Sik Chan, Huanqi Kong, and Guanqing Liang. A comparative study of faithfulness
 512 metrics for model interpretability methods. *arXiv preprint arXiv:2204.05514*, 2022.
- 513 8 Council of European Union. Laying down harmonized rules on artificial intelligence
 514 com(2021)206 final, 2021. URL: [https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206)
 515 [CELEX:52021PC0206](https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX:52021PC0206).
- 516 9 Council of European Union. Artificial Intelligence Act: Council and Parliament Strike a
 517 Deal on the First Rules for AI in the World, 2023. URL: [https://www.consilium.europa.eu/](https://www.consilium.europa.eu/en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/)
 518 [en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-](https://www.consilium.europa.eu/en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/)
 519 [parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/](https://www.consilium.europa.eu/en/press/press-releases/2023/12/09/artificial-intelligence-act-council-and-parliament-strike-a-deal-on-the-first-worldwide-rules-for-ai/).
- 520 10 Philipp Denzel, Stefan Brunner, Yann Billeter, Oliver Forster, Carmen Frischknecht-Gruber,
 521 Monika Ulrike Reif, Frank-Peter Schilling, Joanna Weng, Ricardo Chavarriaga, Amin Amini,
 522 et al. Towards the certification of ai-based systems. In *11th IEEE Swiss Conference on Data
 523 Science (SDS), Zurich, Switzerland, 30-31 May 2024*, 2024. doi:10.21256/zhaw-30439.
- 524 11 DIN, DKE. Artificial intelligence standardization roadmap, 2023. URL: [https://www.dke.de/](https://www.dke.de/en/areas-of-work/core-safety/standardization-roadmap-ai)
 525 [en/areas-of-work/core-safety/standardization-roadmap-ai](https://www.dke.de/en/areas-of-work/core-safety/standardization-roadmap-ai).
- 526 12 EASA and Daedalean. Concepts of Design Assurance for Neural Networks (CoDANN) II.
 527 Technical report, May 2021. URL: <https://www.easa.europa.eu/en/downloads/128161/en>.
- 528 13 GitHub, Inc. GitHub . URL: <https://github.com/>.
- 529 14 GitHub, Inc. GitHub Actions. URL: <https://github.com/features/actions>.
- 530 15 Karan Goel, Nazneen Fatema Rajani, Jesse Vig, Zachary Taschdjian, Mohit Bansal, and
 531 Christopher Ré. Robustness gym: Unifying the NLP evaluation landscape. In *Proceedings
 532 of the 2021 Conference of the North American Chapter of the Association for Computational
 533 Linguistics: Human Language Technologies: Demonstrations*, pages 42–55, Online, June
 534 2021. Association for Computational Linguistics. URL: [https://www.aclweb.org/anthology/](https://www.aclweb.org/anthology/2021.naacl-demos.6)
 535 [2021.naacl-demos.6](https://www.aclweb.org/anthology/2021.naacl-demos.6).
- 536 16 Anna Hedström, Leander Weber, Daniel Krakowczyk, Dilyara Bareeva, Franz Motzkus,
 537 Wojciech Samek, Sebastian Lapuschkin, and Marina Marina M.-C. Höhne. Quantus: An
 538 explainable ai toolkit for responsible evaluation of neural network explanations and beyond.
 539 *Journal of Machine Learning Research*, 24(34):1–11, 2023. URL: [http://jmlr.org/papers/](http://jmlr.org/papers/v24/22-0142.html)
 540 [v24/22-0142.html](http://jmlr.org/papers/v24/22-0142.html).
- 541 17 IEEE. IEEE CertifAIED: the mark of AI ethics, 2022. URL: [https://](https://engagestandards.ieee.org/ieeecertifaiied.html)
 542 engagestandards.ieee.org/ieeecertifaiied.html.
- 543 18 IEEE Standards Association. IEEE Global Initiative on Ethics of Autonomous and Intelligent
 544 Systems, 2023. URL: [https://standards.ieee.org/industry-connections/ec/autonomous-](https://standards.ieee.org/industry-connections/ec/autonomous-systems/)
 545 [systems/](https://standards.ieee.org/industry-connections/ec/autonomous-systems/).
- 546 19 ISO. ISO/IEC JTC 1/SC 42 Artificial Intelligence, 2023. URL: [https://www.iso.org/](https://www.iso.org/committee/6794475.html)
 547 [committee/6794475.html](https://www.iso.org/committee/6794475.html).

- 548 20 Anna Jobin, Marcello Ienca, and Effy Vayena. The Global Landscape of AI Ethics Guidelines.
549 *Nature Machine Intelligence*, 1(9):389–399, 2019. doi:10.1038/s42256-019-0088-2.
- 550 21 Dominik Kreuzberger, Niklas Kühl, and Sebastian Hirschl. Machine Learning Operations
551 (MLOps): Overview, Definition, and Architecture. *IEEE Access*, 11:31866–31879, 2023.
552 doi:10.1109/ACCESS.2023.3262138.
- 553 22 LF Projects, LLC. MLFlow. URL: <https://mlflow.org/>.
- 554 23 LNE. Certification of processes for AI, 2023. URL: [https://www.lne.fr/en/service/](https://www.lne.fr/en/service/certification/certification-processes-ai)
555 [certification/certification-processes-ai](https://www.lne.fr/en/service/certification/certification-processes-ai).
- 556 24 Beatriz M. A. Matsui and Denise H. Goya. Mlops: A guide to its adoption in the context of
557 responsible ai. In *2022 IEEE/ACM 1st International Workshop on Software Engineering for Re-*
558 *sponsible Artificial Intelligence (SE4RAI)*, pages 45–49, 2022. doi:10.1145/3526073.3527591.
- 559 25 Microsoft. InterpretML. URL: <https://github.com/interpretml/interpret>.
- 560 26 Microsoft. Azure devops, 2005. URL: [https://azure.microsoft.com/en-us/products/](https://azure.microsoft.com/en-us/products/devops/#overview)
561 [devops/#overview](https://azure.microsoft.com/en-us/products/devops/#overview).
- 562 27 NIST. NIST Technical AI Standards, 2023. URL: [https://www.nist.gov/artificial-](https://www.nist.gov/artificial-intelligence/technical-ai-standards)
563 [intelligence/technical-ai-standards](https://www.nist.gov/artificial-intelligence/technical-ai-standards).
- 564 28 NIST. AI Risk Management Framework (AI RMF) Knowledge Base, 2024. URL: [https:](https://airc.nist.gov/AI_RMF_Knowledge_Base/AI_RMF)
565 [//airc.nist.gov/AI_RMF_Knowledge_Base/AI_RMF](https://airc.nist.gov/AI_RMF_Knowledge_Base/AI_RMF).
- 566 29 Independent High-Level Expert Group on Artificial Intelligence. Ethics guidelines for
567 trustworthy ai. Technical report, European Commission, 2019. URL: [https://digital-](https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai)
568 [strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai](https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai).
- 569 30 Oxen.ai. oxen. URL: <https://www.oxen.ai/>.
- 570 31 Nicolas Papernot, Fartash Faghri, Nicholas Carlini, Ian Goodfellow, Reuben Feinman, Alexey
571 Kurakin, Cihang Xie, Yash Sharma, Tom Brown, Aurko Roy, Alexander Matyasko, Vahid
572 Behzadan, Karen Hambardzumyan, Zhishuai Zhang, Yi-Lin Juang, Zhi Li, Ryan Sheatsley, Ab-
573 hibhav Garg, Jonathan Uesato, Willi Gierke, Yinpeng Dong, David Berthelot, Paul Hendricks,
574 Jonas Rauber, and Rujun Long. Technical report on the cleverhans v2.1.0 adversarial examples
575 library. *arXiv preprint arXiv:1610.00768*, 2018.
- 576 32 Maximilian Poretschkin et al. Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher
577 Intelligenz (KI-Prüfkatalog), 2021. URL: [https://www.iais.fraunhofer.de/de/forschung/](https://www.iais.fraunhofer.de/de/forschung/kuenstliche-intelligenz/ki-pruefkatalog.html)
578 [kuenstliche-intelligenz/ki-pruefkatalog.html](https://www.iais.fraunhofer.de/de/forschung/kuenstliche-intelligenz/ki-pruefkatalog.html).
- 579 33 Rational Software. Ibm doors, 2018. URL: [https://www.ibm.com/docs/en/engineering-](https://www.ibm.com/docs/en/engineering-lifecycle-management-suite/doors)
580 [lifecycle-management-suite/doors](https://www.ibm.com/docs/en/engineering-lifecycle-management-suite/doors).
- 581 34 Jonas Rauber, Roland Zimmermann, Matthias Bethge, and Wieland Brendel. Foolbox
582 native: Fast adversarial attacks to benchmark the robustness of machine learning models
583 in pytorch, tensorflow, and jax. *Journal of Open Source Software*, 5(53):2607, 2020. doi:
584 10.21105/joss.02607.
- 585 35 SAE. Taxonomy and Definitions for Terms Related to Driving Automation Systems for
586 On-Road Motor Vehicles. SAE J3016, 2021. URL: [https://www.sae.org/standards/content/](https://www.sae.org/standards/content/j3016_202104/)
587 [j3016_202104/](https://www.sae.org/standards/content/j3016_202104/).
- 588 36 Wojciech Samek et al. *Explainable AI: Interpreting, Explaining and Visualizing Deep Learning*.
589 Springer, 2019. doi:10.1007/978-3-030-28954-6.
- 590 37 Anna Schmitz, Michael Mock, Rebekka Gorge, Armin B Cremers, and Maximilian Poretschkin.
591 A global scale comparison of risk aggregation in ai assessment frameworks. *AI and Ethics*,
592 pages 1–26, 2024.
- 593 38 Seldon. Alibi Explain. URL: <https://github.com/SeldonIO/alibi>.
- 594 39 sicara. TF-Explain: Interpretability Methods for tf.keras Models with Tensorflow 2.x. URL:
595 <https://github.com/sicara/tf-explain>.
- 596 40 Siemens. Polarion, 2004. URL: <https://polarion.plm.automation.siemens.com/>.
- 597 41 Ryan Soklaski, Justin Goodwin, Olivia Brown, Michael Yee, and Jason Matterer. Tools and
598 practices for responsible ai engineering. *arXiv preprint arXiv:2201.05647*, 2022.

- 599 42 Guillaume Soudain. First usable guidance for Level 1 machine learning applications: A deliv-
600 erable of the EASA AI Roadmap, 2021. URL: [https://www.easa.europa.eu/en/downloads/
601 134357/en](https://www.easa.europa.eu/en/downloads/134357/en).
- 602 43 The White House. Fact Sheet: President Biden Issues Executive Order on Safe, Secure, and
603 Trustworthy Artificial Intelligence, 2023. URL: [https://www.whitehouse.gov/briefing-room/
605 statements-releases/2023/10/30/fact-sheet-president-biden-issues-executive-
606 order-on-safe-secure-and-trustworthy-artificial-intelligence/](https://www.whitehouse.gov/briefing-room/statements-releases/2023/10/30/fact-sheet-president-biden-issues-executive-
604 order-on-safe-secure-and-trustworthy-artificial-intelligence/).
- 606 44 Trusted-AI LF AI Foundation. Adversarial Robustness Toolbox (ART). URL: [https://
607 github.com/Trusted-AI/adversarial-robustness-toolbox](https://github.com/Trusted-AI/adversarial-robustness-toolbox).
- 608 45 Trusted-AI LF AI Foundation. AI Explainability 360 (AIX360). URL: [https://github.com/
609 Trusted-AI/AIX360](https://github.com/Trusted-AI/AIX360).
- 610 46 Trusted-AI LF AI Foundation. AI Uncertainty Quantification 360 (UQ360). URL: [https:
611 //github.com/Trusted-AI/UQ360](https://github.com/Trusted-AI/UQ360).
- 612 47 Arnaud Van Looveren et al. Alibi detect: Algorithms for outlier, adversarial and drift detection,
613 2019. URL: <https://github.com/SeldonIO/alibi-detect>.
- 614 48 Larysa Visengeriyeva, Anja Kammer, Isabel Bär, Alexander Kniesz, and Michael Plöd. MLOps
615 Principles, 2020. URL: <https://ml-ops.org/content/mlops-principles>.