

A survey on *good* AI: User-Centric AI Design in Healthcare

Andrea Berti^{1,2*}, Valentina Giannini^{3,4}, Simone Mazzetti³,
Maria Antonietta Pascali², Daniele Regge³, Sara Colantonio²

^{1*}Department of Information Engineering, University of Pisa, Via
Caruso 16, Pisa, 56122, Pisa, Italy.

²Institute of Information Science and Technologies (ISTI), National
Research Council (CNR), Via Moruzzi 1, Pisa, 56124, Pisa, Italy.

³Department of Radiology, Candiolo Cancer Institute, FPO-IRCCS,
Strada Provinciale 142, Km 3.95, Candiolo, 10060, Turin, Italy.

⁴Department of Surgical Sciences, University of Turin, via Genova 3,
Turin, 10126, Turin, Italy.

*Corresponding author(s). E-mail(s): andrea.berti@isti.cnr.it;
Contributing authors: gianninivalentina@gmail.com;
simone.mazzetti@ircc.it; maria.antonietta.pascali@isti.cnr.it;
daniele.regge@ircc.it; sara.colantonio@isti.cnr.it;

Abstract

The integration of Artificial Intelligence (AI) in healthcare has the potential to revolutionize patient care by enhancing diagnostic processes, treatment protocols, and overall healthcare delivery. However, the adoption of AI-powered tools and services is contingent upon establishing a robust foundation of trust among healthcare professionals. The ProCAncer-I project, informed by the FUTURE-AI framework, is at the forefront of this effort, promoting a user-centric design philosophy that prioritizes the needs and expectations of end-users, primarily clinicians and radiologists. This paper delves into the co-design methodology adopted by an interdisciplinary team, elucidating the collaborative efforts that underpin the customization of the FUTURE-AI principles to align with the clinical requirements of the project's partners. The introduction sets the stage for a comprehensive discussion on the significance of stakeholder engagement in the design and implementation of trustworthy AI systems within clinical settings.

Keywords: Artificial Intelligence in Healthcare, Artificial Intelligence, Explainability, Trustworthiness, User-centric Design

1 Introduction

The integration of Artificial Intelligence (AI) in healthcare marks a significant transformation, enhancing diagnostic precision, treatment outcomes, and patient care [1–3]. AI tools can leverage large datasets and identify patterns to surpass human performance in several healthcare aspects, offering increased accuracy, reduced costs, and time savings when aiding humans. Trust among clinicians is pivotal for the successful adoption of AI tools [4–6], a principle that the ProCancer-I European project, guided by the FUTURE-AI framework¹, endorses. The framework recognizes the urgency of user-centric design in fostering this trust and advocates for a multidisciplinary approach to design, involving stakeholders and end-users throughout the development process.

Co-design, a creative partnership that spans the entire design journey [7], is central to this approach. It involves stakeholders and end-users in a dynamic process that includes consultations, workshops, focus groups, and dedicated collaborative tools (e.g., Mirò, cards and games, simulated environments) to spur innovation and ensure inclusivity. The ProCancer-I project aimed to assess how the FUTURE-AI principles could be tailored to meet the specific expectations and clinical demands of its partners. To this end, an interdisciplinary team crafted a survey to capture the clinical partners’ perspectives, desires, and expectations for high-quality, trustworthy AI systems.

This team, consisting of experts in radiology, biomedical engineering, computer science, and mathematics, all with extensive experience in cancer imaging and AI, worked together from June to September 2022. They established a common language, defined the survey’s focus, and collaborated on a shared document. The finalized survey, distributed via an online Google form, was distributed among the clinical partners to gather their insights. The findings from this survey and their implications are discussed in the following sections.

2 Survey content and structure

The survey was meticulously designed to cover a comprehensive range of topics crucial for the clinical application of AI. It aimed to gather both quantitative and qualitative insights on the clinicians’ expertise with AI, their views on unreliable AI interventions, preferred methods of interaction with AI systems, performance expectations, and the attributes they consider most important in a trustworthy AI system.

A total of 26 questions were crafted, varying from multiple-choice to open-ended, to capture a wide array of data. The survey began with an introduction explaining its goals (Figure 1) and included a glossary to ensure participants fully understood the terms used (Figure 2). The topics addressed in the survey were:

- clinical expertise and current usage of AI tools (8 questions);
- opinions on unreliable AI-powered interventions (1 question);
- preferred reading modality and interaction with the AI system (3 questions);
- desired balance of sensitivity and specificity for different clinical tasks (4 questions);
- expected success rate for various clinical tasks (2 questions);
- most valued features of trustworthiness (4 questions);

¹<https://future-ai.eu/>

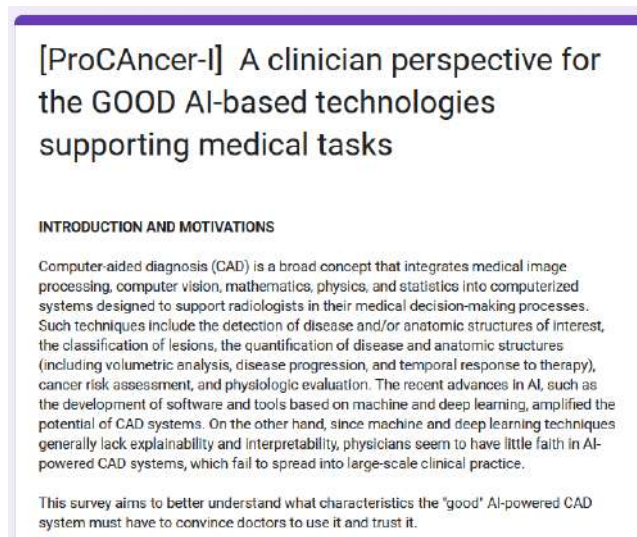


Fig. 1 Entry page and motivations of ProCancer-I survey on trustworthy and good AI

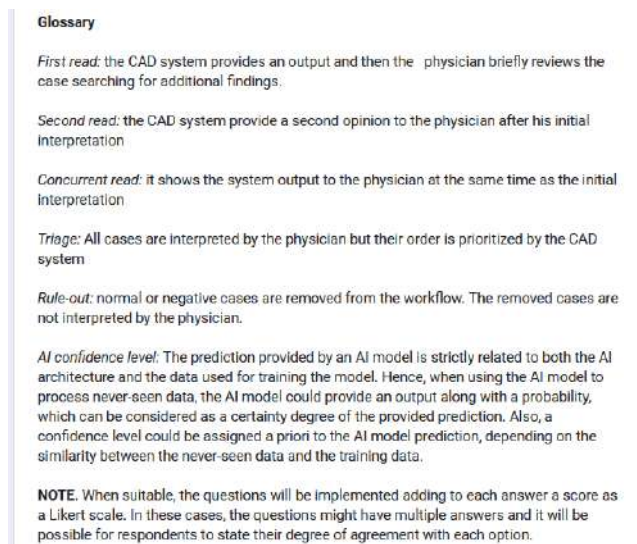


Fig. 2 The glossary included in the introductory section of the ProCancer-I survey

- most valued features of reliability (2 questions);
- preferred elements and format of AI output (2 questions).

The survey was concise to respect the clinicians' limited time, and it concluded with an open-ended question to give the opportunity for additional

feedback. We report in Figures 3 – 6 some of the most relevant questions included in the questionnaire; the complete version of the questionnaire can be found at the following link: <https://docs.google.com/forms/d/e/1FAIpQLScvtWfzJcRg0c7Hcu4vqQjdYOLk6ooLTFXq8xj6XXBjNfsFhw/viewform>

9. Concerning the reading modalities, independently of the clinical application, a CAD system may support your work in several ways. Please assign a number from 1 to 5 based on your preferences. (1: strongly disagree, 5 strongly agree) *

	Strongly disagree	Disagree	Neither agree nor disagree	Agree	Strongly agree
First reader	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Concurrent-reader	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Second reader	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Triage	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Rule-out	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

9.bis None of the previous options. I would prefer.. *

10. How would you like to interact with an AI-powered system? **

- The AI tool should run always in the background, providing an output for all examinations
- The AI tool should run only after request, e.g., by clicking an "AI" button to generate the output
- Other _____

Fig. 3 Questions on preferred reading and interaction modalities

3 Analysis of survey findings

The survey findings provide a rich variety of data, reflecting the diverse perspectives of the participants. We collected a total of 38 responses from October 2022 and March 2023. We carried out the analysis by utilizing Python scientific libraries. Many of the survey questions utilized a Likert scale format, allowing the responses to be analyzed as sentiment scores by assigning an integer value ranging from -2 (complete disagreement) to $+2$ (complete agreement) to the possible answer values. With that notation, a sentiment score of 0 would represent neutrality, a positive score indicates affirmative feedback, and a negative score reflects unfavorable feedback.

Demographic data of the survey participants are concisely presented in Figure 7. Typically, the respondents were European radiologists over the age of 35, predominantly employed in public healthcare settings, with a specialization in abdominal imaging and a basic understanding of AI.

The survey revealed that 56% of participants currently employ AI tools, predominantly in the context of research, education, or for clarification purposes. Conversely,

15. In your opinion, which characteristic should an AI-powered system have to be considered trustworthy? Please assign a number from 1 to 5 based on how much you agree. 1: strongly disagree, 5 strongly agree.

	Strongly disagree	Disagree	Neither agree nor disagree	Agree	Strongly agree
It should be reliable: have a high success rate.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
It should be robust: work in any conditions.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
It should be clearly explained: it should provide motivations for its decision.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
It should be interpretable: its behavior should be clear to me, even though it does not provide me with any explanation.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
It should be certified by a certification body	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
It should be transparent: it should provide the important information related to datasets used to train and validate the model should be provided	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

15.bis Please list other features that would increase your trust in the CAD system (sorting them from the most important to the least important)

16. Which additional information related to the development of the AI-powered system do you consider important for you to trust the system? Please assign a value from 1 to 5 based on the importance of the information. 1: not important, 5: very important.

	Irrelevant	Not important	Neutral	Important	Very important
Scientific peer-reviewed publications reporting methods and performances of the system in the medical community.	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The quality, origin and sample size used to train the system	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
The quality, origin and sample size used to validate the system	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Demographic and clinical characteristics of the population used to train and validate the system (i.e., race, PSA range, percentage of positive cases, gene mutations, etc.)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

16.bis Other than the above (please specify): *

Fig. 4 Questions on most valued trustworthiness features and additional information

41% of respondents indicated they do not utilize AI tools due to lack of access. Only 3%, which corresponds to a single respondent, expressed distrust in AI tools (Figure 8 left). For users who have adopted AI, its primary application has been in the detection of diseases (Figure 8 right).

Regarding the preferred reading modality, the "second reader" option emerged as the most favored, with a sentiment score of 0.68. In contrast, the "rule-out" option was met with disfavor, reflected in a negative sentiment score of -0.22 , which aligns with expectations given that such systems entirely bypass radiologist consultation. Notably, all reading modalities received scores under 1, as can be seen in Figure 10. This outcome, coupled with the fact that nearly all participants use, or would like to use, AI, implies a need for clearer, more task-specific options. This is partly supported by the findings from the question about distrust, where no distrust cases were reported by the majority of respondents (Figure 9). Additionally, open-ended feedback highlighted ongoing concerns regarding lesion characterization among clinicians.

The analysis also showed that the preferred integration modality in the examination workflow was for the AI tool to always run in the background, providing an output for all examinations. This can be seen in Figure 10 right.

17. When could you consider a diagnostic tool sufficiently reliable to be used in clinical practice? Please assign a value from 1 to 5. 1: totally disagree, 5 totally agree.

	Strongly disagree	Disagree	Neither agree nor disagree	Agree	Strongly agree
When I have tested the tool for enough time to demonstrate that its predictions are accurate	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Only if I get an explanation on how it works along with an accurate prediction	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
when I do not understand how it works	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
when it displays visual information about the regions of the images that have been used to get the prediction	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
when the system shows me cases that are similar to the one at hand as comparative examples	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
when a similarity score is provided estimating the proximity of the input data to the training data	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
when I am informed on which radiomics features have determined the output	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

Fig. 5 Questions on most valued reliability features

Regarding the expected detection rates for detection tasks, such as for lesion detection, among positive cases, neither of the proposed options were particularly favored by the respondents. Only one option was distinctly unpopular, i.e., the one regarding the balance between false positives and false negatives. Figure 11 shows the sentiment score and the statistics for this question for positive (up-left) negative cases (down-left) and those on performance metrics for positive (up-right) and negative (down-right) cases. One supplementary free-text answer was also of interest on this topic:

Lesion detection is a tricky subject. False negatives must first be avoided; however, false positives are also very detrimental to workflow when significant in number. After determining the highest possible sensitivity, a minimum degree of specificity must also be maintained.

In terms of trustworthiness, the respondents appreciated all the provided options, with a particular inclination towards aspects such as reliability, certification, and

18. Reporting the results of CAD assessment depends on the provided output. However, considering an agnostic setting, which type of AI output would you prefer? Please assign a value from 1 to 5. 1: totally disagree, 5 totally agree.

	Strongly disagree	Disagree	Neither agree nor disagree	Agree	Strongly agree
Binary classification (e.g., positive vs negative, aggressive vs non-aggressive)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Binary classification + an additional class "indeterminate" for ambiguous cases	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Binary classification + AI confidence level (explaining the certainty degree of the provided output)	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Continuous score from 0 to 100	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>
Color-coded likelihood maps	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>	<input type="radio"/>

18 bis Other than the above (please specify): *

Fig. 6 Questions on preferred way to receive AI models' output

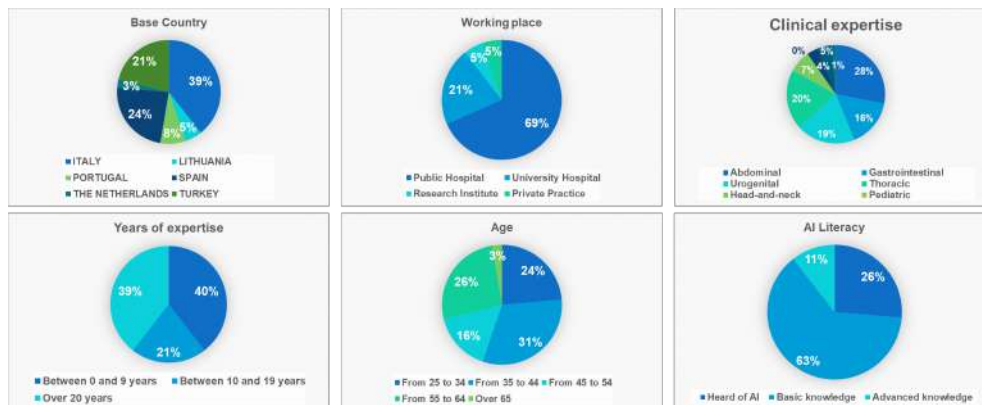


Fig. 7 Demographic and professional profile of respondents (in the "clinical expertise" chart, 0% correspond to pediatric specialisation)

openness, as shown on the left of Figure 12. When it comes to supplementary details that could enhance the credibility of an AI instrument, the participants once again expressed a favorable opinion for all listed options, showing a marked preference for



Fig. 8 Left: the use of AI in clinical workflow. Right: the clinical task for which AI is used



Fig. 9 Possible cases of distrust

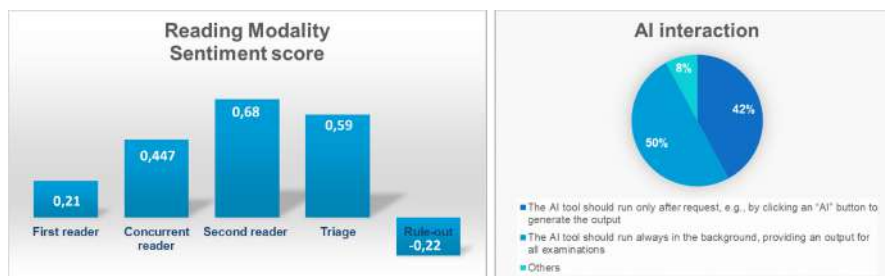


Fig. 10 Left: preferences for the reading modality. Right: preferred workflow integration

revealing details concerning the quality and the provenance of the data employed in training the AI model in question (Figure 12 on the right). This emphasizes the importance of being transparent about the data used.

In terms of reliability, participants placed as the highest valued feature the practical on-field demonstration of the AI tool’s accuracy in everyday usage. This was closely followed by the importance of having explanations for the AI model’s outputs and its precision (Figure 13). The preference for the AI tool to be understandable and its operations to be transparent to radiologists seemed less important. These observations imply that the primary concern is the AI tool’s accuracy. Ideally, the results produced should be explainable. An in-depth access to the AI tool’s mechanics is not deemed essential, provided that the outcomes are dependable and supported by a logical explanation.

Lastly, regarding the way results should be presented, respondents showed a preference for an output also including a confidence score of the prediction. This is reported in Figure 14.

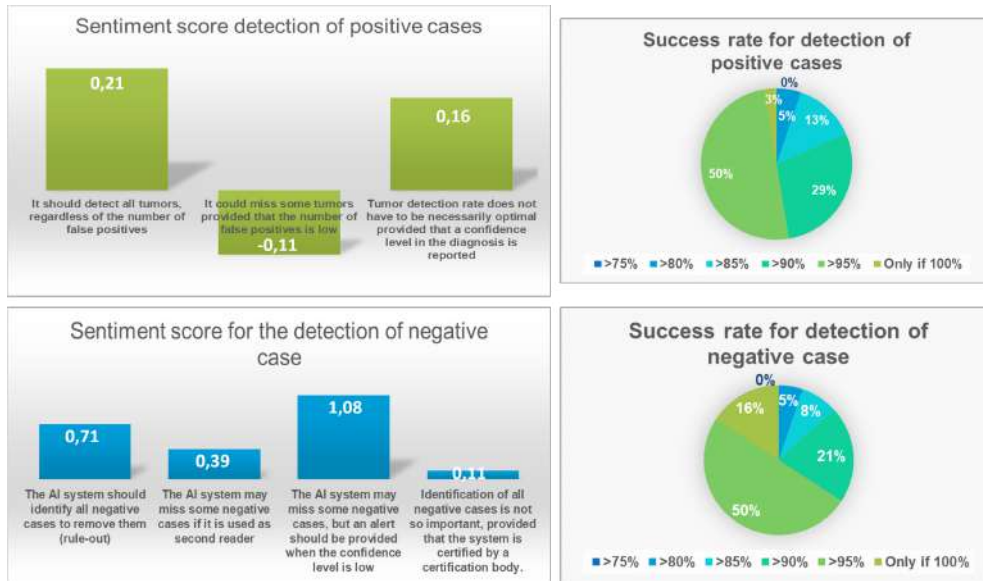


Fig. 11 Expectations in terms of sensitivity and specificity balance and success rates of detection AI tools

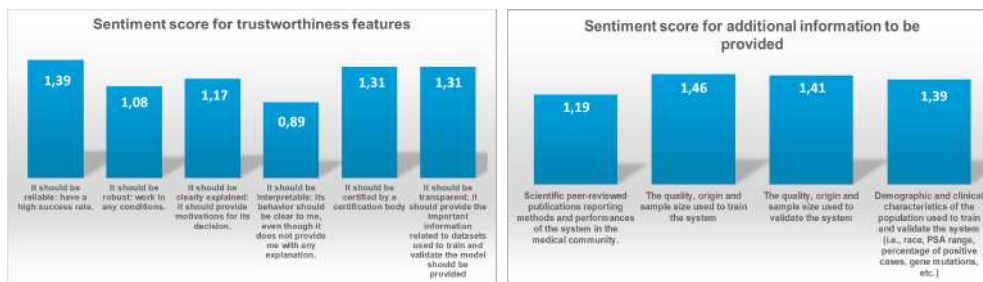


Fig. 12 Desiderata for trustworthiness

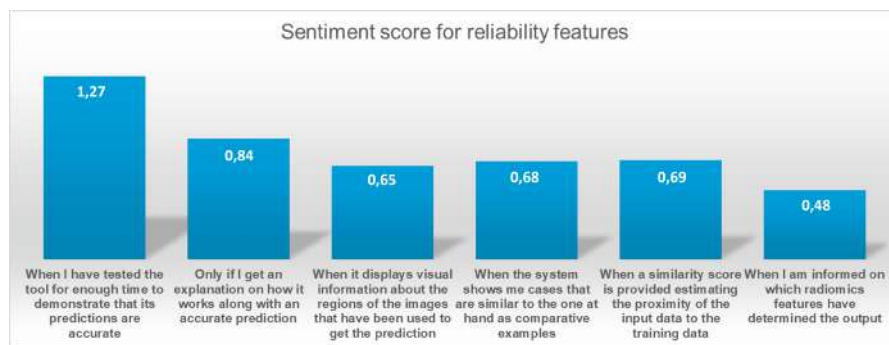


Fig. 13 Desiderata for reliability

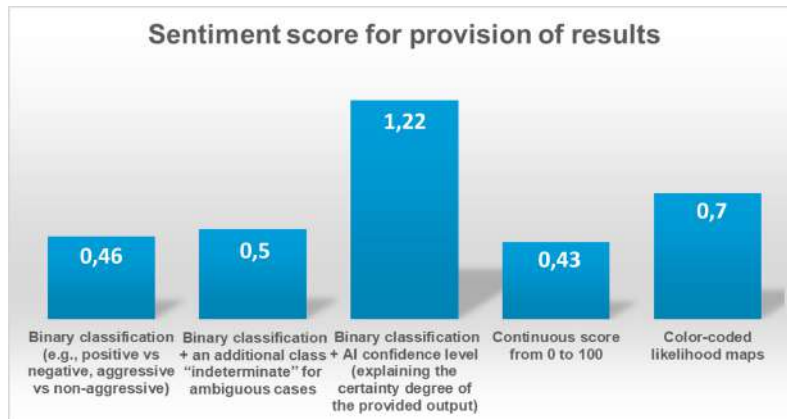


Fig. 14 Desiderata for AI outcome delivery

4 Discussion

This study provides valuable insights into the perspectives of clinicians, specifically radiologists, regarding the integration of AI in healthcare, particularly within cancer imaging. Several key discussion points emerge with important implications for AI development and clinical practice:

- *Human-AI Collaboration*: the strong preference for AI as a "second reader" reinforces the crucial role of human oversight in diagnostic processes. Clinicians view AI as a supportive tool, augmenting their expertise rather than replacing it. This suggests that future AI development should prioritize collaborative models that seamlessly integrate into clinical workflows and empower clinicians with AI-driven insights. The low sentiment towards "rule-out" systems likely reflects concerns about potential errors and medico-legal implications of sole reliance on AI. Future research should explore the optimal balance of human-AI collaboration, investigating how AI can best assist without undermining clinical judgment. Studies comparing the effectiveness and efficiency of different human-AI interaction models are warranted.
- *Trustworthy AI*: the emphasis on reliability, robustness, certification, and transparency as vital features of trustworthy AI systems aligns with broader trends in AI ethics and governance. Clinicians need assurance that AI tools are dependable, accurate, and rigorously validated. Transparency, particularly regarding training data and AI decision-making processes, is paramount for building trust. Future research should explore effective methods for communicating AI limitations and uncertainties to clinicians, including explainability techniques. The preference for on-field accuracy demonstrations highlights the need for real-world evaluations in diverse clinical settings and standardized evaluation frameworks.
- *Explainable AI (XAI)*: the desire for explanations of AI outputs, even without requiring full model interpretability, suggests that clinicians seek to understand the reasoning behind AI-generated results. This is crucial for trust and acceptance. Future research should investigate different XAI techniques to determine which methods are most effective for explaining AI outputs in medical imaging, and how

explanations influence clinicians’ trust and adoption of AI recommendations. This also connects to the reliability concern: if a clinician can understand *why* an AI arrived at a result, they are more likely to trust its reliability.

- *Data Governance*: the high value placed on data transparency, especially regarding training data quality and provenance, highlights the importance of data governance in AI development. Future research should investigate the impact of data bias on AI performance and explore bias mitigation methods in training datasets. Robust data sharing and privacy protection mechanisms are also crucial for developing and evaluating AI models using diverse datasets.
- *Uncertainty Communication*: the preference for AI outputs to include a confidence score and alerts when confidence is low reflects the need for clinicians to assess AI reliability and manage uncertainty. Future research should investigate optimal ways to represent and communicate uncertainty in AI outputs, including deeply exploring different visualization techniques and developing alert systems tailored to specific clinical tasks and integrated into existing workflows.

While this study provides valuable insights into clinician perspectives on trustworthy AI, it also highlights areas for future investigation. As AI in healthcare advances, it will be crucial to consider the evolving legal and regulatory landscape, including the implications of GDPR and the recently approved AI Act. Furthermore, the trade-off between explainability and accuracy in AI models is a complex issue that warrants further exploration. Future research will address these important aspects to ensure the responsible and effective integration of AI-based clinical decision support systems in clinical practice.

In response to input from ProCancer-I clinical partners and result evaluations, the survey underwent revisions. These revisions included the introduction of new queries, such as those identifying the optimal point on the ROC curve, and minor adjustments to current questions, like the one concerning reading and integration methods. This second version of the survey was submitted to the Scientific Boards of the European Society of Oncologic Imaging (ESOI) and the Italian Society of Medical and Interventional Radiology (SIRM), which endorsed its distribution to their members. Additionally, the survey was shared with the AI4HI “AI Validation Working” Group to synchronize efforts with fellow members and establish more defined criteria for assessing trust in clinical validation activities.

We acknowledge that the relatively small sample size of 38 respondents represents a limitation of this study. This limited sample size may hinder the generalizability of our conclusions to the broader population of radiologists. Furthermore, because the respondents were all involved in the ProCancer-I project, a potential bias towards positive attitudes regarding the use of AI in prostate cancer diagnosis and treatment cannot be excluded. However, this potential bias is mitigated by the fact that this work represents a preliminary study—a proof of concept—aimed at developing a robust survey instrument. The primary goal of this initial phase was to refine the survey questions and identify key areas of interest for radiologists involved in AI-related clinical practice. Despite the limited sample size, the responses obtained are representative of the perspectives of the clinical partners involved in the ProCancer-I project,

as we obtained responses from a sufficient number of participants from each partner institution.

5 Conclusion

This study provides a valuable foundation for understanding clinician perspectives on AI in healthcare, specifically within cancer imaging. The findings have significant implications for both AI development and clinical practice.

Regarding the first, developers should prioritize human-centered design, rigorous validation, transparency, explainability, and robust data governance. This includes focusing on collaborative AI models, developing AI for detailed lesion characterization, ensuring data transparency, and providing clear explanations for AI outputs.

Regarding the latter, effective AI integration requires education and training to equip clinicians with the skills to utilize AI tools and interpret AI-generated results. This includes understanding the limitations of AI, interpreting confidence scores, and effectively integrating AI into clinical workflows.

Overall, ongoing dialogue and collaboration between AI developers, clinicians, and stakeholders are essential to ensure ethical, responsible AI implementation aligned with patient and professional needs. This includes developing clear regulatory frameworks and guidelines for AI in healthcare.

This preliminary work has paved the way for the second, more comprehensive, version of the survey that is currently being distributed to a significantly larger and more diverse population of radiologists. Preliminary analysis of the data collected from this expanded survey, which includes a substantially larger number of responses from a more generalized population of radiologists, suggests that the key findings of this initial study are largely consistent, further supporting the validity of our initial observations. This expanded survey will allow for more robust statistical analysis and more generalizable conclusions. We anticipate that the results from the second survey will provide a more comprehensive understanding of radiologists' perspectives on AI in clinical practice and further refine the FUTURE-AI guidelines. This internal survey was instrumental in tailoring the FUTURE-AI guidelines and pinpointing key areas for compliance efforts. The collaborative design activities of the project are still in progress, as the identification of risk sources and points of vulnerability in clinical settings is ongoing. requirements.

Acknowledgements. This study has been partially carried out under the European Union's Horizon 2020 research and innovation program under grant agreement No 952159 (ProCAncer-I) and the Tuscany Region project NAVIGATOR funded and supported by Bando Ricerca Salute Regione Toscana 2018 (DD 15397/2018).

References

- [1] Al Kuwaiti, A., Nazer, K., Al-Reedy, A., Al-Shehri, S., Al-Muhanna, A., Subbarayalu, A.V., Al Muhanna, D., Al-Muhanna, F.A.: A review of the role of artificial intelligence in healthcare. *Journal of personalized medicine* **13**(6), 951 (2023)

- [2] Bajwa, J., Munir, U., Nori, A., Williams, B.: Artificial intelligence in healthcare: transforming the practice of medicine. *Future healthcare journal* **8**(2), 188 (2021)
- [3] Alowais, S.A., Alghamdi, S.S., Alsuhebany, N., Alqahtani, T., Alshaya, A.I., Almo-hareb, S.N., Aldairem, A., Alrashed, M., Bin Saleh, K., Badreldin, H.A., *et al.*: Revolutionizing healthcare: the role of artificial intelligence in clinical practice. *BMC medical education* **23**(1), 689 (2023)
- [4] Henry, K.E., Kornfield, R., Sridharan, A., Linton, R.C., Groh, C., Wang, T., Wu, A., Mutlu, B., Saria, S.: Human–machine teaming is key to ai adoption: clinicians’ experiences with a deployed machine learning system. *NPJ digital medicine* **5**(1), 97 (2022)
- [5] Alanazi, A.: Clinicians’ views on using artificial intelligence in healthcare: oppor-tunities, challenges, and beyond. *Cureus* **15**(9) (2023)
- [6] Choudhury, A.: Factors influencing clinicians’ willingness to use an ai-based clinical decision support system. *Frontiers in Digital Health* **4**, 920662 (2022)
- [7] Robertson, L.J., Abbas, R., Alici, G., Munoz, A., Michael, K.: Engineering-based design methodology for embedding ethics in autonomous robots. *Proceedings of the IEEE* **107**(3), 582–599 (2019)