# Lost in Vagueness: Towards Context-Sensitive Standards for Robustness Assessment under the EU AI Act

**Roberta Tamponi**[1]* **Carina Prunkl**[2,3] **Thomas Bäck**[1] **Anna V. Kononova**[1]

[1] Leiden Institute of Advanced Computer Science, Leiden University

[2] Inria  [3] University of Oxford

## Abstract

Robustness is a key requirement for high-risk AI systems under the EU Artificial Intelligence Act (AI Act). However, both its definition and the methodologies for assessing it remain underspecified, leaving providers with little concrete direction on how to demonstrate compliance. This stems from the Act's horizontal approach, which establishes general obligations applicable across all AI systems, but leaves the task of providing technical guidance to harmonised standards. This paper investigates what it means for AI systems to be *robust* and illustrates the need for context-sensitive standardisation. We argue that robustness is not a fixed property of a system, but depends on which aspects of performance are expected to remain stable ("robustness of what"), the perturbations the system must withstand ("robustness to what") and the operational environment. We identify three contextual drivers—use case, data and model—that shape the relevant perturbations and influence the choice of tests, metrics and benchmarks used to evaluate robustness. The need to provide at least a range of technical options that providers can assess and implement in light of the system's purpose is explicitly recognised by the standardisation request for the AI Act, but planned standards, still focused on horizontal coverage, do not yet offer this level of detail. Building on this, we propose a multi-layered standardisation framework that embeds context-sensitivity in the ongoing standardisation process. Horizontal standards define shared principles and terminology, while domain-specific ones identify relevant risks and map them across the AI lifecycle, guiding the selection of best practices for specific applications. These, together with benchmarks and sandboxes, should be organised in a dynamic repository where providers can propose new informative methods and share lessons learned. Such a system reduces the interpretative burden, mitigates arbitrariness and addresses the obsolescence of static standards, ensuring that robustness assessment is both adaptable and operationally meaningful.

## 1 Introduction

In 2024, the European Union (EU) adopted the Artificial Intelligence Act (AI Act), the world's first comprehensive legal framework for AI technologies [1]. It addresses concerns about safety, reliability and fundamental rights, and ensures AI systems align with European values. The Act builds on the work of the EU High-Level Expert Group [2], which identified three requirements for *trustworthy* AI: lawful, ethical and technically robust. On this basis, the AI Act adopts a risk-based approach: obligations vary depending on the level of risk an AI system poses. Providers of high-risk AI systems are required to conduct conformity assessments before placing their products on the European market. A central challenge is that, despite outlining procedural obligations, the AI Act offers little detail on

---

*Corresponding author: `r.tamponi@liacs.leidenuniv.nl`

how these should be put into practice [3]. This lack of clarity reflects the Act's horizontal approach, which sets general requirements for all AI systems and defers practical implementation to harmonised standards.

Among the obligations, *robustness* is a key requirement [4]. However, both the regulation and the Code of Practice, leave its definition and assessment underspecified, raising questions about what must be demonstrated and how. The notion of robustness is ubiquitous in the field of AI, but its exact meaning varies depending on the context [5]. Researchers either focus on narrow technical definitions or leave the concept too abstract, aiming for generality. This ambiguity reflects the concept's *context-sensitive nature*: the robustness of a system depends on its intended purpose, design and deployment environment, which shape the perturbations and potential failures most relevant for evaluation. Hence, the practical assessment is itself context-sensitive. The choice of robustness tests, metrics, thresholds and mitigation strategies must be tailored to operating conditions, where even minor changes can result in markedly different evaluation procedures. The only international standard on robustness adopted so far by the EU, the CEN/CLC ISO/IEC/TR 24029-1 [6], [2] is highly generic and only offers an overview of the assessment methods for neural networks. Other standards on robustness are under development, but their specificity remains uncertain. The ISO/IEC 24029-3, for instance, is expected to expand the existing series and include tests and metrics [7]. Yet, it seems unlikely that a single extension could provide the level of detail needed to identify, assess and mitigate robustness risks across diverse contexts. Moreover, it focuses only on neural networks, neglecting the broader diversity of AI systems. The risk of overly generic guidance, therefore, remains, leaving providers with excessive interpretative burden and the risk of arbitrary implementation. The AI Act's standardisation request explicitly acknowledges the need to account for different use cases and sectors and to develop vertical specifications where appropriate [8]. European standards need to set clear requirements for AI providers, yet it seems unlikely that horizontal standards alone can adequately fulfil the mandate.

**Statement of Contributions.** We propose a clearer conceptualisation of robustness, distinguishing the *robustness of what* from *the robustness to what*. We further identify three contextual drivers—use case, data and model—which determine the relevant perturbations and the most suitable methods to assess system robustness. Finally, we propose a framework to embed context-sensitivity into the ongoing standardisation process, addressing the limited scope of horizontal standards, which risks leaving too much room for arbitrary choices, and the static nature of standardisation, which struggles to keep pace with technological change.

The paper first outlines the AI Act and its treatment of robustness (Section 2). Section 3 highlights the limits of overly abstract definitions. Section 4 distinguishes two dimensions of robustness: the robustness of what (the performance) and the robustness to what (the types of perturbations the system must withstand). Section 5 examines how the contextual drivers shape relevant perturbations as well as the evaluation methods and concludes with a use case comparison. Section 6 discusses current standardisation challenges and Section 7 proposes a multilayered framework combining horizontal standards with domain-specific provisions, supported by a dynamic repository of practices, benchmarks and sandboxes, which also allows providers to propose new informative methodologies and share experiences.

## 2   Robustness in the AI Act

The Artificial Intelligence Act, adopted in 2024, establishes a product-safety-based regulatory framework that applies different obligations depending on the risk level of AI systems. *Risk* is defined as the combination of the probability of harm and the severity of harm (Art. 3) and systems are classified into four categories: *unacceptable, high, limited,* and *minimal* risk. The AI Act primarily targets high-risk systems–those with potentially significant impacts on safety, fundamental rights, or critical domains. These must obtain the CE marking (Art. 48) through a conformity assessment (Art. 43, Annex VI–VII), demonstrating compliance with Chapter III, Section 2, which sets requirements on risk management (Art. 9), data governance and quality (Art. 10), record keeping (Art. 12), transparency (Art. 13), human oversight (Art. 14), accuracy, robustness and cybersecurity (Art. 15). Providers must maintain technical documentation (Art. 11), implement a quality management system

---

[2]"CEN/CLC ISO/IEC" indicates that the technical report (TR) by ISO (International Organisation for Standardisation) and IEC (International Electrotechnical Commission) has been adopted as a European standard.

(Art. 17) and, after deployment, monitor performance using tools defined during the development. The provisions are entering into force gradually, starting in August 2024, allowing time for the development of harmonised standards, providers' adaptation and the establishment of supervisory authorities.

Robustness is thus a core requirement. Article 15 states that AI systems must be designed to ensure an appropriate level of accuracy, robustness and cybersecurity throughout their lifecycle (Art. 15.1) and to be as resilient as possible to errors, faults and inconsistencies (Art. 15.3). Systems that continue to learn after being placed on the market have to include mitigation measures to minimise the risk of biased outputs arising from feedback loops (Art. 15.4). These provisions are vague as they only set general objectives without providing parameters, metrics, or acceptability thresholds. The regulation itself acknowledges that benchmarks and measurement methodologies are still lacking and indicates that the Commission is working with stakeholders to develop them (Art. 15.2). This is coherent with the will of creating a horizontal regulation that applies to all AI systems, but highlights the need for complementary, detailed standards. Yet, despite the ongoing standardisation process, major uncertainties remain about the level of detail these standards will provide, and, in their current form, both the definition of robustness and the related provisions remain insufficiently specified.

## 3    Limitations of an abstract definition of robustness

The notion of robustness in AI remains conceptually ambiguous [5] and has taken multiple meanings across different domains and contexts [9]. The existing harmonised standard on the robustness assessment of neural networks [6] describes robustness as the ability of an AI system to *maintain its level of performance under any circumstances*. While this definition may be adopted for regulatory purposes such as the AI Act, access to standards is restricted by copyright and they are not publicly available. As a result, researchers often rely on alternative definitions. For example, Nobandegani et al. [10] define it as the insensitivity of a model's performance to miscalculations of its parameters. Freiesleben and Grote [5] describe it as a model's capacity to sustain stable predictive performance in the face of variations in input data. Across subfields, the term has taken on multiple meanings. In computer vision, the term has been used to denote raw performance on held-out test sets, generalisation within and across domains, maintaining performance on manipulated inputs and naturally-induced image corruptions and resistance to adversarial attacks [11]. In natural language processing (NLP), La Malfa et al. [12] frame robustness of text classification in terms of the capacity to maintain predictions under word substitutions, formalised as a maximal safe radius for a given input text. Hendrycks et al. [13], focusing on out-of-distribution generalisation, understand robustness as the ability to handle unforeseen real-world distributional shifts.

Despite the variety of formulations, the underlying principle to which they all refer is *the stability of performance*. This broad intuition can be captured in the high-level definition: *a model is robust when small changes in input cause small changes in output*. While this is an intuitive starting point, it remains too abstract to be operationally useful. What is a small change in input? Which output? To answer these, a more comprehensive definition can be: *a system is considered robust to some kind of small perturbations in the input data when the system's performance under the perturbed state remains similar to that of the unperturbed state*. But then, which perturbations? All of them or only some? Which performances must be stable? How stable? Going deeper with the necessary clarifications, which typically depend on the context in which the system is developed and deployed, leads us to case-specific definitions. Therefore, *robustness is not a fixed notion of a system but depends on the system's different assumptions and goals*. As a result, a system considered robust in one context may not be robust in another.

## 4    Robustness dimensions

Building on this context-sensitive understanding, we propose a structured way of interpreting robustness that supports practical evaluation. In line with Freiesleben and Grote [5], who frame robustness as a relation between targets, modifiers, domains and tolerance thresholds, we share the view that robustness must be context-specific. Our focus, however, is on linking abstract definitions to practical evaluation. Starting from the principle of the stability of performance, we distinguish between two core aspects: the system's output that is required to remain stable (*robustness of what*) and the types of changes in input it is required to withstand (*robustness to what*).

### 4.1 Robustness of what?

To assess robustness in practice, we must first specify what is expected to remain stable, which, in machine learning, is the system's *performance*, as robustness concerns the system's ability to maintain adequate predictive capabilities under non-ideal or varying conditions. In line with the terminology used in current international standardisation efforts, we adopt the definition provided in BS EN ISO/IEC 25059 [14], which states that, within the field of AI, performance means *how well a certain AI system performs the intended tasks*. Performance can be measured with relevant metrics that depend on the system's purpose. CEN/CLC ISO/IEC/TR 24029-1 [6] proposes an exhaustive list which includes accuracy, precision, recall, specificity, $F_1$ score, ROC and AUC. These can be used independently or in combination, but they are not interchangeable. Notably, there are also numerous task-specific metrics, such as BLEU, TER and METEOR for machine translation, or mean average precision for ranked retrieval. To be meaningful indicators of performance, metrics need to be chosen based on the system's purpose and the impact of the different types of errors. For example, in a child-safety video filter, where the retrieved set is the non-harmful content, high precision is preferred: excluding some acceptable videos (low recall) is better than showing harmful content. In shoplifter detection, on the other hand and depending on operational preferences, false alerts may be tolerated if this ensures that most shoplifters are detected, so recall may be prioritised over precision [15]. Since such measures often trade off, the system's function and associated risks must be understood before identifying which aspects of performance are most crucial to maintain under perturbations. Robustness evaluation assumes a fully trained and tested system with established performance under standard conditions. This baseline is then used to assess how perturbations affect performance and to define tolerable deviation thresholds. These thresholds depend on deployment context, with a lower tolerance level due to higher associated risks (e.g., clinical decision-making).

### 4.2 Robustness to what?

After clarifying the role of performance in robustness evaluation, the next step is to identify against what the model must be robust. Robustness assessment concerns how input perturbations affect the system's predictive capabilities and whether the model can maintain its intended behaviour. Possible threats identified in BS EN ISO/IEC 25059 [14] include unseen, biased, adversarial or invalid data, external interference and environmental conditions; the AI Act similarly points to errors, faults, inconsistencies and unexpected situations (Recital 75). Some minor perturbations in input are unavoidable, while others might lead to more severe problems. A report of the latter must be documented under the risk management obligations of Art. 9, which requires the identification, estimation and evaluation of foreseeable risks. Identifying which input changes are critical depends on the system's intended use; therefore, the robustness goal must be clearly specified. Generic instructions, such as "robust to dissimilar inputs", are too vague, since a model may generalise well across different distributions in the same application domain, but fail entirely on inputs of unrelated domains [6]. Assessing robustness requires a precise goal, defined by narrowing down the relevant perturbations that the system can face. The following section examines the contextual drivers that shape the identification of these perturbations.

## 5 Context sensitivity of robustness assessment

Not all types of perturbations are equally relevant for every system, so it is necessary to identify those most critical for the intended application. For self-driving cars, protecting against adversarial attacks is crucial, but it is not necessarily relevant for an epidemiologist using AI models to forecast virus spread [5]. Relevant perturbations cannot be defined universally, but depend on system design and deployment environment. We identify three contextual drivers—use case, data and model architecture—that guide the selection of relevant perturbations and, in turn, inform the tests, metrics and benchmarks through which robustness is evaluated in a context-sensitive way. Figure 1 provides a schematic overview of context-sensitive robustness assessment.

### 5.1 Contextual drivers of robustness

**Use case.** When developing an AI system, its intended use case needs to be specified. This includes its application domain, the task and the deployment environment in which it is expected to operate. The AI Act (Annexes I and III) gives several examples of domains where AI use can be considered
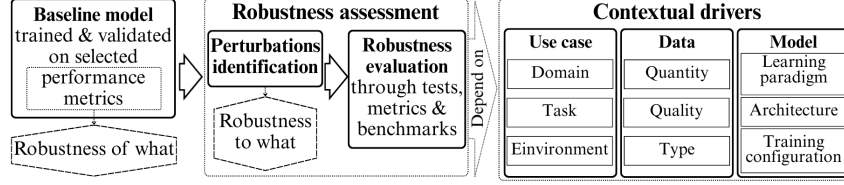
Figure 1: Conceptual pipeline for context-sensitive robustness assessment.

high-risk, such as healthcare, education, employment, access to services and law enforcement. Within each domain, AI may serve different tasks, for example, real-time health monitoring, school grading, needs-based resource allocation, etc. The deployment environment refers to the real-world conditions in which the AI system is expected to operate, which is typically dynamic, partially observable and/or subject to noise or user interaction. Addressing these different aspects means delineating the problem that the AI system is meant to solve and how, by clarifying the goals and the requirements that need to be met [3, 16]. In healthcare, for example, systems may require robustness to co-occurring conditions such as gastrointestinal or neurological disorders, or robustness to variations in computed tomography protocols, including scanner model, patient size, or radiation dose [17]. Understanding the use case helps identify the system's potential failures and limitations, guiding both the perturbations to be addressed and the criteria for assessing them.

**Data.** Data quantity, quality and type determine applicable models, relevant perturbations and suitable evaluation metrics. A main challenge while building an AI system is the lack of training data: even basic problems need thousands of examples, while complex ones may require millions [15]. Shallow learning models like decision trees or support vector machines are often effective when data is limited, while neural networks require larger amounts of data. Another challenge is the presence of non-representative and poor-quality data. A system can learn and detect underlying patterns only if the training data include enough relevant features that are not full of errors, outliers and noise [15]. Robustness to data drift depends on how performance degrades over time due to shifts in the patterns or in the environment [18]. Low-quality data are detrimental to adversarial training, as they cause robust overfitting and robustness overestimation against weaker adversaries such as projected gradient descent (PGD) [19]. The use of large and diverse datasets increases generalisability, making a system more robust to variations and noise sources in data [20]. Finally, the data type: even when robustness is framed in terms of distribution shift, its manifestation is data-specific. In images, these perturbations may appear as noise, blur or weather effects, each testable through specific techniques, while for text they can be probed at character, word, or sentence level, reflecting typos, substitutions, or stylistic changes, each requiring tailored evaluation methods [21].

**Model.** Different models have different sensitivities to data variations and robustness depends on learning paradigm, architecture and training configuration choices. In terms of learning paradigm, basic classification methods, such as decision trees or rule-based systems, may suffice for structured tasks like automated coding of medical records, while highly sensitive applications, such as medical image analysis or patient outcome prediction, require models that can handle complex, high-dimensional data, like deep learning models [3, 16]. At the architectural level, linear classifiers like Naive Bayes might be preferred due to their efficiency and interpretability, but if the system is expected to deal with non-linear perturbations, then models such as k-Nearest Neighbours might be more appropriate [22]. Moreover, some models' architectures can better mitigate the accuracy/robustness trade-off than others, balancing the performance on clean data with the architecture's ability to contrast adversarial perturbations [23, 24]. Also training configuration has to balance this trade-off. The inclusion of batch normalisation has been shown to improve model accuracy, but in some models, it increases the vulnerability to adversarial attacks [25, 26]. Hyperparameters such as learning rate, batch size and regularisation coefficients are likewise critical for robustness and optimisation strategies like grid or random search can substantially enhance it by fine-tuning these settings [20, 27, 28]

5

## 5.2 Robustness evaluation

Use case, data and model can guide the selection of relevant perturbations. Unfortunately, there is no unambiguous way of defining them, since they may originate from diverse sources and the deployment environment can introduce unpredictable forms of uncertainty that are difficult to anticipate or quantify. The literature commonly organises robustness into two macro-categories: adversarial robustness and natural robustness [9, 29]. *Adversarial robustness* describes the ability of a model to maintain stable performance when exposed to intentional manipulations of the input data, designed to mislead it [30]. These perturbations are non-random and precisely constructed to subtly change the input and maximise the probability that the model will produce an incorrect prediction [31]. *Natural robustness* is also called distribution shift or out-of-distribution data and concerns the robustness of the model to perturbations that occur spontaneously in the deployment environment, such as noise or visual alterations caused by external factors [11]. Natural perturbations introduce a deviation between the distribution of the test data and the one on which the model was originally trained [32].

Once the relevant types of perturbations that may pose a risk to the system have been identified, it is time to test whether the system can remain functional. To test robustness means submitting the model to various types of perturbations to identify vulnerabilities and weaknesses [33]. Various testing techniques have been investigated. For adversarial attacks, for example, generative adversarial networks (GANs) [34] can be used to synthesise challenging test cases. In the metamorphic testing [35, 36], the model's behaviour is evaluated under meaningful input transformations. Mutation testing [37, 38] introduces small faults into the test data to evaluate whether the model can detect them. To evaluate the test, metrics are needed to quantify how much a model's performance degrades (or resists) when exposed to perturbations. CEN/CLC ISO/IEC/TR 24029-1 [6] overview on the robustness of neural networks assessment distinguishes *statistical methods*, where robustness is defined as a limited performance drop relative to an unperturbed reference set, *formal methods*, which provide mathematical proofs but rely on restrictive assumptions, and *empirical methods*, which rely on experimentation, observation and expert judgement.

Benchmarking an AI system can also help assess some degree of its robustness, as it provides shared and standardised conditions to compare models under controlled perturbations in a reproducible way. The AI Act encourages the development of benchmarks to support the evaluation of high-risk AI systems, but leaves open what counts as an appropriate benchmark and who should be responsible for their creation or maintenance. Some benchmarks focus on specific types of robustness. RobustBench standardises adversarial robustness evaluation across models [39], while WILDS targets robustness to distribution shifts, including domain generalisation and subpopulation shift [32]. Some are tailored to a learning paradigm: OpenAI Gym [40] supports robustness testing in reinforcement learning. While others focus on optimisation methods, for example, in heuristic optimisation, IOHprofiler evaluates algorithm selection and configuration under perturbed scenarios [41, 42]. Benchmarks can thus help establish initial trust in an AI solution even before deployment, but their relevance depends on the type of robustness being tested and the task at hand.

## 5.3 Comparing use cases

Even within similar contexts, what makes a system robust may change drastically, and with it, its assessment. Table 1 compares two studies on AI-based medical image classification for disease detection. The first investigates convolutional neural networks (CNN) classifiers for renal cell carcinoma [26]. The second examines CNN for tuberculosis-related chest radiographs classification [43]. Use case 1 trains on data from The Cancer Genome Atlas and tests on data from the University of Aachen, while use case 2 trains on CXR from Shenzhen Hospital (China) and tests both on a held-out Shenzhen set and on the NIH ChestX-ray8 dataset (USA), which includes tuberculosis-related features but also other lung abnormalities. Robustness in the first case is tested against multiple adversarial attacks that differ in nature and strength (white-box vs. black-box, including ensemble and transformation-based attacks), complemented by adversarial training, dual batch normalisation and evaluation of attack success rates on selected image subsets, while in the second case it is assessed through domain shift, testing cross-population generalisability from Chinese to US data. In both cases, robustness is measured by comparing the AUROC (ability to discriminate between true and false positives) before and after perturbations, with use case 1 also reporting attack success rates to capture the degree of misclassification under adversarial inputs.

Table 1: Comparison of two use cases where slightly different design choices lead to different perturbations and therefore very different methods for robustness evaluation.

| Details | Use case 1 | Use case 2 |
|---------|-----------|-----------|
| Use case | Detection of renal cell carcinoma (RCC) | Detection of tuberculosis-related features |
| Data | Whole-Slide Images (WSIs) of Hematoxylin and Eosin stained tissue (RCC samples) taken from two different datasets for training and testing | Chest radiograph (CXR) taken from two datasets: (1) confirmed tuberculosis (2) tuberculosis-associated features |
| Model | CNN (ResNet) | CNN (Inception V3) |
| Perturbation | Adversarial attacks: FGSM, PGD, FAB, Square Attacks, AutoAttack, AdvDrop | Distribution shift |
| Robustness test | 1) Multiple adversarial attacks using the 6 attacks, each with three different attack strengths 2) Adversarially robust training 3) Dual Batch Normalisation training 4) Selection of 450 tiles from the RCC set to test the degree of misclassification | Train the model with one dataset and test it on one with a different distribution |
| Evaluation metric | 1) AUROC before and after the perturbations 2) Attack Success Rate | AUROC before and after distribution shift |

This comparison shows that even within similar contexts, disease detection through CNN on image data, what defines the system as robust can change significantly. Small shifts in assumptions—such as the deployment environment—lead to very different evaluations. In digital pathology workflows, where adversarial manipulations are a plausible risk, robustness must be tested against gradient-based perturbations, minimal decision-boundary changes, random or combined attacks and high-frequency manipulations. On the other hand, the CXR model expects to face the challenge of distribution drift, requiring robustness to differences in population between training and deployment data. Similar variation occurs even within the same class of perturbations: for instance, in large language models for medical question-answering, adversarial robustness can be tested through MedFuzz attacks that introduce misleading patient characteristics [44].

## 6 Current standardisation directions

Harmonised standards may be used to demonstrate compliance with the AI Act, and while not mandatory, they grant presumption of conformity (Art. 40) and are fundamental for the practical implementation of the Regulation. The European Commission is working with the European Standardisation Organisations (ESO), formed by the European Committee for Standardisation (CEN), the European Committee for Electrotechnical Standardisation (CENELEC) and the European Telecommunications Standards Institute (ETSI), to develop technical specifications for the AI Act. In May 2023, CEN and CENELEC formally accepted the standardisation request M/593. The JTC 21 (Joint Technical Committee), established by CEN and CENELEC, is the dedicated body for developing harmonised standards in support of the Act, complementing and adapting existing international standards to the Act requirements. At the time of writing, their work programme includes 25 ongoing activities and 15 published standards, [3] which include both CEN/CENELEC developments and ISO/IEC standards formally adopted as European ones. The deadline for the deliverables was set for the 30th of April 2025, but they are now expected to be completed by early 2026 [45].

Although harmonised standards are meant to provide the structured guidance needed to implement the AI Act in practice, the Joint Research Centre (JRC) report *Harmonised Standards for the European AI Act* emphasises that standards should consist of a small number of horizontal requirements, applicable across different AI systems and sectors, complemented by sector-specific provisions only when strictly necessary [4]. At the same time, the report expects these horizontal requirements to be "sufficiently prescriptive and clear" to support practical application for specific use and the identified

---

[3]See the CEN/CLC/JTC 21 Work programme and Published Standards.

risks, guiding providers in selecting appropriate tests and measures. Therefore, the report demands standards to be both generic and context-sensitive, a hardly feasible ambition. By contrast, the Standardisation Request M/593 [8] is much more straightforward and specifies in concrete terms what European standards should deliver. The request indeed calls for European standards to take into account the risks common (horizontal) to AI systems in general. Yet, notwithstanding their horizontal nature, it also makes clear that standards may—and where relevant should—provide specifications for particular domains and use cases, taking into account the intended purpose and context of use of those systems. In addition, standards must reflect the generally acknowledged state of the art and provide concrete, verifiable technical specifications, including design and development requirements, as well as verification, validation and testing procedures. Finally, even if it is not possible to cover every specific intended purpose, standards are expected to set out at least a range of technical options that providers can assess and implement in light of the purpose of their system, together with guidance on how such options should be applied. Therefore, the request explicitly acknowledges the need for a context-sensitive approach and this makes it difficult to argue that horizontal standards alone could fulfil the mandate. Horizontal standardisation may be advantageous for efficiency in development and maintainability, but when the requirements to standardise vary significantly with context—as in the case of robustness—it risks becoming overly generic and detached from actual practice.

The current work on robustness standards focuses mainly on horizontal coverage, with limited evidence of vertical specifications. The ISO/IEC 24029 series addresses robustness at a high level: Part 1 provides an overview of approaches for evaluating neural networks, Part 2 describes formal methods but has not yet been adopted as a European standard and a planned Part 3 is expected to cover methodologies for assessing adversarial robustness and distribution shift. Other planned deliverables include the AI trustworthiness framework, which treats robustness and other several dimensions, and two specifications for computer vision and NLP models. [4] Yet robustness is still addressed almost exclusively for neural networks, treated as a uniform block without differentiation by architecture, application domain, or operational constraints, leaving other types of AI systems—such as supervised, unsupervised, generative and general-purpose models—largely uncovered. To be effective, robustness standards should set clear requirements that allow providers to select relevant perturbations, tests, metrics and thresholds. Although Part 3 appears promising, its scope and level of specificity remain to be seen. It is unlikely that this part alone will be able to provide all the necessary specifications that account for different contexts. As a result, neither the current nor the planned robustness standards appear sufficient to meet the standardisation request, highlighting the need for a layered approach that links horizontal standards to more detailed specifications.

## 7 Context-sensitive standardisation through a multi-layered approach

Horizontal standards provide a useful common core, but they must explicitly incorporate contextual hooks and minimal evidence requirements to adequately address the context-sensitive nature of robustness. Completely changing the current standardisation process may not be feasible, but its course can be adjusted. Horizontal standards could be designed as a starting point that establishes high-level principles, shared terminology and general best practices, while linking to more detailed standards that address domain-specific requirements. A multi-layered approach would ensure coherence through common principles set in the horizontal documents and adaptability through tailored provisions for specific risks and contexts, improving auditability and allowing organisations to follow clear directions without having to interpret abstract requirements. In what follows, we propose a framework, summarised in Figure 2, that directly addresses the shortcomings of current standards. It aims to reduce ambiguity by going beyond horizontal provisions and by identifying domain-specific risks across the AI lifecycle, guiding the choice of best practices for compliance with the AI Act context-sensitive requirements. Given the static nature of standards, which risks guiding to obsolete methodologies, as technology evolves faster than bureaucratic processes, the framework proposes a context-sensitive and dynamic repository of recommended practices where stakeholders can propose new methods and experiences. This would enable the ESOs to track technological developments, update practices more efficiently and provide a shared platform that supports providers both in the design and the validation of their systems. This proposal also supports the objectives of the so-called HAS assessment process checklist. [5] In particular, it better satisfies two central criteria: that the

---

[4] Deliverables listed in the "Complete overview of the work programme" document on JTC21.

[5] A procedure to evaluate if a deliverable complies with standardisation requests and EU legal requirements.
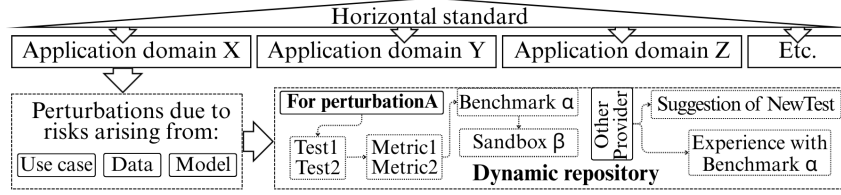
Figure 2: Multi-layered framework for context-sensitive robustness evaluation under the AI Act.

mandatory parts of a standard are written clearly, leaving no room for arbitrary choices, and that risk assessment is complete, with evidence that all relevant risks have been considered. The goal is not to replace existing processes but to create a more structured framework where horizontal and context-sensitive standards complement each other, bringing clarity at the top level while respecting the diversity of needs on the ground.

## 7.1 Domain-specific standardisation

Robustness standards need to be domain-specific. AI is a broad umbrella term that includes a large number of subfields [46]. The AI Act identifies several high-risk application areas (see Section 5.1) and adapting horizontal rules to such diverse applications is problematic, since each domain has unique characteristics [47]. Different contexts require distinct ethical considerations and the requirements for AI systems change depending on the underlying technology and domain in which they are deployed. Domain-specific standardisation enables legal obligations to be translated into measurable criteria and better aligned with sectoral ethical priorities [48]. For example, regulations that limit the use of patient data, suitable in general AI applications, may need to be modified in healthcare to allow the use of de-identified records [49], so that the ethical principles that guide the industry, such as autonomy, beneficence, non-maleficence and justice, are respected [50]. Domain-specific standards do not entail fragmentation, as long as they are built on horizontal ones and would ensure completeness, alignment with general requirements and prevent ambiguities or arbitrary interpretations.

## 7.2 Context-sensitive perturbation taxonomy

Soler Garrido et al. [51] highlight the need for guidance on how to set acceptable thresholds for different robustness methods, taking into account the context of use and risks of specific AI systems. Building on this, for each domain, there should be an exhaustive perturbation taxonomy. As discussed in Section 5, perturbation identification is context-sensitive. Therefore, guidelines should map potential perturbations across the contextual drivers—use case, data, model—illustrate their root causes and suggest mitigation strategies. Ultimately, at each stage of implementation, providers decide the direction the system should take based on the perturbations identified up to that point. A similar proposal was advanced by Schnitzer et al. [52], who proposed a taxonomy of AI hazards linked to lifecycle stages. Although mainly focused on DNN-related hazards, it offers a solid foundation that could be adapted both to robustness deliverables and to other requirements.

## 7.3 Dynamic repository of best practices

To address robustness tests, metrics and thresholds, it is necessary to refer to state-of-the-art practices, understood as widely accepted good practice at a given time [53]. To support this, our framework proposes the development of a repository of best practices, maintained by the ESOs and informed by stakeholder input. The repository should link these best practices and their limitations to the relevant risks to robustness—and to other requirements—that may arise during development and deployment, thereby supporting their evaluation and mitigation, including mitigation measures against biased outputs arising from feedback loops. Given the rapid evolution of technology and the risk of standards obsolescence, the repository should allow providers to share experiences and propose new informative methodologies, complementing recommended practices until official updates are made. The proposed methodology would inform the ESOs of the areas requiring updates, making the process faster, while shared experiences would support providers in selecting effective methods for their use cases and aligning system development with regulatory expectations. Since identifying all possible risks is unrealistic, standards will supply a range of technical options that providers can assess and implement

9

in light of the purpose of their system, while the repository supports providers in selecting those most appropriate for their context. By remaining flexible, outcome-focused and adaptable to different system contexts, it ensures technological neutrality and adherence to performance-based principles. A similar initiative is the Artificial Intelligence Measurement and Evaluation (AIME) [54], launched by the U.S. agency NIST, which develops metrics and methods for robustness and performance in areas such as language, vision and robotics. The program does not offer a centralised repository, but in collaboration with public and private sectors, NIST has developed a voluntary framework to manage AI-related risks for individuals, organisations and society. [6] The OECD has also developed a Catalogue of Tools & Metrics for Trustworthy AI, [7] which collects practices and metrics for fairness, transparency, explainability, robustness and safety. While it offers a valuable foundation, it is not up to date nor designed for context-sensitivity, but it could serve as a reference for a dynamic European AI repository.

## 7.4 External validation

Context-sensitivity should also be extended to the external validation tools referenced in the AI Act: benchmarks and sandboxes. The repository should indicate which benchmarks, used as post-development evaluation tools, are most suitable for different tasks and categories of high-risk systems. While a single benchmark can cover multiple systems (see Section 5.3), interpretation of results must remain context-sensitive and should be interpreted with care [55], therefore, sharing experiences would help providers better contextualise outcomes and highlight practical limitations. Moreover, benchmarks risk becoming obsolete within months [56], making the repository crucial for collecting and updating proposals for new benchmarks. Since benchmarks can fail to account for data or domain shifts, regulatory sandboxes introduced by the AI Act (Art. 57) provide a valuable complement to ensure that systems are fit for deployment. Sandboxes are controlled environments, established and supervised by national competent authorities, that facilitate technical testing under real or close-to-real conditions, providing resources, data and testing infrastructure [57]. Yet, like best practices and benchmarks, sandboxes face context-sensitivity due to the different needs of systems across sectors, where a one-size-fits-all solution is not feasible. As Due et al. [58] note, stakeholders already stress the difficulty of designing a sandbox that can cover all participants' needs, reinforcing the case for domain-specific sandboxes, especially for high-risk AI systems, even if the Act mandates horizontal national ones. Including context-sensitive sandboxes in the repository would guide use-case-specific testing and experience sharing. In this way, the repository becomes a practical tool to support both model design and model evaluation, enabling AI systems to be iteratively improved, compliant and aligned with context-specific requirements.

## 8 Conclusions

Robustness has a context-sensitive nature. Its assessment should follow a structured approach that considers the system's intended use, the data available and the model architecture that would best fit the goal. The analysis of each of these contextual drivers narrows down the perturbations the system should be robust against and enables the selection of appropriate tests, evaluation and mitigation metrics. Shifts in contextual assumptions can lead to substantial differences in the assessment process. Since only harmonised standards confer presumption of conformity under the AI Act, this variability must be accounted for in their ongoing development. The current approach aims to design standards applicable to various types of AI systems across sectors. However, some requirements, such as robustness, are highly context-sensitive and horizontal standardisation may result in high-level directions that are operationally meaningless, leaving too much space for interpretation that could cause superficial compliance.

Standards could be designed with a multi-layered approach that establishes high-level principles, but then points to specific methodologies in each domain. To ensure robustness, a perturbation taxonomy is needed and it should be mapped into the system's lifecycle, guiding providers in the selection of the perturbations to test, evaluation methods and benchmarks. These should be listed in a repository designed for temporal adaptability to prevent technological obsolescence. The repository should also enable providers to share experiences and suggest new methods. Such contributions, while

---

[6]See NIST AI Risk Management Framework

[7]See OECD Catalogue

only informative, could complement recommended practices and inform future updates. In this way, robustness can be treated as a context-sensitive requirement, assessed through precise yet evolving standards that ensure compliance with the AI Act.

## References

[1] European Parliament and Council. Artificial intelligence act. *EU's Official Journal*, Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence (Artificial Intelligence Act), 2024.

[2] European Commission. Building trust in human-centric artificial intelligence. *Communication from the commission to the European Parliament, the Council, the European Economic and Social Committee and the Committee of the Regions: COM (2019) 168 final 8.4. 2019*, 2019.

[3] Luciano Floridi, Matthias Holweg, Mariarosaria Taddeo, Javier Amaya, Jakob Mökander, and Yuni Wen. Capai-a procedure for conducting conformity assessment of ai systems in line with the eu artificial intelligence act. *Available at SSRN 4064091*, 2022.

[4] Josep Soler Garrido, Sarah De Nigris, Elias Bassani, Ignacio Sanchez, Tatjana Evas, Antoine-Alexandre André, and Thierry Boulangé. Harmonised standards for the european ai act. Technical Report JRC139430, Joint Research Centre (JRC), European Commission, 2024.

[5] Timo Freiesleben and Thomas Grote. Beyond generalization: a theory of robustness in machine learning. *Synthese*, 202(109), 2023.

[6] CEN/CLC ISO/IEC/TR 24029-1. Artificial intelligence (AI)—assessment of the robustness of neural networks — part 1: Overview (iso/iec tr 24029-1:2021). Technical report, European Committee for Standardisation and European Committee for Electrotechnical Standardisation, 2023.

[7] CEN-CENELEC JTC 21. About the joint technical committee, 2025. URL `https://jtc21.eu/about/`.

[8] Standardisation Request M/593. Commission Implementing Decision of 22 May 2023 on a standardisation request to the European Committee for Standardisation and the European Committee for Electrotechnical Standardisation in support of Union policy on artificial intelligence. Official Journal of the European Union, C/2023/3259, 2023. Available at: `https://ec.europa.eu/transparency/documents-register/detail?ref=C(2023)3215&lang=en`.

[9] Andrea Tocchetti, Lorenzo Corti, Agathe Balayn, Mireia Yurrita, Philip Lippmann, Marco Brambilla, and Jie Yang. AI robustness: a human-centered perspective on technological challenges and opportunities. *ACM Computing Surveys*, 57(6):1–38, 2025.

[10] Ardavan Salehi Nobandegani, Kevin da Silva Castanheira, Timothy O'Donnell, and Thomas R Shultz. On robustness: An undervalued dimension of human rationality. In *CogSci*, page 3327, 2019.

[11] Nathan Drenkow, Numair Sani, Ilya Shpitser, and Mathias Unberath. A systematic review of robustness in deep learning for computer vision: Mind the gap?, 2021.

[12] Emanuele La Malfa, Min Wu, Luca Laurenti, Benjie Wang, Anthony Hartshorn, and Marta Kwiatkowska. Assessing robustness of text classification through maximal safe radius computation. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2949–2968. Association for Computational Linguistics, 2020.

[13] Dan Hendrycks, Steven Basart, Norman Mu, Saurav Kadavath, Frank Wang, Evan Dorundo, Rahul Desai, Tyler Lixuan Zhu, Samyak Parajuli, Mike Guo, Dawn Xiaodong Song, Jacob Steinhardt, and Justin Gilmer. The many faces of robustness: A critical analysis of out-of-distribution generalization. *2021 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages 8320–8329, 2020.

[14] BS EN ISO/IEC 25059. Software engineering — systems and software quality requirements and evaluation (SQuaRE) — quality model for AI system. Technical report, International Organisation for Standardisation (ISO), 2024.

[15] Aurélien Géron. *Hands-on machine learning with Scikit-Learn, Keras, and TensorFlow: Concepts, tools, and techniques to build intelligent systems*. " O'Reilly Media, Inc.", 2022.

[16] Francois Chollet and François Chollet. *Deep learning with Python*. Simon and Schuster, 2021.

[17] Alan Balendran, Céline Beji, Florie Bouvier, Ottavio Khalifa, Theodoros Evgeniou, Philippe Ravaud, and Raphaël Porcher. A scoping review of robustness concepts for machine learning in healthcare. *npj Digital Medicine*, 8(1):38, 2025.

[18] Joaquin Quiñonero-Candela, Masashi Sugiyama, Anton Schwaighofer, and Neil D Lawrence. *Dataset shift in machine learning*. Mit Press, 2022.

[19] Chengyu Dong, Liyuan Liu, and Jingbo Shang. Data quality matters for adversarial training: An empirical study. *arXiv preprint arXiv:2102.07437*, 2021.

[20] Haseeb Javed, Shaker El-Sappagh, and Tamer Abuhmed. Robustness in deep learning models for medical diagnostics: security and adversarial challenges towards robust AI applications. *Artificial Intelligence Review*, 58(12), 2025. doi: 10.1007/s10462-024-11005-9.

[21] Jielin Qiu, Yi Zhu, Xingjian Shi, Florian Wenzel, Zhiqiang Tang, Ding Zhao, Bo Li, and Mu Li. Benchmarking robustness of multimodal image-text models under distribution shift, 2022.

[22] Christopher D. Manning, Prabhakar Raghavan, and Hinrich Schütze. Introduction to information retrieval, 2008.

[23] Zhun Deng, Cynthia Dwork, Jialiang Wang, and Yao Zhao. Architecture selection via the trade-off between accuracy and robustness, 2019.

[24] Dimitris Tsipras, Shibani Santurkar, Logan Engstrom, Alexander Turner, and Aleksander Madry. Robustness may be at odds with accuracy. In *7th International Conference on Learning Representations, ICLR 2019, New Orleans, LA, USA, May 6-9, 2019*. OpenReview.net, 2019.

[25] Firuz Juraev, Mohammed Abuhamad, Simon S Woo, George K Thiruvathukal, and Tamer Abuhmed. Impact of architectural modifications on deep learning adversarial robustness. In *2024 Silicon Valley Cybersecurity Conference (SVCC)*, pages 1–7, 2024. doi: 10.1109/SVCC61185.2024.10637362.

[26] Narmin Ghaffari Laleh, Daniel Truhn, Gregory Patrick Veldhuizen, Tianyu Han, Marko van Treeck, Roman D. Buelow, Rupert Langer, Bastian Dislich, Peter Boor, Volkmar Schulz, and Jakob Nikolas Kather. Adversarial attacks and adversarial robustness in computational pathology. *Nature Communications*, 13(1):5711, 2022.

[27] Sunita Roy, Ranjan Mehera, Rajat Kumar Pal, and Samir Kumar Bandyopadhyay. Hyperparameter optimization for deep neural network models: a comprehensive study on methods and techniques. *Innovations in Systems and Software Engineering*, pages 1–12, 2023.

[28] Christian Arnold, Luka Biedebach, Andreas Küpfer, and Marcel Neunhoeffer. The role of hyperparameters in machine learning models and how to tune them. *Political Science Research and Methods*, 12(4):841–848, 2024.

[29] Houssem Ben Braiek and Foutse Khomh. Machine learning robustness: A primer. In *Trustworthy AI in Medical Imaging*, pages 37–71. Elsevier, 2025.

[30] Qinkai Zheng, Xu Zou, Yuxiao Dong, Yukuo Cen, Da Yin, Jiarong Xu, Yang Yang, and Jie Tang. Graph robustness benchmark: Benchmarking the adversarial robustness of graph machine learning. In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*, 2021.

[31] Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. In *2nd International Conference on Learning Representations, ICLR 2014, Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*, 2013.

[32] Pang Wei Koh, Shiori Sagawa, Henrik Marklund, Sang Michael Xie, Marvin Zhang, Akshay Balsubramani, Weihua Hu, Michihiro Yasunaga, Richard Lanas Phillips, Irena Gao, et al. Wilds: A benchmark of in-the-wild distribution shifts. In *International conference on machine learning*, pages 5637–5664. PMLR, 2021.

[33] Carlos Lassance, Vincent Gripon, Jian Tang, and Antonio Ortega. Structural robustness for deep learning architectures. In *2019 IEEE Data Science Workshop (DSW)*, pages 125–129. IEEE, 2019.

[34] Ian J Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial nets. *Advances in neural information processing systems*, 27, 2014.

[35] Yuchi Tian, Kexin Pei, Suman Jana, and Baishakhi Ray. Deeptest: Automated testing of deep-neural-network-driven autonomous cars. In *Proceedings of the 40th international conference on software engineering*, pages 303–314, 2018.

[36] Zhi Quan Zhou and Liqun Sun. Metamorphic testing of driverless cars. *Communications of the ACM*, 62(3):61–67, 2019.

[37] Weijun Shen, Jun Wan, and Zhenyu Chen. MuNN: Mutation analysis of neural networks. In *2018 IEEE International Conference on Software Quality, Reliability and Security Companion (QRS-C)*, pages 108–115. IEEE, 2018.

[38] Lei Ma, Fuyuan Zhang, Jiyuan Sun, Minhui Xue, Bo Li, Felix Juefei-Xu, Chao Xie, Li Li, Yang Liu, Jianjun Zhao, et al. Deepmutation: Mutation testing of deep learning systems. In *2018 IEEE 29th international symposium on software reliability engineering (ISSRE)*, pages 100–111. IEEE, 2018.

[39] Francesco Croce, Maksym Andriushchenko, Vikash Sehwag, Edoardo Debenedetti, Nicolas Flammarion, Mung Chiang, Prateek Mittal, and Matthias Hein. Robustbench: a standardized adversarial robustness benchmark. In *NeurIPS 2021 Datasets and Benchmarks Track (Round 2)*, 2020.

[40] Greg Brockman, Vicki Cheung, Ludwig Pettersson, Jonas Schneider, John Schulman, Jie Tang, and Wojciech Zaremba. Openai gym, 2016.

[41] Carola Doerr, Hao Wang, Furong Ye, Sander Van Rijn, and Thomas Bäck. Iohprofiler: A benchmarking and profiling tool for iterative optimization heuristics. *CoRR*, abs/1810.05281, 2018.

[42] Hao Wang, Diederick Vermetten, Furong Ye, Carola Doerr, and Thomas Bäck. Iohanalyzer: Detailed performance analyses for iterative optimization heuristics. *ACM Transactions on Evolutionary Learning and Optimization*, 2(1):1–29, 2022.

[43] Seelwan Sathitratanacheewin, Panasun Sunanta, and Krit Pongpirul. Deep learning for automated classification of tuberculosis-related chest x-ray: dataset distribution shift limits diagnostic performance generalizability. *Heliyon*, 6(8), 2020.

[44] Robert Osazuwa Ness, Katie Matton, Hayden Helm, Sheng Zhang, Junaid Bajwa, Carey E Priebe, and Eric Horvitz. Medfuzz: Exploring the robustness of large language models in medical question answering, 2024. Submitted to ICLR 2025.

[45] Hadrien Pouget, Koen Holtman, and Tekla Emborg. Standard setting overview | eu artificial intelligence act, 2025. URL https://artificialintelligenceact.eu/standard-setting-overview/. Updated: 21 July 2025.

[46] Inga Ulnicane, William Knight, Tonii Leach, Bernd Carsten Stahl, and Winter-Gladys Wanjiku. Framing governance for a contested emerging technology: insights from ai policy. *Policy and Society*, 40(2):158–177, 2021.

[47] Sandeep Reddy. Navigating the ai revolution: the case for precise regulation in health care. *Journal of medical Internet research*, 25:e49989, 2023.

[48] Mélanie Gornet. Too broad to handle: can we" fix" harmonised standards on artificial intelligence by focusing on vertical sectors?, 2024. HAL open archive.

[49] Jessica Morley, Caio C.V. Machado, Christopher Burr, Josh Cowls, Indra Joshi, Mariarosaria Taddeo, and Luciano Floridi. The ethics of ai in health care: a mapping review. *Social science & medicine*, 260:113172, 2020.

[50] Tom L Beauchamp and James F Childress. *Principles of biomedical ethics*. Edicoes Loyola, 1994.

[51] Josep Soler Garrido, Delia Fano Yela, Cecilia Panigutti, Henrik Junklewitz, Ronan Hamon, Tatjana Evas, Antoine-Alexandre André, and Salvatore Scalzo. Analysis of the preliminary ai standardisation work plan in support of the ai act. Technical report, Joint Research Centre (JRC), European Commission, 2023.

[52] Ronald Schnitzer, Andreas Hapfelmeier, Sven Gaube, and Sonja Zillner. Ai hazard management: A framework for the systematic management of root causes for ai risks. In *International Conference on Frontiers of Artificial Intelligence, Ethics, and Multidisciplinary Applications*, pages 359–375. Springer, 2023.

[53] Standardisation request M/593 Annexes. Annexes to the Commission Implementing Decision of 22 May 2023 on a standardisation request to the European Committee for Standardisation and the European Committee for Electrotechnical Standardisation in support of Union policy on artificial intelligence. Official Journal of the European Union, C/2023/3259, 2023.

[54] AIME Planning Team. Artificial intelligence measurement and evaluation at the national institute of standards and technology. *National Institute of Standards and Technology*, 2021.

[55] Lena Maier-Hein, Matthias Eisenmann, Annika Reinke, Sinan Onogur, Marko Stankovic, Patrick Scholz, Tal Arbel, Hrvoje Bogunovic, Andrew P Bradley, Aaron Carass, et al. Why rankings of biomedical image analysis competitions should be interpreted with care. *Nature communications*, 9(1):5217, 2018.

[56] Laura Weidinger, Deb Raji, Hanna Wallach, Margaret Mitchell, Angelina Wang, Olawale Salaudeen, Rishi Bommasani, Sayash Kapoor, Deep Ganguli, Sanmi Koyejo, et al. Toward an evaluation science for generative ai systems, 2025.

[57] Tambiama Madiega and Anne Louise Van De Pol. Artificial intelligence act and regulatory sandboxes. *European Parliamentary Research Service*, 6, 2022.

[58] Stig A. Due, Hira Shah, Thiago Moraes, Nathan Genicot, and Martin Canter. Sandboxing artificial intelligence: Balancing innovation, regulation, and stakeholder needs. https://www.fari.brussels/research-and-innovation/publication/sandboxing-artificial-intelligence, 2025. FARI Institute.