

# Interdisciplinary Remarks on the Assessment List on Trustworthy AI (ALTAI) Applied to a P5 Medicine Tool

Denise Amram<sup>1\*</sup>, Arianna Cignoni<sup>2</sup>, Tommaso Banfi<sup>2</sup>, Gastone Ciuti<sup>2</sup>

<sup>1</sup>Lider-Lab, Dirpolis Institute, Scuola Superiore Sant'Anna, Pisa, Italy; <sup>2</sup>Department of Excellence in Robotics & AI, Scuola Superiore Sant'Anna, Pisa, Italy

## ABSTRACT

The authors apply the Assessment List on Trustworthy AI (ALTAI) to a possible AI-based tool aiming at supporting the melanoma cancer diagnosis. They take the opportunity to provide an interdisciplinary analysis of the proposed self-assessment tool in light of its possible mandatory application in R&D&I. The presented empirical exercise highlights some pros and cons of the adopted checklist, stimulating further remarks on the EU regulatory initiatives on AI.

Finally, we try to understand the improvement of procedures, medical knowledge and treatment collected and improved during these months, that allowed for a lower mortality rate in the referring period.

**Keywords:** Artificial intelligence; Assessment; ALTAI checklist; Healthcare sector

## INTRODUCTION

In the last months, the EU accelerated the legislative process on Artificial Intelligence (AI). In July 2020, the High-Level Expert Group on Artificial Intelligence (AI HLEG) presented the final Assessment List for Trustworthy Artificial Intelligence (ALTAI). While, in October 2020, the European Parliament adopted a "Provisional text on the Framework of ethical aspects of artificial intelligence, robotics and related technologies", where the risk-based approach is confirmed as main strategy for the further AI legislative initiatives that the EU Commission is working on (hereinafter "Provisional Resolution") [1].

## LITERATURE REVIEW

From this perspective, a deep analysis of the ALTAI methodology and structure becomes crucial to address in a responsible and proactive way the current compliance challenges for those who develop AI-based systems. ALTAI consists of a series of questions that may steer the AI designers (rectius the AI-controller, or developers according to the mentioned Provisional Resolution) towards a multidisciplinary evaluation path aimed at addressing the seven ethical-legal-safety compliance challenges emerging by the Guidelines on Trustworthy AI, adopted by the EU Commission in April 2019 [2,3].

According to them, an AI system becomes trustworthy whereas it is lawful (i.e., compliant with the applicable legal framework), ethical (i.e., compliant with the applicable ethical framework) and robust (i.e., compliant with the applicable safety standards). The interplay

between these three pillars is determined by the following seven grounds [4,5].

### Human agency and oversight

It includes both the ethical and the legal dimension as it refers to fundamental rights protection aimed at maintaining the balance between human control and technical progress in terms of human agency and oversight. Human beings shall be protected both as individuals and groups, taking into account inclusiveness, fairness, non-discrimination and vulnerabilities protection as paramount interests.

### Technical robustness and safety

It refers to the system resilience to attacks and security, including fall back plans and the compliance with the highest levels of general safety, accuracy, reliability, and reproducibility.

### Privacy and data governance

This profile establishes a bridge with the most effective compliance process by design and by default introduced for personal data processing by the EU Reg. n. 2016/679 on General Data Protection Regulation (GDPR) aiming at guaranteeing the respect for confidentiality, quality, and integrity of data.

### Transparency

This is a principle established to guarantee the traceability, explainability, communication of methods, goals, and results of the given AI system.

**Correspondence to:** Denise Amram, Lider-Lab, Dirpolis Institute, Scuola Superiore Sant'Anna, Pisa, Italy, E-mail: denise.amram@santannapisa.it

**Received:** May 13, 2021; **Accepted:** May 27, 2021; **Published:** June 3, 2021

**Citation:** Amram D, Cignoni A, Banfi T, Ciuti G (2021) Interdisciplinary Remarks on the Assessment List on Trustworthy AI (ALTAI) Applied to a P5 Medicine Tool. *J Clin Res Bioeth.* 12:374.

**Copyright:** © 2021 Amram D, et al. This is an open-access article distributed under the terms of the Creative Commons Attribution License, which permits unrestricted use, distribution, and reproduction in any medium, provided the original author and source are credited.

## Diversity, non-discrimination, and fairness

This ground refers to the interdisciplinary safeguards to be implemented in order to avoid a misuse or an unfair use of AI in terms of bias, accessibility, and universal design.

## Societal and environmental well-being

The AI system shall be put in the market as a sustainable solution under the environmental, social, and societal perspectives, considering the democratic values rooted within the EU framework.

## Accountability

This is the main principle that enables the compliance process in terms of proactively responsibilities allocation through a risk-based approach that includes auditability, minimization and reporting of negative impact, trade-offs, and redress.

In order to facilitate to undertake the ALTAI checklist, a web-based tool that provides results in terms of level of compliance and list of recommendations has been developed. This could be followed by anyone who is responsible to design and develop an AI system to confidentially perform a self-assessment. This preparatory tool can be interpreted as a pilot for possible future obligations on the topic, as the risk-based approach addressed through an impact assessment seems to be suggested by the above-mentioned EU Parliament Resolution as well.

In the following paragraphs, we will assess the ALTAI checklist on possible design of AI-tools applied to healthcare, considering the paramount role that the health-data debate has among the European strategy for data and the related investments to encourage the upscaling of cross-border exchange of health data and their reuse to improve the detection, diagnosis, and treatment of diseases [6].

In particular, we will provide some methodological remarks emerging from the application of the ALTAI to an AI-based tool developed through a Predictive, Preventive, Personalised, and Participatory (P4) cancer medicine approach, in order to contribute to the debate from a bottom up and interdisciplinary perspective. In fact, within this empirical assessment, our efforts focused on both explaining the tool design from a scientific viewpoint, covering engineering, ethical-legal, and medical aspects and on introducing technical and organizational enablers to allow the real participation of unexpert users in the prevention of specific pathologies.

According to the main digitalization challenge of healthcare, we addressed our remarks to boost outcomes personalization features-based on the user/patient ones-starting from the contribute that high level datasets of users'/patients' health-data could bring to the enhancement of the healthcare sector. Sensitive health-data processing, indeed, plays a crucial role in the development and deployment of these specific tools that impact both on patients as a vulnerable group, and on individuals. In the event that the AI-based technologies would support wrong evaluations or would provide non-compliant data management strategies, they could lead to misjudgments and adverse outcomes. Consequences might be envisaged not only in terms of physical harm for a given patient, but they could also affect the psychological dimension of a given patient/end user. For this reason, our remarks deal with a further challenge, as hereinafter we will refer to the P5 medicine, including "psycho-cognitive" aspects arising from the commented technologies [7].

In this paper, therefore, we will deal with common issues – ranging from ethical-legal to technical ones-to properly address the compliance activities related to the technical, as well as ethical-legal, challenges emerging in the context of AI-based tools development for the implementation of an effective P5 medicine. At the same time, we will suggest possible solutions to cover gaps and lacks emerging from the available ethical-legal assessment tools.

## PRELIMINARY ACTIVITIES TO PERFORM A RESPONSIBLE ALTAI

ALTAI is a method to drive the self-assessment and it is based on 63 questions divided into the above-illustrated seven key-requirements. Interdisciplinary expertise is required to answer since the very first questions. A close yes/no answer is not sufficient, in fact, to perform the analysis and achieve a conscious and resolute opinion about the level of trustworthiness of the designed tool. The ALTAI is a method to drive the self-assessment [9].

For example, the first block of questions (Q1-Q6.2) is focused on assessing the impact of the designed application on fundamental rights. Therefore, a deep knowledge of what is meant for fundamental rights protection and how to assess the corresponding impact shall be introduced in the evaluation workflow. Furthermore, to properly answer the second block (Q7-Q18.5), a deep knowledge of cybersecurity and safety standards is required as, again, yes/no answers are not sufficient to the whole trustworthy evaluation process. In particular, the AI-designer shall justify the reasons that bring to choose a given technical measure as well as to implement a safeguard instead of another one under several grounds, including human safety, animal protection, environment, security, and misuse. Within the same context, it is required to identify a "fallback plan" aimed at ensuring the maintenance of an acceptable level of risks also if something goes wrong. Those evaluations are functional also in a prognostic perspective to answer questions included in the sixth block (Q49-Q53), related to the societal and environmental wellbeing. Therefore, the trustworthiness combines the need to prevent harms as well as to address the current multi-faced challenges of societal empowerment [10-12].

The third block (Q19-Q29.3) recalls the GDPR compliance and the results of the data protection assessment performed under article 35 GDPR for the given personal data processing shall be confirmed within the more comprehensive framework of the AI-based system. In this context, the Data Protection Officer (DPO) involvement is not only suggested, but also an organizational measure to be assessed in order to reach the trustworthy standard of the developed ecosystem. As a preliminary activity, the AI designer/developer has to verify whether or not a data protection officer and/or a privacy expert shall be consulted and appointed. This profile is strictly connected with the first block as the impact on the other fundamental rights than data protection reasonably stands in a cause-effect relationship with the protection of the confidentiality, availability, and integrity of personal data. For instance, a data breach concerning an AI-application for patients could also infringe their health, or private life, dignity, etc., but it could also identify possible grounds of discrimination according to the end users' vulnerability. This last profile on fairness is addressed by the fifth block (Q41-Q48) of the check-list assessment [12].

The fourth block (Q30-Q40.2) on transparency consists of a series of questions related to the development of each AI-based system, considering it as the result of a series of human decisions taken by

the AI controller. In a binding legal framework that defines roles and responsibilities, he/she may assume the role of AI-controller, similarly to the data controller as defined by article 5 GDPR. AI-based systems, in fact, firstly include the identification of methods for data acquisition regarding the function/algorithm that must be applied to a previously determined dataset. Secondly, the AI-controller shall define what tasks shall the AI perform (the so-called required actions) as well as the final purposes (goals) of the automated decision making/reasoning activities [13].

The last block (Q54-Q63) will assess the overall level of accountability of the process, suggesting the implementation of organizational measures aimed at monitoring the process and also providing solutions in case of issues.

A first step to be performed by the AI controller/developer is to engage an interdisciplinary team aimed at strengthening a dialogue and share best practices for the ALTAI purposes. Furthermore, an independent advisor might support the assessment in order to interpret notions and corresponding adjectives (e.g., what is meant for “meaningful interactions and appropriate human oversight and control” in Question n. 44, “an adequate working definition of fairness” in Question n. 44, “wide range of individual preferences” in Question n. 45).

Secondly, the identified specific parameters and standards shall drive the overall assessment and provide a robust and coherent internal framework of reference. This is true also for non-technical requirements, whose harmonized application shall emerge from the individual answers given to each block. For instance, if we identify a ground of vulnerability in respect to the impact on fundamental rights, the same analysis shall be reproduced within the assessment of fairness as well as within the governance-related issues and social impact ones [14,15].

Once that the technical and organizational measures to reach the acceptable level of trustworthiness of the AI-ecosystem have been implemented, a continuous evaluation system shall be maintained in order to ensure upgrades both in terms of performance and enhancement of fundamental rights. To this end, proper mechanisms of check and balance could be introduced within codes of conducts, encouraged by article 40 GDPR for personal data processing. They currently provide a compliance support for small and medium enterprises, addressing common issues and challenges in terms of self-regulation and best practices not only for personal data processing, but for the development of AI-based systems as well [16].

These preliminary remarks shall orient the AI-controller/developer towards a legal attentive design of the given application.

In the following paragraphs, we will present the results of a discussion between an interdisciplinary group, including scholars in law, biomedical engineering, and computer science, on how to develop a trustworthy AI-based tool aimed at early detecting the melanoma skin cancer.

The digitalization of the healthcare services, in fact, is boosted where existing data flows can be re-usable to train (such as in our case-study) an algorithm that could process information in order to predict a decision. These tools shall be enabled within a robust data governance ecosystem aimed at enabling the healthcare service among stakeholders and, at the same time, ensuring the exercise of their rights. Data collection, processing, and storage shall, therefore, allow mechanisms of re-training of the algorithms,

while providing a specific prediction/decision making result for the user/users. The interoperability for data formats and the enhancement of data security shall be combined with acceptable levels of pseudonymisation and anonymisation for the training, as well as on the linking and de-linking of records for the given query.

These remarks become crucial to enable innovative solutions, care delivery, including the opportunity for patients to control and administer care themselves. Considering that the P5 medicine includes keywords as participative, preventive and personalised, highlighting how the participation of the patient is paramount to prevent adverse outcomes in pathologies development, our aim is to assess the ALTAI while assessing a given technology in a mutual exchange of technical and organizational good practices. Under the purpose to overcome possible practical issues emerging from this innovative mechanism of assessment and, at the same time, to properly address the ethical-legal compliance in R&D&I sectors, we will highlight weaknesses and strengthens of the proposed structured evaluation system [17].

## BUILDING UP A TRUSTWORTHY AI-BASED TOOL FOR P5 MEDICINE

Our analysis starts from the need to develop an ethical-legal by design and by default AI-based tool to support the early-detection of melanomas. Melanoma has the highest mortality rate among skin cancers, and it can grow from the early stage (called melanoma insitu) to the latest stage (metastatic melanoma) in a period ranging from 8 to 12 months. An early diagnosis is essential to improve the survival rate and reduce treatments costs. Since this pathology appears on the skin surface, it can be detected by monitoring changes of the skin itself. This condition has peculiar morphological attributes and an expert clinician is needed to make a diagnosis. Nevertheless, melanoma lesions are not easy to detect in its early stage when these lesions present borderline features, even to an expert eye. AI technology may aid clinicians in the melanoma diagnosis, especially in its earlier stage. The use of the smartphone that has become a large-scale deployable tool, paired with this kind of technology may enable a common user to take an active and participative role into the skin cancer prevention [18-22].

As anticipated, despite of the internal allocation of technical tasks, the AI-controller/developer shall identify a series of roles aimed at giving advice and being responsible of some specific tasks within the ALTAI context. The first preliminary issue is the allocation of human resources in the design process. From a law and policy making perspective, indeed, incentive mechanisms shall be identified in order to overcome the lack of effectiveness of a (still) unbinding approach. In this regard, the GDPR legislative model seems particularly effective, as it includes a series of binding obligations framed within a structured illustration of principles, roles, and enforcing tools. In addition, a frame of optional mechanisms and safeguards that each legal system may decide to regulate or not. The second transversal issue is related to establishment of a multi-level governance system. This is functional to identify a monitoring process and coordinate the engagement of the previously identified roles.

In the analyzed P5 medicine tool, a technical board that includes medical doctors, biomedical engineers, software engineers, data protection and ethics experts shall be established to interpret the different interdisciplinary key-requirements that the ALTAI assessment provides. This technical board shall be able to arise issues and identify solutions in light of a general mutual purpose

to co-create a new technology enhancing fundamental rights protection. Such a technical board could also host different expertise for external advice as well as include stakeholders and end-users' opinions to assess different solutions.

Furthermore, the multilevel governance shall comply with the applicable legal framework that means for example that the GDPR governance in terms of appointment of joint controllers and data processors shall be structured according to the security governance determined by the specific standards followed by the developer. In addition, further engagements are envisaged by each block of questions. This may contribute to share responsibilities and to prove the overall accountability within the process.

The framework becomes more complex in case of public/private stakeholders as well as in case of cross-border relationships between the identified players since the national compliance process may present some gaps/overlapping profiles. In this regard, the possible legislative misalignment shall be addressed and covered during the assessment by specific agreements between the involved parties.

The third step is end-users centered. In fact, compliance activities shall deal with the main features that characterize those persons or groups of persons whose data are processed (i.e., the data subjects under the GDPR) and those who are the addressees of the prediction/automated decision-making process. The two categories might sometimes overlap, but they usually do not. The ALTAI

process shall deal with all possible end-users both to protect and enhance their rights. To this end, the technical board shall open to collect feedback as well to include a validation step with end-users to collect feedback not only on the technical level related to the automated decision-making, but also on its usability. In fact, in our example an AI-based tool will be addressed both to clinicians and patients with evident differences in terms of awareness, risks, benefits, and impact on corresponding rights. For instance, health protection will be assessed in terms individual fundamental right for each patient and in the collective dimension as far as the clinicians are concerned. At the same time, the supportive role of the AI shall be properly addressed in order to do not make the human decision too overconfident or, on the contrary, providing a replacement distress among professionals. Once that these profiles have been addressed as priorities by the AI-controller/developer, the ALTAI might be filled.

### ALTAI GAPS AND STRENGTHS

Questions developed within the ALTAI checklist have been interpreted to design an AI-based tool aimed at early-detecting melanoma by a decision-making system that processes images, in order to highlight gaps and strengths of the checklist.

The chart below shows the results of the ALTAI analysis and possible comments that may either address good practices or issues (Table 1).

**Table 1:** Results of the ALTAI analysis and possible comments.

Human Agency and oversight	Answer	Comments
Fundamental rights	<p>Patients fundamental rights involved are dignity health, data protection.</p> <p>Clinicians fundamental rights involved are dignity and work-life.</p> <p>The decision-making process could interact with patients stimulating their awareness towards the risk of melanoma, suggesting contacting a clinician.</p> <p>The tool could recommend the patient to get clinical advice and therefore it could interfere patient's decision to get or not to get a clinical advice.</p> <p>Individual vulnerabilities shall be addressed in the information section before using the tool.</p> <p>The clinicians may be supported in the pre-screening activities, but diagnosis shall be performed under the current clinical practice.</p>	<p>This key-requirement arises the need to identify the list of fundamental rights.</p> <p>The list could be included within the Terms &amp; Conditions of the given system/device/tool to accomplish an extensive information duty.</p> <p>A possible survey could be implemented to identify whether or not the end-user is vulnerable as well as her/his attitudes towards the results of the decision-making system (e.g., someone who assumes drugs could be temporary vulnerable and the use of the tool may cause undesirable consequences).</p> <p>Communication interfaces shall be addressed in a clear and user-friendly format.</p>
Human Agency	<p>The AI system shall support the cancer prevention actions of the healthcare system, revising internal processes in light of such a support in the pre-screening activities. No risks of overreliance and overconfidence in the AI system, as it gives just a pre-screening information. Clinical diagnosis shall in any case be provided through gold standard methodologies.</p>	<p>Disclaimer on the system characteristics shall be included both in the Terms and Conditions and in the handbook.</p> <p>Awareness and training campaigns shall be promoted among clinicians.</p>
Human oversight Technical robustness and safety [23]	Answer	Comments
Resilience to attack and security	<p>A monolithic system ensures a safer ecosystem. Regarding AI models, specific technical methodologies may be implemented to ameliorate model robustness against external perturbations and attacks. A further layer of security might be implemented in the model itself by (1) using ad hoc methods enabling the use on encrypted data as input to certain models, and (2) adopting models and libraries optimized for privacy sensitive applications (e.g., Opacus) [24].</p>	<p>Stress tests shall be scheduled to assess the overall resilience of the system to attacks and breaches</p>



Fallback plan and general safety	A fallback plan aligned to the general backup policy is applied. As the learnability of a specific problem is uncertain a-priori, a rigorous scientific methodology should be applied to the evaluation of a solution performance. To gauge the level of uncertainty associated to a specific prediction of a specific model, and hence to mitigate the effect of identifiable unreliable predictions, specific methodological approaches may be used (e.g. as reported in previous studies [25-28])	Adding redundancy in the software and hardware components of the system may also be a viable and widely adopted strategy to boost the overall system safety and ensure prompt recovery from anomalies, also enabling internal and automatic system diagnostic e.g. as done in avionics systems or, regarding AI models, using ensemble of models.
Accuracy	In general, a specific AI algorithm architecture is used to solve a certain task and needs a specific metric to correctly evaluate its performances. The AI-controller/developer is responsible to identify the proper statistical method for the specific case of implementation. For instance, a single shot detector model is used for object detection problems and its performances are evaluated using mean average precision metric [29]. As an example, in the specific case of classification task, the tool accuracy is evaluated as the ratio between the sum of True Positives (TP) and True Negatives (TN) values scored by the algorithm on the training data and the total number of training samples. Other metrics, such as sensitivity (ratio between True Positives and all positive samples) and specificity (ratio between True Negative and all Negative samples), are computed to better understand the learning level of the implemented algorithms.	In the current application, true Positives represent the number of positive samples (i.e., patients with melanoma) that are scored by the tool correctly affected by the pathology. True Negatives represent the number of negative samples (i.e., patients not affected by melanoma) the tool scores as healthy. In this field of application (i.e., melanoma detection), sensitivity is usually maximized at the expense of specificity.
Reliability and reproducibility	A specific monitoring system has been implemented to assess the algorithm performance. Code shall be extensively documented and it could be versioned using an open-source implementation for version management based on Git.	Extensive and detailed documentation of the scientific methodology followed to build and validate a specific AI solution should always be produced to ensure that independent reproduction and scrutiny of the technical implementation is feasible. When applicable, the source code of the specific implementation should also be made available at least for external technical review [30].
<b>Privacy, and data governance</b>	<b>Answer</b>	<b>Comments</b>
Respect for privacy and data protection	End-users' consent is the legal basis of the data processing provided with the tool. As a consequence, technical and organizational measures shall be implemented to inform, and to let the data subject exercise her/his rights. For instance, the interface will allow to enable/disable both for algorithm continuous training and the decision-making process anytime. A data protection impact assessment is performed and risks for the availability, confidentiality, integrity of data shall be mitigated in terms of acceptability. The DPOs shall be involved in the process. Within the step of training of the AI-tool: Data governance and ownership shall be governed in compliance with the applicable legal framework. Therefore, data shall be anonymous or, in case of pseudonymised ones in the frame of a collection provided directly from trials in the clinics, flows shall be encrypted, pseudonymisation techniques shall be applied and the involvement of the competent ethics committees shall be included in the process.	Within this section, the AI-assessment deals with all the GDPR compliance activities. These activities include the identification of the data governance, the security measures, and the identification of roles and responsibilities that ensure that all data processing are provided in light of the principles of lawfulness, fairness, and transparency, accountability, purposes limitation, data minimisation, and accuracy. For those AI-controllers/developers that are not processing personal data, they shall be compliant with the EU Regulation (EU) 2018/1807 of the European Parliament and of the Council of 14 November 2018 on a framework for the free flow of non-personal data in the European Union, that entered into force on 28 May 2019 [31].
Quality and integrity of data	The development of the tool is framed in a detailed privacy governance policy, following a tailored management/authorization path.	Data pooling activities shall pay attention to the possibility that to process information by crossing anonymous and pseudonymised datasets may incur the risk of re-identification or, at least, the necessity to regulate data flows through specific data sharing agreements.
Access to data	Dedicated servers, with limited access, host the development of the tool and anonymised/pseudonymised data.	The internal and external governance as well as technical measures shall be identified, implemented and maintained for the entire life-cycle of the data according to the results of the impact assessment.

Transparency [30]	Answer	Comments
Traceability	<p>User is guided in the data gathering process (e.g., using viewfinder to centre the skin lesion, checking the blur of the image), ensuring a correct data quality. This is essential to guarantee the correct pre-processing phase, before feeding the algorithm using the user's data.</p> <p>Decision-making outputs can be monitored using thresholds. As an example, in case of probability as output, data that gives as output probabilities within certain ranges, rising concern on the reliability of the algorithm output, are forwarded to a clinician, as described in the Communication requirement. These data paired with the clinician response (i.e., the ground truth or target) enlarge the original dataset that it is used to re-train the AI algorithm. Here, metrics described in the Accuracy section are used to evaluate the algorithm training, monitoring its performances.</p>	<p>Traceability is functional to address the transparency of the AI-system in order to verify the quality of the decision-making outputs. AI algorithms are strongly dependent on the data format and quality used to train them to boost their learning ability and produce the desired result (in this case, detecting a malignant skin lesion). The format and quality must be similar to the ones the algorithm uses during the training phase. AI algorithms can be considered "black boxes", since they process data in ways that are not easily audited or understood by human. This limits the traceability of what happens within the algorithm. Nevertheless, strategies can be implemented to monitor the data quality and the algorithm performances (e.g. blockchain technologies) [31].</p>
Explainability [32]	<p>A viable inclusion of computer-aided design system in melanoma diagnosis pipeline was shown by Waterhouse et al. in 2019. The AI-decision may influence the screening activities, according to the user's experience, including clinicians. As an alternative, other recent methods tend to "open" black boxes models such as AI algorithms. Here the authors provided a classification of approaches to interpret the decision-making process based on the problem definition and the black box type [33,34].</p>	<p>The experience plays a crucial role in trusting an AI-based decision. Unexpert users tend to follow the results output by the AI tool. This is also true for clinicians. Indeed, it has been shown that clinicians with least experience in general follow the AI-tool decisions if it contradicted their initial diagnosis, even if they were confident. However, faulty AI can mislead both experienced and least experienced clinicians. This is an aspect to consider when deploying these tools. This could be overcome providing additional features to the clinician (e.g., for the current application: asymmetry index, border irregularity index) that can aid their final decision.</p>
Communication	<p>The interface shall be developed in accordance with the given task and the involved roles (clinicians or end users).</p>	<p>It has been demonstrated that for clinicians it is better to show multiclass probability when dealing with multiclass diagnosis. Conversely, malignancy probability can be used to manage binary decisions (e.g., whether or not do a biopsy, whether or not go visit an expert clinician). As described in the "explainability" requirement, unexpert user tend to follow the AI-results, even in the case of wrong suggestion. Melanoma growth happens within months, thus missing the detection of this lesion type can lead to adverse outcomes. The solely probability could be ambiguous, especially in case of intermediate values. In the effort of maximize the sensitivity of the tool, giving clear responses to the user, the clinician himself/herself can be involved in this detection process. In case of probability/uncertainty, the lesion data could be forwarded to a clinician for an additional evaluation. The clinician decision is then used as the final response to the user as well as to improve the AI algorithm performances. Indeed, the new data, paired with the clinician response (i.e., the ground truth or target), enlarges the original training dataset and it is used to re-train the AI algorithm [35-37].</p>
Diversity and non-discrimination	Answer	Comments
Unfair bias avoidance	<p>Biases are possible. Biases are mainly related to the personal characteristics that compose the training datasets. Melanoma incidence rate indeed is higher in those people with lower phototypes, while is almost rare in those with higher ones (e.g., black people). Datasets will thus be biased in the skin colour. Data gathering and assumptions made when training an algorithm can also introduce biases, thereby distorting the final output. These are important aspects to consider when dealing with automatic algorithm to prevent injustice and discriminations. As an example, this study showed that using the health costs as a proxy of health needs introduced a racial bias in the system, and the algorithm consider black patients healthier than equally sick white patients. This led to a reduction of 28.8% in the number of black patients to be considered with high priority of health needs [38,39].</p>	<p>There are three types of bias: Productive bias, bias that someone would qualify as unfair, and discrimination bias. Algorithms cannot be unbiased. Bias in machine learning guarantees the success in modelling a certain distribution of data, according to the "No Free Lunch (NFL) Theorem", thus solving the task of concern. The choice of cost function that the model must minimize to converge to the solution, the purpose and the use of limited training, and the test data are examples of productive bias. Also, the assumption that training data distribution will be the same of the one of the test data represents a production since in practice is often violated. However, solutions can be implemented in order to understand, mitigate or account for bias, as summarized in this survey by Ntoutsis. The discriminatory level of ML models can be limited adding constraints to the model optimization problem. Thus, it must be found a trade-off between constraints and accuracy (since it can worsen when adding to many constraints) [40,41].</p>

Societal and environmental wellbeing	Answer	Comments
Sustainable and environmentally friendly AI	One main weakness of using AI methods is tied to the high computational demand in terms of both hardware and energy. Hence, electronic and consumables replacement/disposal are involved. Moreover, the type of model-task to be solved may require the use of especially complex and computationally expensive models (e.g. text based applications). In general, more complex models require higher expense in both economical, energetic and availability terms.	A limitation of AI tools lies on the expensive hardware that is needed, especially for the training process of deep learning models. The use of GPUs is essential to increase the complexity of the model architecture, accelerating its learning process. The more GPUs are used, the more complex the task to be solved (and hence the implemented model) can be. The price of one good GPU is around \$1,800 each and computational cloud resources may be rented. Moreover specific technical methods may be adopted to reduce energy expenditure while only marginally affecting final performance [42,43].
Social impact	The AI system contributes to the P5 medicine purposes.	As illustrated in the previous paragraphs, the effectiveness of these AI-tools is strongly connected to the clinician-patient relationship. The empowerment of patient's awareness on cancer diseases is aligned with the public prevention purposes to enhance current society. This has also impactful consequences on the improvement of the healthcare systems and services.  Reproducibility and transparency of the scientific methods applied to develop a given AI-based system boost innovation towards the mentioned societal challenges.
Society and democracy	The AI system contributes to the P5 medicine purposes.	Multiple and bottom-up solutions, that engage stakeholders, including vulnerable groups, may contribute to overcome societal barriers, reaching a more inclusive society. Trustworthiness of the given AI-based system constitutes a pre-condition to meet these challenges and enhance democratic values.
Accountability	Answer	Comments
Auditability	As stated in the Fallback plan and general safety requirement, adding redundancy to the AI method and continuous technical assessment could create a system of internal audit	Also to schedule an independent audit to better test under common standards the developed technology could provide an accountable measure. The activity might be time-consuming and expensive, however the results are usually functional to align knowledge, competence, and skills useful to improve one's approach towards innovation.
Minimising and reporting negative impact	In the current clinical scenario, considering the possible different level of education, sensitivity of the end-users, other stakeholders might be involved in case of pre-identified vulnerabilities.	A survey could be implemented in the AI-tool interface in order to limit the access and use in case of specific vulnerabilities. The same solution could be adopted for minors. The tool can be used by a legal representative in case of incapacity of the end-user. In other cases, the same AI-tool interface may suggest the user to not use the tool alone.
Documenting trade-offs	Within the tool development, an efficient management shall ensure a continuous monitoring of technical activities and compliance ones.	A trade-off analysis is part of the R&D&I life-cycle. The decision-making process provided by the AI-developer is continuously addressed to assess consequences to losing one quality, aspect or amount of something in return for gaining another quality, aspect or amount that shall in any case considered as trustworthy as the first one. Any decision shall be documented in order to be able to intervene if conditions change.
Ability to redress	Tests will evaluate the accuracy of the automated decision making. In case of adverse impact occurs, specific mechanisms will ensure adequate redress.	The user could be invited to self-evaluate the accuracy in terms of reproducibility of the obtained result. A survey could be implemented in the AI-tool interface to evaluate it and to report any possible adverse answer.

The chart summarizes those answers provided during an ALTAI session aiming at evaluating the level of trustworthiness of a given technology in order to discuss possible technical and organizational measures to be implemented in an AI-based tool to reach an acceptable level of trustworthiness. A possible limit of the self-assessment approach consists of the fact that AI-controllers/developers may encounter difficulties in explaining, in plain language, how a given tool works and its consequent functionalities. This is particularly evident in case of code interpretation, even between AI-developers, but also confirmed for different domain experts. This test of the ALTAI checklist constitutes a unique exercise of semantic alignment, awareness

development, and interdisciplinary training, strongly impacting on possible standardization of skills and competence involved in the AI compliance processes.

As shown in the "comments" column of the chart, the checklist becomes a very useful tool for self-assessment only if it is accompanied by a series of good practices aimed at:

- i) Developing a common syllabus useful to align competences and skills among the different experts involved in the assessment.
- ii) Tailoring the ground of analysis to the specific sectors where the AI-ecosystem shall be performed (e.g., health sector, workplace,

mobility.

iii) Addressing specific actions to firstly identify individual and group vulnerabilities, and then to overcome the relative barriers.

iv) Addressing specific actions to pursue the mitigation actions beyond the development step in order to transfer also in the market context the trustworthy knowledge and know-how developed during the assessment.

v) The technical board is updated to the highest standards applicable to the given sector/market where the AI-tool is placed.

In addition, the self-assessment approach shall be promoted through policy-making incentives in order to ensure its application also without mechanisms of enforcement. In this regard, the Provisional Resolution specifies that they could constitute “a good starting point but cannot ensure that developers, deployers and users act fairly and guarantee the effective protection of individuals” and that in any cases, technology-neutral as well as specific standards shall be developed where appropriate. In particular, as far as artificial intelligence, robotics and related technologies are concerned, it suggests providing “mandatory compliance with legal obligations and ethical principles as laid down in the regulatory framework for AI” to be performed through “an impartial, regulated and external ex-ante assessment based on concrete and defined criteria”. Research and innovation, therefore, are those key-sectors where the ALTAI checklist or other impact assessment methodologies could find application.

In this context, the illustrated steps will be included in every AI-related project life-cycle by design, becoming a significant component of the research integrity and reproducibility within the scientific methodology [23-47].

## NEW REGULATORY CHALLENGES FOR AI IN THE HEALTHCARE SYSTEM

In this paper, we provided an assessment of the ALTAI checklist considering the possible issues emerging while providing the evaluation for an AI-based tool applied to cancer medicine and, specifically, to early detecting the melanoma skin cancer.

Firstly, the interdisciplinary approach that characterized the development of the check-list shall be applied also in the executive phase and maintained in the entire life-cycle of the technology development. The methodological outcome of the analysis is needed both to interpret and then to accomplish to the technical and organizational measures to be implemented as a consequence of the risk-based analysis.

Secondly, considering that the above-mentioned Provisional Resolution identifies, as high risk ones, the applications whose “development, deployment and use entail a significant risk of causing injury or harm to individuals or society, in breach of fundamental rights and safety rules as laid down in Union law”, the role of AI-based systems in P5 medicine shall require the inclusion of standardized notions, tailored risks, and mechanisms to ensure either the coherence among the sectorial legislative frameworks or the opportunity to re-assess their impact anytime the scientific progress could affect one of their fields of application.

Thirdly, the frontiers of P5 medicine are significantly affected by the ongoing debate on the regulatory framework for AI. The opportunity to get access to health data and to re-use them for algorithms training purposes, as well as for providing further

information on a given patient, is crucial for innovation in terms of effectiveness and sustainability of the developed solutions. In addition, the chance to manage health data despite of the means, time, site of their collection is a strategy that could enhance the competitiveness of the related industrial sector and, at the same time, it could promote inclusiveness and awareness among citizens.

To manage personal and non-personal data in compliance with the ethical-legal framework, therefore, constitutes a pre-condition to implement innovative tools and solutions in healthcare, that could enhance diagnosis and treatment and, at the same time, human dignity through the empowerment of the patient control over her data. A new regulatory framework on AI, indeed, shall provide specific instructions on how to develop per se compliant hosting infrastructures and it shall address the proper consistency mechanisms to ensure that any data processing could be mapped, assessed, and enabled for pre-determined purposes. Human dignity and fundamental rights shall be protected by design and by default with concrete action points that could be understandable, applicable, and enforceable in every R&D&I context.

## CONCLUSION

Thus, to promote a system of certification could represent a valuable solution to balance the need of procedures standardization, monitoring, and control of the level of compliance. Moreover, a multilevel system of enforcement, based on the accountability principle, and then on liability, could facilitate the establishment of a fruitful dialogue between developers and users, aiming at consolidating trust and awareness among citizens and stakeholders towards such a technological, ethical-legal revolution that AI brought in our society.

## REFERENCES

1. AIA. Independent high-level expert group on artificial intelligence: The Assessment List for Trustworthy Artificial Intelligence (ALTAI) for self-assessment. 2020.
2. EP. Framework of ethical aspects of artificial intelligence, robotics and related technologies. 2020.
3. Hagendorff T. The Ethics of AI Ethics: An Evaluation of Guidelines. *Minds Mach.* 2020; 30: 99-120.
4. EU Commission. High-Level Expert Group on AI, ‘Ethics Guidelines for Trustworthy Artificial Intelligence. 2020.
5. Comandè G. Unfolding the legal component of trustworthy AI: A must to avoid ethics washing. SSRN. 2020.
6. Kumar A, Braud T, Tarkoma S, Hui P. Trustworthy AI in the Age of Pervasive Computing and Big Data. *IEEE.* 2020.
7. EU Commission. Communication from the commission to the european parliament, the council, the european economic and social committee and the committee of the regions. 2020.
8. Gorini A, Pravettoni G. P5 medicine: a plus for a personalized approach to oncology. *Nat Rev Clin Oncol.* 2010; 8: 444
9. Davola A, Black B, Gulson K, Rockwell G, Selinger E. From shortcut to sleight of hand: why the checklist approach in the EU guidelines does not work. *AI pulse.* 2019.
10. Brüggemeier GA, Ciacchi C, Comandè G. *Fundamental Rights and Private Law in the European Union.* Cambridge University Press. 2010.
11. Khakurel J, Penzenstadler B, Porras J, Knutas A, Zhang W. The Rise of Artificial Intelligence under the Lens of Sustainability. *Technologies.* 2018; 6:10.



12. Scherer MU. Regulating artificial intelligence systems: Risks, challenges, Competencies, and strategies. *Harv J Law & Technol.* 2016; 29(2): 354-398
13. FRA. Handbook on European data protection law. 2018.
14. Amram D. The Role of the GDPR in designing the European Strategy on Artificial Intelligence: Law-making potentialities of a recurrent synecdoche. *Opinio Juris in Comparatione.* 2020.
15. Mittelstadt B. Principles alone cannot guarantee ethical AI. *Nat Mach Intell.* 2019; 1: 501-507.
16. Malgieri G. Automated decision-making in the EU Member States: The right to explanation and other “suitable safeguards” in the national legislations. *Comp Law and Secu Rev.* 2019;35 (5):105327.
17. Malgieri G (eds.). Guide to the processing and security of personal data. *Isole24Ore.* 2019.
18. Pesapane F, Volonté C, Codari M, Sardanelli F. Artificial intelligence as a medical device in radiology: ethical and regulatory issues in Europe and the United States. *Insig imag.* 2018;9(5):745-753.
19. Chartrand G, Cheng PM, Vorontsov E, Drozdal M, Turcotte S, Pal CJ, et al. Deep learning: A primer for radiologists. *Radiographics.* 2017;37(7):2113-2131.
20. Krittanawong C, Zhang H, Wang Z, Aydar M, Kitai T. Artificial intelligence in precision cardiovascular medicine. *J Am Coll Cardiol.* 2017;69(21):2657-2664.
21. World Health Organization. Cancer mortality database. 2018.
22. Liu W, Dowling JP, Murray WK, McArthur GA, Thompson JF, Wolfe R, et al. ‘Rate of Growth in Melanomas: Characteristics and Associations of Rapidly Growing Melanomas’. *Arch Dermatol.* 2006;142(12):1551-1558.
23. American Joint Committee on Cancer. Melanoma research alliance. Melanoma staging.
24. Matsumoto M, Secrest A, Anderson A, Saul MI, Ho J, Kirkwood JM, Ferris LK. Estimating the cost of skin cancer detection by dermatology providers in a large health care system. *J Am Acad Dermatol.* 2018; 78:701-709.
25. Papageorgiou V, Apalla Z, Sotiriou E, Papageorgiou C, Lazaridou E, Vakirlis S, et al. The limitations of dermoscopy: false positive and false negative tumours. *J Eur Acad Dermatology Venereol.* 2018; 32(6):879-888.
26. Sculley D, Holt G, Golovin D, Davydov E, Phillips T, Ebner D, et al. Hidden technical debt in machine learning systems. *Adv Neu Info Proce Sys.* 2015;28:2503-2511.
27. Gilad-Bachrach R, Dowlin N, Laine K, Lauter K, Naehrig M, Wernsing J. Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy. In *International Conference on Machine Learning* 2016;pp. 201-210.
28. David BS, Hrubes P, Moran S, Shpilka A, Yehudayoff A. Learnability can be undecidable. *Nat Machi Intel.* 2019;1(1):44-48.
29. Tagasovska N, Lopez-Paz D. ‘Single-Model Uncertainties for Deep Learning’. *Neu Infor Proce Systems* 2019;32: 1811.00908.
30. Gal Y, Ghahramani Z. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *int Conf Machine learning.* 2016;pp. 1050-1059.
31. Gal Y. Uncertainty in deep learning. *University of Cambridge.* 2016;1(3):4.
32. Walters WP. Code Sharing in the open science era. *J Chem Info Model.* 2020;60(10):4417-4420.
33. Eglén SJ, Marwick B, Halchenko YO, Hanke M, Sufi S, Gleeson P, et al. Toward standard practices for sharing computer code and programs in neuroscience. *Nature neuroscience,* 2017;20(6):770-773.
34. Regulation (EU) 2018/1807 of the European Parliament and of the Council of 14 November 2018 on a framework for the free flow of non-personal data in the European Union. *Official J Euro Union.* 2018.
35. Tschandl P, Rinner C, Apalla Z. Human-computer collaboration for skin cancer recognition. *Nat Med.* 2020; 26: 1229-1234.
36. Marechaux JL. Towards advanced artificial intelligence using block chain technologies. 2019.
37. Rader E, Cotter K, Cho J. Explanations as mechanisms for supporting algorithmic transparency. *CHI 18.* 2018; 1-13.
38. Waterhouse DJ, Fitzpatrick CRM, Pogue BW. A roadmap for the clinical implementation of optical-imaging biomarkers. *Nat Biomed Eng.* 2019; 3: 339-353.
39. Guidotti R, Monreale A, Ruggieri S, Turini F, Giannotti F. A Survey of Methods for Explaining Black Box Models. *ACM Comput Surv.* 2018 ; 51(5) : 42.
40. Obermeyer Z. Dissecting racial bias in an algorithm used to manage the health of populations. *AAAS.* 2019; 366(6464): 447-453.
41. Canziani A. An Analysis of Deep Neural Network Models for Practical Applications. 2016.
42. Chen PC, Liu Y, Peng L. How to develop machine learning models for healthcare. *Nat Mater.* 2019; 18(5): 410-414.
43. Ntoutsi E, Fafalios P, Gadiraju U. Bias in data driven artificial intelligence systems: An introductory survey. *Wires Data Min Knowl.* 2020; 10(3): e1356.
44. Yang T, Chen Y, Emer J, Sze V. A method to estimate the energy consumption of deep neural networks. *IEEE.* 2017.
45. Chouldechova A. Fair prediction with disparate impact: A study of bias in recidivism prediction instruments. 2017
46. Lo Piano S. Ethical principles in machine learning and artificial intelligence: Cases from the field and possible ways forward. *Humanit Soc Sci Commun.* 2020; 7(9).
47. EU Commission. Reproducibility of Scientific Results in the EU. 2020.