



**HAL**  
open science

## Towards Trustworthy-AI-by-Design Methodology for Intelligent Radiology Systems

Clotilde Brayé, Jérémy Clech, Arnaud Gotlieb, Nadjib Lazaar, Patrick Malléa

► **To cite this version:**

Clotilde Brayé, Jérémy Clech, Arnaud Gotlieb, Nadjib Lazaar, Patrick Malléa. Towards Trustworthy-AI-by-Design Methodology for Intelligent Radiology Systems. Journée Santé et IA @PFIA\_2023, Plate-Forme Intelligence Artificielle - PFIA, Jul 2023, Strasbourg, France. lirmm-04160772

**HAL Id: lirmm-04160772**

**<https://hal-lirmm.ccsd.cnrs.fr/lirmm-04160772>**

Submitted on 12 Jul 2023

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

# Towards Trustworthy-AI-by-Design Methodology for Intelligent Radiology Systems

Clotilde Brayé<sup>1,2,3</sup>, Jérémy Clech<sup>1</sup>, Arnaud Gotlieb<sup>2</sup>, Nadjib Lazaar<sup>3</sup>, Patrick Malléa<sup>1</sup>

<sup>1</sup> NEHS DIGITAL, France

<sup>2</sup> Simula Research Laboratory, Norway

<sup>3</sup> Université de Montpellier, CNRS, LIRMM, France

July 12, 2023

## Résumé

*Cet article présente nos travaux en cours sur l'élaboration d'une méthodologie généralisable pour développer des systèmes radiologiques intelligents et fiables. Notre méthodologie, appelée "trustworthy-AI-by-design", vise à comprendre comment développer des systèmes d'IA pour la radiologie qui répondent aux futures exigences européennes en matière d'IA digne de confiance. Dans cet article, nous deployons et testons la méthodologie "trustworthy-AI-by-design" sur deux cas d'utilisation distincts : 1) la détection de la maladie COVID-19 sur des tomographies thoraciques en utilisant de l'apprentissage profond et 2) un système d'alerte et de prédiction de non-présentation des patients à leur rendez-vous radiologique qui s'appuie sur des modèles d'apprentissage automatique. Nos premières analyses montrent que les différents critères éthiques d'une IA digne de confiance nécessitent des méthodes de conception distinctes et doivent être pris en compte tout au long du cycle de développement de l'IA.*

## Mots-clés

*Trustworthy AI, Ethical AI, Intelligent Radiology Systems*

## Abstract

*This paper discusses our work-in-progress to develop a comprehensive and unified methodology for creating smart and reliable intelligent radiology systems (IRS) that meet the upcoming European requirements on trustworthy AI. Specifically, we present our so-called trustworthy-AI-by-design methodology by showcasing two distinct use cases of IRS: 1) Deep Learning-based COVID detection in Chest Tomography scans and 2) AI-supported radiological no-shows alerting systems. Our initial analysis highlights that different key requirements of trustworthy AI necessitate distinct design methods and issues that must be addressed throughout the entire AI development cycle, from initial conception to final design.*

## 1 Introduction

AI is a breakthrough in healthcare, particularly in radiology. It brings new features and automation that support diagnosis and improve patient care. However, the use of AI

in radiology must be carefully managed due to the risks it poses, which can even be life-threatening for patients. The European Union is taking steps to establish an ethical and legal framework for AI, with trustworthy AI as its guiding principle. The European Commission has introduced ethical guidelines [1] and a law proposal "The AI Act" [2] to regulate the development and use of AI.

The forthcoming European legislation on AI will significantly alter the design, development, and evaluation of AI systems which possess ethical risks. All AI stakeholders will have to comply with the ethical and legal constraints set by the European regulations. Therefore, in Radiology, it is vital to grasp the construction of smart and reliable Intelligent Radiology Systems (IRS). In this paper, we present an initial release version of our structured and comprehensive Trustworthy-AI-by-Design (TAID) methodology for IRS. Our objective is to evaluate how to exploit TAID to design AI-based IRS which complies with the upcoming European regulation in a systematic manner. To evaluate our TAID methodology, we developed two distinct use cases: 1) A deep learning-based COVID detection system in Chest Tomography scans, by using an existing and released European database named 'FIDAC' of images [3] and 2) an AI-supported no-shows alerting and predicting system for radiology appointments named 'NOSHOW'.

The remainder of the paper is structured as follows: In Sec.2, we provide the necessary background required to understand the paper. Sec.3 outlines the initial principles and steps of our TAID methodology. Sec.4 presents the two use cases and their initial results. In Sec.5, we discuss the current limits of our proposed methodology. Finally, Sec.6 concludes this article by summarizing the main findings of our work and identifying further work.

## 2 Background

This section gives an overview of the current state of the European and national AI regulation and knowledge about trustworthy AI in radiology. It also introduces the principles of usual risk management and assessment procedures.

## 2.1 European and National Regulations on Artificial Intelligence

The ethical guidelines for trustworthy AI [1] as provided by the High-Level Expert Group (HLEG) mandated by the EC includes four main principles, namely a) respect of human autonomy; b) prevention of harm; c) fairness; and d) explicability. These four ethical principles are then translated into seven trustworthy AI requirements which are: (TAIR<sub>1</sub>) Human agency and oversight; (TAIR<sub>2</sub>) technical robustness and safety; (TAIR<sub>3</sub>) privacy and data governance; (TAIR<sub>4</sub>) transparency; (TAIR<sub>5</sub>) diversity, non-discrimination and fairness; (TAIR<sub>6</sub>) societal and environmental well-being; (TAIR<sub>7</sub>) accountability.

Besides the ethical guidelines, the European Parliament has designed an upcoming legal framework proposition named "The AI Act" [2] which is based on a four-category-risk classification for any AI-based system. Each risk category has its own requirements proportional to the level of risk, from the highest to the lowest level: (i) Unacceptable risk; (ii) High-risk AI systems; (iii) Limited risk; (iv) Low and minimal risk [4].

Amongst the seven key requirements of Trustworthy AI, human agency and oversight is critical in healthcare and medical radiology. Although AI solutions demonstrate an impressive breadth of applications and excellent performance, it is crucial for humans to maintain control over the decision-making process. This requirement aligns with the accountability, non-discrimination and fairness, data governance and privacy and societal well-being requirements for medical diagnostics and decisions, in order to guarantee the fundamental principle of equal access to healthcare to everyone in the population. Additionally, transparency is essential to explain how AI solutions reach specific conclusions and provide means for verifying automated system results in radiology. Lastly, technical robustness and safety are crucial requirements to evaluate the absence of risk in using AI-based devices. The CE marking of medical devices already applies this requirement.

In 2022, the French Ministry of Solidarity and Health released recommendations for implementing ethics-by-design for AI-based healthcare solutions [5]. These recommendations align with the European key requirements for trustworthy AI, including human agency and oversight, technical robustness, transparency, and more. Additionally, the ministry advises performing a risk management assessment based on ISO standard 14971 [6] for the application of risk management to medical devices.

## 2.2 Artificial Intelligence in Radiology

In Radiology, AI has been prominent for over 20 years with numerous AI-based solutions available on the market [7]. These include automatic measures such as breast density level using computer vision and logistic regression models [8], auto segmentation for 3D complex tumours based on image recognition and support vector machines models [9], and computer-aided detection for breast cancer using image recognition and random forests / deep learning models [10].

Recent advances in AI have also enabled its use in personalized medicine, such as personalized prognosis and treatment planning using radiogenomics and deep learning [11], and in patient medical context, such as summarizing medical events using natural language processing of medical reports [12].

The widespread adoption of AI-based solutions in daily clinical practice remains limited due to issues related to awareness and concerns about trustworthiness. Radiologists do not trust the results because the tools do not explain well, responsibility is diluted, and accuracy cannot be guaranteed. Developing trustworthy AI-based radiology information systems can help overcome these barriers and facilitate AI-based solutions' integration in radiology practice.

## 2.3 Risk Management for Medical Devices

The ISO standard 14971 [6] provides a risk management process for medical devices. The objective of this standard is to assist manufacturers in preventing or minimizing dangerous situations for patients or users of the medical devices under development. A multidisciplinary team is required to conduct each risk management activity. The life cycle of the medical device must adhere to the following steps:

**1. Risk Analysis:** Identify potential hazardous situations in which the medical device may be involved, including its intended use and reasonably foreseeable misuse, as well as identifying hazards and hazardous situations that may arise, and their associated risks;

**2. Risk Evaluation:** Risk assessment is conducted for each hazardous situation by a multidisciplinary team according to predefined risk acceptability criteria. The team estimates the plausibility (formerly known as probability of occurrence in ISO 14971) and severity of the potential damage. The risk is then quantified as a combination of plausibility and severity;

**3. Risk Control:** For each risk that is deemed unacceptable, the team must identify and implement risk control measures. Once this is done, the residual risk - i.e., the risk that remains after the implementation of the control measures - must be evaluated, along with the effectiveness of the control measures in reducing risk. This process should be repeated until the residual risks are deemed acceptable.

## 3 The TAID Methodology

TAID is a three-steps methodology:

1. Health Purpose Definition (HPD);
2. Iterative Risk Mgt. on Trustworthy AI (IRM-TAI):
  - (a) Identification of risks associated with patients and IRS users;
  - (b) Mapping identified risks to the Trustworthy AI regulation's key requirements;
  - (c) Evaluating the criticality of Trustworthy AI requirements and implementing appropriate mitigation measures to reduce criticality as acceptable;

### 3. Testing and Validation of Trustworthiness (TVT).

TAID includes an iterative risk management process following the principles of ISO 14971. However, it differs slightly from ISO 14971 as it is primarily focused on the seven key requirements of Trustworthy AI in the European regulation.

#### 3.1 Health Purpose Definition (HPD)

The first TAID step is to define the application scope of the system, including its main users such as patients, radiologists, generalist doctors, device assistants, etc., as well as the IRS access policy. This step is crucial in guiding the overall assessment process of the IRS system's use of AI. The outcome of this step is a set of recommendations and documentation on how to use and exploit the AI-based IRS system, with a specific focus on ensuring conformity with European and national AI regulations, and identifying the most critical key requirements for Trustworthy AI.

#### 3.2 Iterative Risk Management on Trustworthy AI (IRM-TAI)

As the second TAID step, we perform a detailed risk identification process based on the plausibility<sup>1</sup> and severity of each identified risk.

Let  $R = \{r_1, \dots, r_n\}$  be the set of identified risks, and  $V_p(r_i)$  the plausibility of risk  $r_i$  and  $V_s(r_i)$  be the severity of  $r_i$ . Both  $V_p$  and  $V_s$  can take only four possible values: *Residual* (1), *Limited* (2), *Significant* (3), and *Maximum* (4). For each risk  $r_i$ , assigning the values  $V_p(r_i)$  and  $V_s(r_i)$  is realized through a "College of Experts" (CoE) based their experience of previous AI-development projects. The CoE brings together a multidisciplinary team consisting, among others, of AI and medical experts. Then, we calculate the *quantification of each risk*, denoted as  $Q(r_i)$ , by using the formula from ISO 14971:

$$Q(r_i) = V_p(r_i) \times V_s(r_i)^2 \quad (\text{Risk Quantification})$$

It is worth noticing that each risk can be associated with one or more trustworthy AI key requirements  $\text{TAIR}_i$ . So we define the *criticality of a risk*  $r_i$ , denoted  $C(r_i)$  as its weighted quantification by the number of requirements it is associated with:

$$C(r_i) = \frac{k_i}{k} * Q(r_i) \quad (\text{Risk Criticality})$$

where  $k_i$  is the number of trustworthy AI requirements associated with risk  $r_i$  and  $k$  is the total number of trustworthy AI requirements.

After the criticality of each risk has been computed, the next TAID step is to evaluate if the criticality is acceptable, based on the two following conditions<sup>2</sup>:

$$\begin{cases} \forall r_i \in R, Q(r_i) < \theta & (\text{cond 1}) \\ \frac{1}{n} * \sum_{i=1}^n C(r_i) \leq \eta & (\text{cond 2}) \end{cases}$$

<sup>1</sup>ISO standard 14971 defines the plausibility of a risk as the probability of its occurrence and discretizes it through different levels

<sup>2</sup>These conditions correspond to an interpretation of the ISO standard 14971 which is not concerned with Trustworthy AI

where both  $\theta$  and  $\eta$  are arbitrarily selected thresholds that are determined by the CoE. Note that (cond 1) is concerned with each individual risk while (cond 2) is related to all risks. If any of these two conditions is violated, then mitigation actions must be taken to reduce either the risk plausibility, or its severity, or both. Note also that mitigation actions can be taken to reduce the number of trustworthy AI requirements associated to a risk  $r_i$ . This can reduce the criticality of  $r_i$  without modifying its quantification. From there, the quantification and criticality are recomputed and the conditions are evaluated again. This iterative process continues until the two conditions are respected, and we can proceed to the third step of the methodology.

#### 3.3 Testing and Validation of Trustworthiness (TVT)

The third TAID step involves testing the trustworthiness of the AI-based IRS. For that purpose, different types of tests are involved. Roughly speaking, we distinguish three different test campaigns:

- 1. Unit and Integration Tests (UIT):** These tests validate each IRS component individually for the seven key requirements  $\text{TAIR}_i$  and test the integration of the components. By defining expected behavior and Trustworthy AI requirement properties in test scripts, some tests can be automated.
- 2. System Tests (SYT):** These tests are more complex as they require human operators to execute detailed test plans that activate all parts of the IRS. With system tests, it is possible to verify high-level properties through advanced scenarios such as the fairness of the IRS when providing access to patient healthcare, scenarios where patient data privacy breaches are simulated, or scenarios where opaque decisions are made without providing enough transparency.
- 3. Acceptance Tests (ACT):** These tests demonstrate to a third-party authority, for certification or audit purposes, that the seven key requirements  $\text{TAIR}_i$  have been properly evaluated and validated. The failure of any of these tests leads to a revision of the IRS implementation and re-entry into the development cycle to improve the IRS with respect to Trustworthy AI.

## 4 Evaluation

This section presents the results of applying TAID to two distinct use cases. We selected these use cases because they involve vastly different processes and can serve as a useful basis for identifying the limitations of our proposed methodology. We begin by presenting the use cases and subsequently present our findings and results from applying the TAID methodology.

### 4.1 Use Cases

This section provides an overview of the two use cases selected for our study. The first use case involves radiology image processing and serves as an illustrative example of image processing systems. The second use case pertains to patient pathways and access to radiology.

**The 'FIDAC' Use Case.** 'FIDAC' aims to automatically detect COVID-19<sup>3</sup> on CT-scans using Deep Learning (DL). To accomplish this, a large database called *French Images Database Against Coronavirus (FIDAC)* [3] is used to train DL models that can classify and segment radiology images. These models are intended to be integrated into IRS to support the medical decisions to orient patients towards different treatments. The detailed usage of these models has been extensively described in the literature [13].

**The 'NOSHOW' Use Case.** In 2023, the National Academy of Medicine issued a press release designating patient no-show as a public health issue [14]. 'NOSHOW' aims to use ML to estimate the likelihood of a patient attending a radiology appointment. The goal is to implement mitigation measures based on the predictions generated by the AI system, such as sending repeated messages to the patient, making phone calls, or overbooking. AI-based solutions are believed to optimize patient access to radiology appointment[15].

## 4.2 Results of Applying TAID

Both use cases present different levels of risk and ethical implications. 'FIDAC' is considered a high-risk system according to the AI Act terminology (cf. Sec.2), as it involves the diagnosis of patients, which directly impacts their survival. On the other hand, 'NOSHOW' is classified as a low-risk system. Disruptions on access to the healthcare system can have serious consequences for patients, but not as critical as diagnosis failure. In our evaluation we did not yet consider the safety aspect of the AI system. We will consider it in a latter phase as safety is mainly related to patient information flows and we do not yet have access to the technical specifications of these flows. TAID was applied to both use cases. Tab.1 results from the first iteration of the IRM-TAI step<sup>4</sup>, forming the initial risk matrix, while Tab.1 results from the last iteration of the IRM-TAI step, forming the residual risk matrix. Tab.1 summarizes the identified risks, along with their trustworthy AI requirements, their initial and residual criticality. In both cases, the  $\theta$  threshold was arbitrarily fixed to 27 and  $\eta$  to 5. For the sake of comprehension, we detail the result of TAID on two selected risks, namely 1) F\_R1: Personal data breaches on 'FIDAC'; and 2) N\_R2 Lack of explicability of the prediction on 'NOSHOW'.

**Reducing F\_R1 Criticality.** 'FIDAC' involves sensitive data (i.e., patient information) such as gender, age, diagnosis, CT-scan. Then, F\_R1 has a significant plausibility and a maximum severity. According to the risk quantification equation (in Sec.3.2), we have:  $Q(F\_R1) = 48$ . We associate this risk to TAIR<sub>1</sub>, TAIR<sub>2</sub>, TAIR<sub>3</sub>, then its initial criticality is  $C(F\_R1) = 13.71$ .

To address the risk of personal data breaches, we used

<sup>3</sup>COVID-19 is a severe acute respiratory disease that can cause lung infection

<sup>4</sup>We did not consider available AI risk ontologies yet because the goal was to prove the effectiveness of the TAID methodology. However, we plan to use it based on the state of the art [16]

an anonymous database which was created by using the privacy-by-design methodology given in [17]. This methodology is based on GDPR principles and it reduces the evaluation of the risk to a residual plausibility and a residual severity. Then, the residual quantification is  $Q(F\_R1) = 1$  and the residual criticality is  $C(F\_R1) = 0.28$ .

**Reducing N\_R2 Criticality.** 'NOSHOW' can help the medical staff in radiology appointment management. As such, the staff shall take an informed decision regarding maintaining or withdrawing the appointment; the AI system shall be robust and transparent, not discriminatory and fair. Additionally, the decision process shall be fully traceable. Hence, we initially evaluated the lack of explicability (i.e., N\_R2) as having a significant plausibility and a significant severity. N\_R2 quantification is then  $Q(N\_R2) = 27$ . The CoE associate this risk to TAIR<sub>1</sub>, TAIR<sub>2</sub>, TAIR<sub>4</sub>, TAIR<sub>5</sub>, TAIR<sub>7</sub>. The initial criticality is then equal to  $C(N\_R2) = 19.28$  (according to the Risk Criticality eq. in Sec.3.2).

We reduce this risk by using only white-box ML models such as Random Forest and XGBoost, instead of opaque models based on DL. The results of white-box models are easier to explain because their computations are based on decision trees. The literature shows that these white-box ML models offer us acceptable performance [18]. In addition, we plan to use XAI libraries such as `kernelSHAP` or `LIME` to increase the level of explanations by evaluating the impact of each feature on the final classification. The XAI libraries will bring more information for the user and will be helpful to detect any possible discrimination [19]. Thanks to these actions, Risk N\_R2 is reduced to residual plausibility and significant severity, which makes its residual quantification equal to  $Q(N\_R2) = 9$ . The trustworthy AI requirements associated are TAIR<sub>4</sub> and TAIR<sub>7</sub>, then the residual criticality equals to:  $C(N\_R2) = 2.57$ . We proceed in the same way for each risk, resulting in the residual risk matrix given in Fig.1 and the residual criticality values shown in Tab.1.

## 5 Discussion and Limits of TAID

While TAID offers a comprehensive and effective framework for managing risks associated with AI systems, it does have some limitations and challenges that should be considered. Here are a few potential limitations to keep in mind:

**1. Threshold selection:** As mentioned earlier, TAID recommends decreasing the criticality of Trustworthy AI requirements up to a certain minimum threshold that is arbitrarily selected. This means that the threshold may not necessarily be compliant with any legal or regulatory requirements in place. Additionally, the chosen threshold may not be appropriate for all types of AI systems and use cases.

**2. Generalization:** The risk mitigation actions that are specific to each use case may not necessarily be generalizable to new AI systems or applications. Each AI system has unique characteristics, datasets, and challenges, and it may be difficult to determine which risk mitigation actions will

UC	Risks	Risk Description	TAIR <sub>i</sub> Initial	TAIR <sub>i</sub> Residual	Initial Criticality	Residual Criticality
FIDAC	F_R1	Personal data breaches	TAIR <sub>2</sub> , TAIR <sub>3</sub>	TAIR <sub>3</sub>	13.71	0.14
FIDAC	F_R2	Lack of explicability of the prediction	TAIR <sub>1</sub> , TAIR <sub>4</sub> , TAIR <sub>7</sub>	TAIR <sub>4</sub> , TAIR <sub>7</sub>	11.57	5.14
FIDAC	F_R3	Model attacks	TAIR <sub>2</sub>	TAIR <sub>2</sub>	2.57	2.57
FIDAC	F_R4	Wrong patient care	TAIR <sub>1</sub> , TAIR <sub>2</sub> , TAIR <sub>6</sub>	TAIR <sub>1</sub> , TAIR <sub>2</sub> , TAIR <sub>6</sub>	20.57	7.71
FIDAC	F_R5	Differences of performance depending on age or gender	TAIR <sub>2</sub> , TAIR <sub>3</sub> , TAIR <sub>5</sub>	TAIR <sub>2</sub> , TAIR <sub>3</sub> , TAIR <sub>5</sub>	7.71	3.86
NOSHOW	N_R1	Personal data breaches	TAIR <sub>2</sub> , TAIR <sub>3</sub>	TAIR <sub>3</sub>	7.71	0.14
NOSHOW	N_R2	Lack of explicability of the prediction	TAIR <sub>1</sub> , TAIR <sub>2</sub> , TAIR <sub>4</sub> , TAIR <sub>5</sub> , TAIR <sub>7</sub>	TAIR <sub>4</sub> , TAIR <sub>7</sub>	19.28	2.57
NOSHOW	N_R3	Model attacks	TAIR <sub>2</sub>	TAIR <sub>2</sub>	2.57	2.57
NOSHOW	N_R4	Patient categorisation	TAIR <sub>1</sub> , TAIR <sub>4</sub> , TAIR <sub>5</sub>	TAIR <sub>1</sub> , TAIR <sub>4</sub> , TAIR <sub>5</sub>	20.57	7.71
NOSHOW	N_R5	Excessive patient reminders	TAIR <sub>1</sub> , TAIR <sub>2</sub> , TAIR <sub>3</sub> , TAIR <sub>4</sub> , TAIR <sub>5</sub>	TAIR <sub>2</sub> , TAIR <sub>4</sub> , TAIR <sub>5</sub>	8.57	3.42
NOSHOW	N_R6	Disorganization of the center	TAIR <sub>1</sub> , TAIR <sub>2</sub>	TAIR <sub>2</sub>	5.14	1.28
NOSHOW	N_R7	Deterioration of the facility's image	TAIR <sub>2</sub> , TAIR <sub>7</sub>	TAIR <sub>2</sub> , TAIR <sub>7</sub>	2.57	1.14
NOSHOW	N_R8	Inability of the facility to complete the planned medical exam	TAIR <sub>1</sub>	TAIR <sub>1</sub>	6.86	0.14
NOSHOW	N_R9	Equal access to healthcare	TAIR <sub>1</sub>	TAIR <sub>1</sub>	6.86	0.14

Table 1: Risks Associated to the Seven Trustworthy-AI Key Requirements and the Criticality for 'FIDAC' and 'NOSHOW'

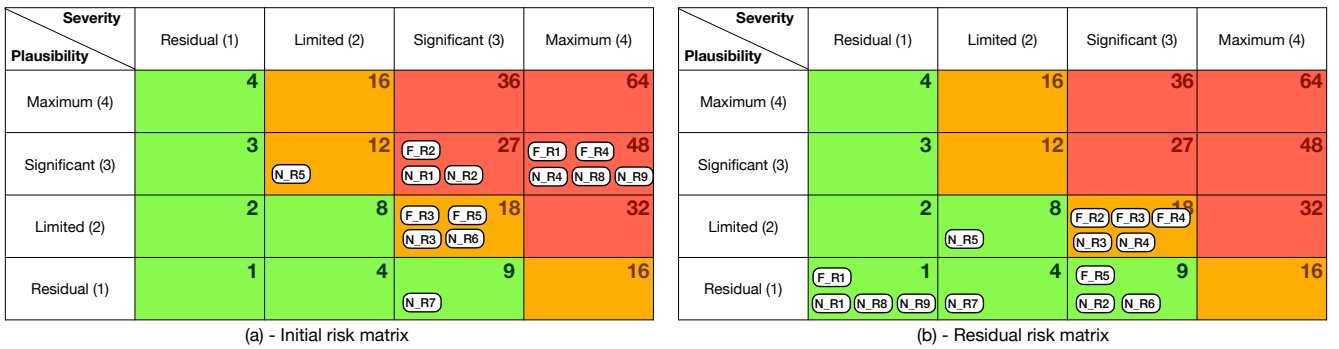


Figure 1: Risk Matrices for the 'FIDAC' and 'NOSHOW' Use Cases.

be most effective for each new system.

**3. Limited impact:** While TAID can help to reduce the criticality of Trustworthy AI key requirements, it may not eliminate risks completely. There may still be residual risks associated with the use of AI systems, even after mitigation actions have been taken.

**4. Explainability challenges:** Although TAID recommends using explainable AI techniques to improve the transparency of AI systems, there is still ongoing research into how best to interpret and understand the output of opaque models. This means that there may be limitations to how effectively AI systems can be made transparent and understandable, even with the use of explainability methods. Overall, while TAID offers a comprehensive framework for managing AI-related risks, its limitations and challenges must be understood. Continued research and development is needed to ensure that AI systems can be used responsibly while minimizing risks to individuals and society.

## 6 Conclusion and Future Work

In this paper, we introduced TAID, a methodology that addresses all the seven Trustworthy AI requirements during the design phase, enabling the development of trustworthy-by-design AI systems for radiology. We acknowledge that the risk mitigation actions taken are specific to each use case and may not be generalized to other systems. However, it is reassuring to note that the assessment of all seven requirements is similar for both use cases (e.g., mitigation actions are similar for both use cases regarding level of autonomy and bias identification). Our future work includes formulating the iterative risk management procedure as a multi-criteria minimization problem, aiming to minimize the risk for each requirement as part of a larger problem. This will involve finding the right balance between risk reduction for each requirement and model performance.

## References

- [1] High-Level Expert Group on Artificial Intelligence. *Ethics guidelines for trustworthy AI*. Publications Office of the European Union, 2019.
- [2] European Commission. *Proposal for a Regulation Laying down Harmonised Rules on Artificial Intelligence (AI Act)*. 2021. [https://eur-lex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75ed71a1.0001.02/DOC\\_1&format=PDF](https://eur-lex.europa.eu/resource.html?uri=cellar:e0649735-a372-11eb-9585-01aa75ed71a1.0001.02/DOC_1&format=PDF).
- [3] Loic Bussel, Jean-Michel Bartoli, Samy Adnane, Jean-François Meder, Patrick Malléa, Jeremy Clech, Marc Zins, and Jean-Paul Bérégi. French imaging database against coronavirus (FIDAC): A large COVID-19 multi-center chest CT database. *Diagnostic and Interventional Imaging*, 103(10):460–463, 2022.
- [4] Tambiama André Madiega. Artificial intelligence act. *European Parliament: European Parliamentary Research Service*, 2021.
- [5] GT3. *Recommandations de bonnes pratiques pour intégrer l'éthique dès le développement des solutions d'Intelligence Artificielle en Santé : mise en œuvre de « l'éthique by design »*. 2022. [https://esante.gouv.fr/sites/default/files/media\\_entity/documents/ethic\\_by\\_design\\_guide\\_vf.pdf](https://esante.gouv.fr/sites/default/files/media_entity/documents/ethic_by_design_guide_vf.pdf).
- [6] International Organization for Standardization. *Medical devices — Application of risk management to medical devices*. ISO standard no. 14971:2019 edition. <https://www.iso.org/obp/ui/#iso:std:iso:14971:ed-3:v1:en>.
- [7] Jean-Philippe Masson. *L'imagerie médicale en France: un atout pour la santé, un atout pour l'économie*. ISBN 978-2-9558316-0-1, 2016.
- [8] Amy T Wang, Celine M Vachon, Kathleen R Brandt, and Karthik Ghosh. Breast density and breast cancer risk: a practical review. *Elsevier*, 89(4):548–557, 2014.
- [9] Kai Roman Laukamp, Frank Thiele, Georgy Shakirin, David Zopfs, Andrea Faymonville, Marco Timmer, David Maintz, Michael Perkuhn, and Jan Borggrefe. Fully automated detection and segmentation of meningiomas using deep learning on routine multiparametric MRI. 29(1):124–132, 2019.
- [10] Serena Pacilè, January Lopez, Pauline Chone, Thomas Bertinotti, Jean Marie Grouin, and Pierre Fillard. Improving breast cancer detection accuracy of mammography with the concurrent use of an artificial intelligence tool. *Radiological Society of North America*, 2(6):e190208, 2020. <https://pubs.rsna.org/doi/full/10.1148/ryai.2020190208>.
- [11] Sanjay Saxena, Biswajit Jena, Neha Gupta, Suchismita Das, Deepaneeta Sarmah, Pallab Bhattacharya, Tanmay Nath, Sudip Paul, Mostafa M. Fouda, Manudeep Kalra, Luca Saba, Gyan Pareek, and Jasjit S. Suri. Role of artificial intelligence in radiogenomics for cancers in the era of precision medicine. *Cancers*, 14(12):2860, 2022. Number: 12 Publisher: Multidisciplinary Digital Publishing Institute.
- [12] Rimma Pivovarov and Noémie Elhadad. Automated methods for the summarization of electronic health records. *Journal of the American Medical Informatics Association*, 22(5):938–947, 04 2015.
- [13] Sakshi Ahuja, Bijaya Ketan Panigrahi, Nilanjan Dey, Venkatesan Rajinikanth, and Tapan Kumar Gandhi. Deep transfer learning-based automated detection of COVID-19 from lung CT scan slices. *Appl. Intell.*, 51(1):571–585, 2021.
- [14] Académie nationale de médecine and Conseil national de l'Ordre des Médecins. Rendez-vous non honorés, communiqué commun, 2023. <https://www.academie-medecine.fr/wp-content/uploads/2023/01/23.1.27-Communique-RV-non-honores.pdf>.
- [15] Le Roy Chong, Koh Tzan Tsai, Lee Lian Lee, Seck Guan Foo, and Piek Chim Chang. Artificial intelligence predictive analytics in the management of outpatient MRI appointment no-shows. *American Journal of Roentgenology*, 215(5):1155–1162, 2020.
- [16] Delaram Golpayegani, Harshvardhan J. Pandit, and Dave Lewis. AIRO: an ontology for representing AI risks based on the proposed EU AI act and ISO risk management standards. In Anastasia Dimou, Sebastian Neumaier, Tassilo Pellegrini, and Sahar Vahdati, editors, *Towards a Knowledge-Aware AI - SEMANTiCS 2022 - Proceedings of the 18th International Conference on Semantic Systems, 13-15 September 2022, Vienna, Austria*.
- [17] Jérémy Clech, Arnaud Gotlieb, Florence Sève, Frédérique Didout, and Patrick Malléa. Méthodologie d'anonymisation dès la conception d'un jeu de données en imagerie médicale. In *Conférence Nationale d'Intelligence Artificielle Année 2022*, page 25, 2022.
- [18] Luiz Henrique Américo Salazar, Wemerson Delcicio Parreira, Anita Maria da Rocha Fernandes, and Valderi Reis Quietinho Leithardt. No-show in medical appointments with machine learning techniques: A systematic literature review. *Information*, 13(11), 2022.
- [19] Aditya Jain, Manish Ravula, and Joydeep Ghosh. Biased models have biased explanations. *arXiv preprint arXiv:2012.10986*, 2020.