



Deposited via The University of York.

White Rose Research Online URL for this paper:

<https://eprints.whiterose.ac.uk/id/eprint/211773/>

Version: Published Version

---

**Article:**

Lawton, Tom, Morgan, Phillip, Porter, Zoe et al. (2024) Clinicians risk becoming 'liability sinks' for artificial intelligence. *Future Healthcare Journal*. 100007. ISSN: 2514-6653

<https://doi.org/10.1016/j.fhj.2024.100007>

---

**Reuse**

This article is distributed under the terms of the Creative Commons Attribution (CC BY) licence. This licence allows you to distribute, remix, tweak, and build upon the work, even commercially, as long as you credit the authors for the original work. More information and the full terms of the licence here:

<https://creativecommons.org/licenses/>

**Takedown**

If you consider content in White Rose Research Online to be in breach of UK law, please notify us by emailing [eprints@whiterose.ac.uk](mailto:eprints@whiterose.ac.uk) including the URL of the record and the reason for the withdrawal request.



## Clinicians risk becoming ‘liability sinks’ for artificial intelligence

Tom Lawton<sup>a,b,\*</sup>, Phillip Morgan<sup>c</sup>, Zoe Porter<sup>b</sup>, Shireen Hickey<sup>a</sup>, Alice Cunningham<sup>a</sup>, Nathan Hughes<sup>b</sup>, Ioanna Iacovides<sup>d</sup>, Yan Jia<sup>b</sup>, Vishal Sharma<sup>a</sup>, Ibrahim Habli<sup>b</sup>

<sup>a</sup> Improvement Academy, Bradford Institute for Health Research, Bradford Royal Infirmary, Duckworth Lane, Bradford BD9 6RJ, UK

<sup>b</sup> Assuring Autonomy International Programme, University of York, Heslington, York YO10 5DD, UK

<sup>c</sup> York Law School, University of York, Heslington, York YO10 5DD, UK

<sup>d</sup> Department of Computer Science, University of York, Heslington, York YO10 5DD, UK

### The problem

Artificial Intelligence (AI) is often touted as healthcare’s saviour, but its potential will only be realised if developers and providers consider the whole clinical context and AI’s place within it. One of many aspects of that clinical context is the question of liability.

Analysis of responsibility attributions in complex, partly automated socio-technical systems has identified the risk that the nearest human operator may bear the brunt of responsibility for overall system malfunctions.<sup>1</sup> As we move towards integrating AI into healthcare systems, it is important to ensure that this does not translate into clinicians unfairly absorbing legal liability for errors and adverse outcomes over which they have limited control.

In the current, standard model of AI-supported decision-making in healthcare, electronic data is fed into an algorithm, typically a machine-learned model, which integrates the acquired information to arrive at a recommendation which is output to a human clinician. The clinician can consider this recommendation alongside information from other sources, including examination of and discussion with the patient, and either accept the recommendation as-is, or replace it with a decision they make themselves (Fig. 1). For example, in a system recommending treatment for diabetes, the system may recommend – based on coded electronic data – that it is appropriate to start insulin; though after considering patient context and wishes the clinician may choose to override this. Due to differences in regulatory approval processes, the positioning of such systems as clinical support rather than diagnostic makes them cheaper and quicker to get to market. Additionally, given recent guidance from the National Health Service in England, which clarifies that the final decision must be taken by a healthcare professional,<sup>2</sup> this model looks set to become the norm across the UK healthcare system.

But the standard model may have a negative impact on the clinician, who must choose between accepting the AI recommendation, or substituting their own decision - which, despite probably being AI-influenced, involves largely reverting to a traditional (non-AI) approach. They risk

no longer doing what they are best at, including exercising sensitivity to patient preferences and context, but in effect acting as a sense-check on, or conduit for, the machine. There has been substantial discussion of the cognitive and practical challenges humans face when monitoring automation, such as the additional load of maintaining effective oversight, ensuring sufficient understanding to identify a fault in the system, changes to the way they evaluate information sources, and automation bias.<sup>3,4</sup> For instance the clinician may lack knowledge about the training dataset of the diabetes recommendation system and be unaware that it is less accurate for patients from some ethnic backgrounds; meanwhile, its influence may make the clinician more likely to question their own evaluation. At the same time, the guidance states that the clinician may be held legally accountable for a decision made using the support of AI.<sup>2</sup> Analogous to the way a ‘heat sink’ takes up unwanted heat from a system, the human clinician risks being used here as a ‘liability sink’, where they absorb liability for the consequences of the AI’s recommendation whilst being disenfranchised from its decision-making process, and also having difficult new demands placed on them.

A similar situation exists in driver assistance and self-driving systems for cars, where despite the AI being in direct control of the vehicle, in some jurisdictions it seems the human in the driving seat is already being used as a liability sink. For example, a driver activating self-driving mode typically has to accept that they will take over manual control immediately when required. But in many Tesla collisions Autopilot aborted control less than one second prior to the first impact.<sup>5</sup> This does not give the driver enough time to resume control safely - and yet in practice, for jurisdictions that adopt fault based systems of liability for motor vehicle accidents such the UK, it is likely that they would be liable for the accident. As the most obvious ‘driver’ close to where AI is used in a clinical setting, the clinician could easily end up being held similarly liable for harmful outcomes from AI-based decision-support systems, and carrying this stress and worry, but having limited practical control over their development and deployment, or understanding of how the AI recommendations are reached.<sup>6</sup>

This article reflects the opinions of the author(s) and should not be taken to represent the policy of the Royal College of Physicians unless specifically stated.

\* Corresponding author.

E-mail address: [tom.lawton@bthft.nhs.uk](mailto:tom.lawton@bthft.nhs.uk) (T. Lawton).

Social media:  (T. Lawton)

<https://doi.org/10.1016/j.fhj.2024.100007>



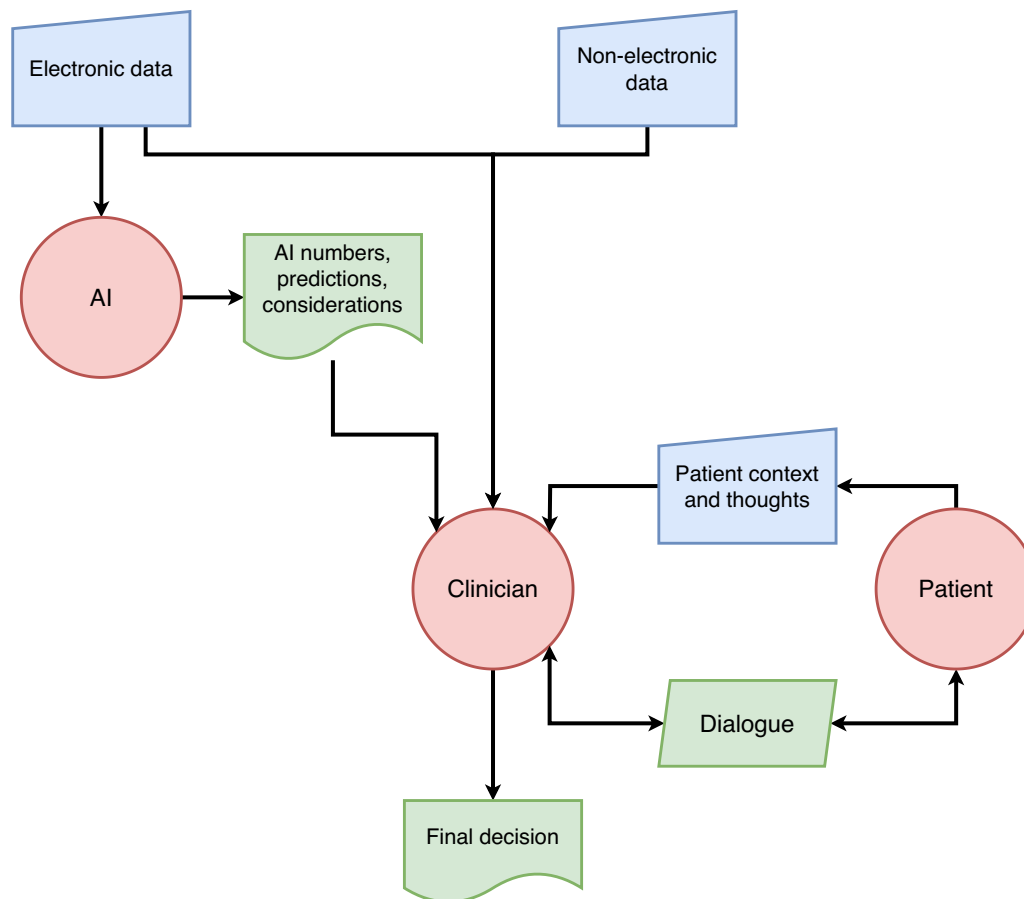


Fig. 2. AI model with alternative outputs to inform patient/clinician dialogue.

icians to seek such contribution. Thus, in practical terms with systems of this type the clinician remains liable for acting on the recommendations or decisions of an AI they do not and cannot fully understand. Facing the stress and worry of the consequences of using it, many clinicians may refuse to accept the risk, and simply turn off the machine.

### Alternative models

Pooling risk might prevent the clinician becoming a liability sink, but the prevalent model may have other drawbacks for the clinician, the patient, and the system as a whole. Fig. 1 shows that the entire input of the patient and clinician into the decision is restricted to either accepting the AI's recommendation, or - for this case - rejecting it (effectively switching it off and returning to standard practice, albeit likely influenced by the AI's recommendation). This is at odds with the goal of patient-centred decision-making,<sup>23</sup> as the AI cannot easily incorporate patient context and ideas, concerns, and expectations itself - this context is only added by the clinician choosing to accept or replace the AI's output. It may also be frustrating for the clinician by eroding their ability to do what they do best: integrating clinical science and patient context in a dialogue to come to a shared decision.

Fortunately, this is not the only possible approach. Rather than restructuring systems in a clinical setting around an AI designed to work this way, it may be preferable to explore alternative models which give greater focus to the patient and clinician.<sup>24</sup> In some of these models the AI may not even give a decision or recommendation, but instead show predictions of the effect of different decisions (e.g. treatment options), or highlight data that is most relevant to the AI model in its decision making. In this way, the explanation of an explainable-AI system may be more useful than the decision or recommendation itself.<sup>25,26</sup> Fig. 2 shows a model where these alternative outputs from the AI system in-

form a dialogue between the clinician and patient, leading to a decision. Whilst the model refers to complex AI systems which cannot be directly interpreted by anyone, including the clinician, this is clearly analogous to existent non-AI systems such as automated ECG analysis where the inner workings are not easily available to the clinician.

In the diabetes example, the outputs may be predictions for each treatment option of outcomes such as blood test results or other end-points such as the risk of a heart attack, forming the basis for the dialogue and subsequent decision, and the AI may not output a direct treatment recommendation at all. Most current AI radiology systems are similar - providing information to a reporting clinician to highlight areas and possible diagnoses without directly completing the report. There are vendors attempting to take the clinician out of the loop, but presently systems can only take on a small proportion of the workload.<sup>27</sup> As these truly autonomous systems advance, without a nearby clinician liability sink, they may well test some of the legal issues discussed above.

In Fig. 3, a more advanced AI system communicates directly with the patient and a three-way dialogue proceeds before a decision emerges. A year ago, dialogue with an AI capable of explaining itself to patients might have been considered fanciful, but advances in Large Language Models employed in tools like ChatGPT have made them seem very plausible. A diabetes system built this way might be capable of eliciting the patient's thoughts and concerns about the difficulties of starting insulin. It could provide a tailored approach that does not lose the patient voice, and provide an explanation to the clinician in more the manner of discussion with a multidisciplinary team member. Other models can be conceived along these lines, bringing the patient and clinician back into the decision-making focus.

With both models in Figs. 2 and 3, the clinician retains the final decision as recommended by NHS England. Are they still a liability sink for the AI? The models may not remove liability, but we would argue

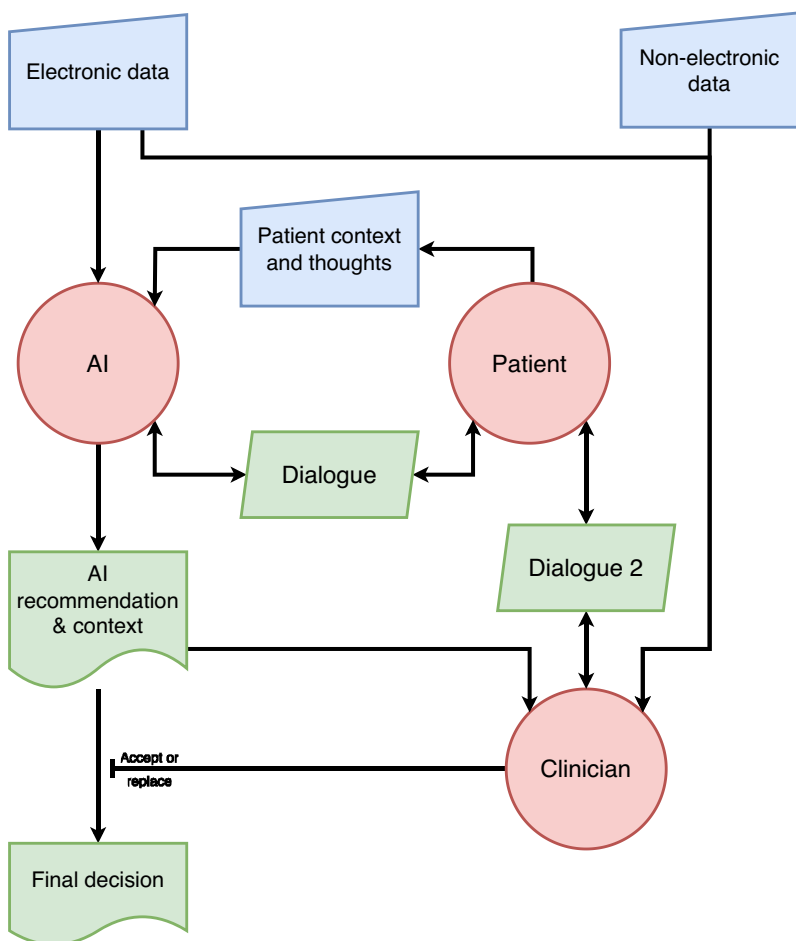


Fig. 3. Advanced AI model capable of sustaining dialogue with the patient.

that the clinician’s role here is much more traditional, as they are integrating a variety of data and opinions in a manner of working that has become familiar with the advent of the multidisciplinary team.<sup>28</sup> Clinicians should feel much more comfortable in accepting liability for a decision where they have genuine understanding and agency, and the socio-technical system as a whole will be much more acceptable to both clinicians and patients as it retains compatibility with patient-centred care.

The question remaining in this setup, however, is the assignment of liability for defective AI advice or information. As these models return the clinician to a more traditional role, the current legal position becomes more appropriate: treating the AI as a standard medical device. This could be dealt with via product liability, suitably adjusted to take into account the problems within such regimes as applied to AI systems, such as proof of causation, and the failure discussed above of the PLD<sup>19</sup> to cover unembodied software. The European Union has recognised this need, and published reform proposals for the PLD. If we do not want clinicians to become liability sinks, similar reforms may be needed in the United Kingdom.

In summary, AI systems being developed using current models risk using clinicians as ‘liability sinks’, absorbing liability which could otherwise be shared across all those involved in the design, institution, running, and use of the system. Alternative models can return the patient to the centre of decision-making, and also allow the clinician to do what they are best at, rather than simply acting as a final check on a machine.

**Summary**

- The benefits of AI in healthcare will only be realised if we consider the whole clinical context and the AI’s role in it.

- The current, standard model of AI-supported decision-making in healthcare risks reducing the clinician’s role to a mere ‘sense check’ on the AI, whilst at the same time leaving them to be held legally accountable for decisions made using AI.
- This model means that clinicians risk becoming ‘liability sinks’, unfairly absorbing liability for the consequences of an AI’s recommendation without having sufficient understanding or practical control over how those recommendations were reached.
- Furthermore, this could have an impact on the ‘second victim’ experience of clinicians.
- It also means that clinicians are less able to do what they are best at, specifically exercising sensitivity to patient preferences in a shared clinician-patient decision-making process.
- There are alternatives to this model that can have a more positive impact on clinicians and patients alike.

**Declaration of competing interest**

TL has received an honorarium for a lecture on this topic from Al Sultan United Medical Co and is head of clinical artificial intelligence at Bradford Teaching Hospitals NHS Foundation Trust, and a potential liability sink

All other authors report no conflicts of interest

**CRediT authorship contribution statement**

**Tom Lawton:** Conceptualization, Funding acquisition, Writing – original draft, Writing – review & editing, Formal analysis, Visualization. **Phillip Morgan:** Writing – original draft, Formal analysis, Visualization, Writing – review & editing. **Zoe Porter:** Conceptualization,

Funding acquisition, Writing – original draft, Writing – review & editing, Formal analysis, Visualization. **Shireen Hickey:** Writing – original draft, Formal analysis, Visualization, Writing – review & editing. **Alice Cunningham:** Formal analysis, Visualization, Writing – review & editing. **Nathan Hughes:** Formal analysis, Visualization, Writing – review & editing. **Ioanna Iacovides:** Formal analysis, Visualization, Writing – review & editing. **Yan Jia:** Formal analysis, Visualization, Writing – review & editing. **Vishal Sharma:** Formal analysis, Visualization, Writing – review & editing. **Ibrahim Habli:** Conceptualization, Funding acquisition, Writing – original draft, Writing – review & editing, Formal analysis, Visualization.

## Acknowledgements

This work was supported by The MPS Foundation Grant Programme. The MPS Foundation was established to undertake research, analysis, education and training to enable healthcare professionals to provide better care for their patients and improve their own wellbeing. To achieve this, it supports and funds research across the world that will make a difference and can be applied in the workplace. The work was also supported by the Engineering and Physical Sciences Research Council (EP/W011239/1).

## References

- Elish MC. Moral crumple zones: cautionary tales in human-robot interaction. *Engag Sci Technol Soc.* 2019;5:40–60 Mar 23.
- NHS England. Information governance guidance: artificial intelligence [Internet]. NHS England - transformation directorate; 2022 [cited 2022 Nov 3]. Available from: <https://transform.england.nhs.uk/information-governance/guidance/artificial-intelligence/>.
- Bainbridge L. Ironies of automation. In: Johannsen G, Rijnsdorp JE, eds. *Analysis, Design and Evaluation of Man-Machine Systems* Pergamon; 1983:129–135.
- Parasuraman R, Sheridan TB, Wickens CD. A model for types and levels of human interaction with automation. *IEEE Trans Syst, Man, Cybernet - Part A: Syst Humans.* 2000;30(3):286–297 May.
- Engineering Analysis 22-002 [Internet]. National highway traffic safety administration, office of defects investigation; 2022 [cited 2022 Nov 3]. Available from: <https://static.nhtsa.gov/odi/inv/2022/INOA-EA22002-3184.PDF>.
- Habli I, Lawton T, Porter Z. Artificial intelligence in health care: accountability and safety. *Bull World Health Organ.* 2020;98(4):251–256 Feb.
- Wu AW, Steckelberg RC. Medical error, incident investigation and the second victim: doing better but feeling worse? *BMJ Qual Saf.* 2012;21(4):267–270 Apr 1.
- Sirriyeh R, Lawton R, Gardner P, Armitage G. Coping with medical error: a systematic review of papers to assess the effects of involvement in medical errors on healthcare professionals' psychological well-being. *Qual Saf Health Care.* 2010;19(6) Dec 1e43–e43.
- Engel KG, Rosenthal M, Sutcliffe KM. Residents' responses to medical error: coping, learning, and change. *Acad Med.* 2006;81(1):86–93 Jan.
- Gianfrancesco MA, Tamang S, Yazdany J, Schmajak G. Potential biases in machine learning algorithms using electronic health record data. *JAMA Intern Med.* 2018;178(11):1544–1547 Nov 1.
- McDermid JA, Jia Y, Porter Z, Habli I. Artificial intelligence explainability: the technical and ethical dimensions. *Philosoph Trans R Soc A: Mathemat, Phys Eng Sci.* 2021;379(2207):20200363 Aug 16.
- Chesterman S. Artificial intelligence and the limits of legal personality. *ICLQ.* 2020;69(4):819–844.
- Smith H, Fotheringham K. Artificial intelligence in clinical decision-making: rethinking liability. *Med Law Int.* 2020;20(2):131–154 Jun 1.
- Wilsher v Essex area health authority [1987]QB 730 (CA). 1987.
- Junior v McNicol. Times law reports, March 26 1959. 1959.
- Armitage M, editor. Chapter 10: persons professing some special skill. In: Charlesworth & Percy on Negligence. 15th ed. London: Sweet & Maxwell; p. 10–147. (Common Law Library).
- Commission E, Justice DG for, Consumers. *Liability for Artificial Intelligence and Other Emerging Digital Technologies.* Publications Office; 2019.
- Morgan, Phillip. Chapter 6: tort law and artificial intelligence – vicarious liability. In: Lim E, Morgan P, editors. *The Cambridge Handbook of Private Law and Artificial Intelligence.* Cambridge University Press;
- Council Directive 85/374/EEC of 25 July 1985 on the approximation of the laws, regulations and administrative provisions of the Member States concerning liability for defective products [Internet]. OJ L Jul 25, 1985. Available from: <http://data.europa.eu/eli/dir/1985/374/oj/eng>.
- Burton S, Habli I, Lawton T, McDermid J, Morgan P, Porter Z. Mind the gaps: assuring the safety of autonomous systems from an engineering, ethical, and legal perspective. *Artif Intell.* 2020;279:103201 Feb 1.
- Heywood R. Systemic negligence and NHS hospitals: an underutilised argument. *King's Law J.* 2021;32(3):437–465 Sep 2.
- Abbott R. *The Reasonable Robot: Artificial Intelligence and the Law* [Internet]. Cambridge: Cambridge University Press; 2020 [cited 2023 Feb 22]. Available from: <https://www.cambridge.org/core/books/reasonable-robot/092E62F0087270F1ADD9F62160F23B5A>.
- Bjerring JC, Busch J. Artificial intelligence and patient-centered decision-making. *Philos Technol.* 2021;34(2):349–371 Jun 1.
- Birch J, Creel KA, Jha AK, Plutynski A. Clinical decisions using AI must consider patient values. *Nat Med* [Internet]. 2022. Jan 31 [cited 2022 Feb 1]; Available from: <https://www.nature.com/articles/s41591-021-01624-y>.
- Jia Y, McDermid JA, Lawton T, Habli I. The role of explainability in assuring safety of machine learning in healthcare. *IEEE Trans Emerg Top Comput.* 2022 1–1.
- Mittelstadt B, Russell C, Wachter S. Explaining explanations in AI. *Proceedings of the Conference on Fairness, Accountability, and Transparency* [Internet]. New York, NY, USA: Association for Computing Machinery; 2019:279–288. doi:10.1145/3287560.3287574.
- Park CMAI. Workload reduction by autonomous reporting of normal chest radiographs. *Radiology.* 2023;307(3):e230252 May.
- Epstein NE. Multidisciplinary in-hospital teams improve patient outcomes: a review. *Surg Neurol Int.* 2014;5(Suppl 7):S295.