

Making AI Infused Products and Services more Legible.

Franziska Pilling PhD Candidate, Lancaster University, f.pilling@lancaster.ac.uk, ORCID: 0000-0001-5324-0575.

Haider Ali Akmal Research Associate, Lancaster University, h.a.akmal@lancaster.ac.uk, ORCID: 0000-0001-9578-3578

Joseph Lindley Lecturer in Design Research, Lancaster University, j.lindley@lancaster.ac.uk, ORCID: 0000-0002-5527-3028

Adrian Gradinar Lecturer in Smart Home Futures, Lancaster University, a.gradinar@lancaster.ac.uk ORCID: 0000-0001-5416-6207

Paul Coulton Professor of Speculative Design, Lancaster University, p.coulton@lancaster.ac.uk, ORCID: 0000-0001-5938-4393

Abstract

The increasing availability of large data sets has initiated a resurgence in Artificial Intelligence (AI) research. Today AI is integrated into a wide variety of so-called smart products and services to personalize user experiences. Smart Technologies are typically designed for ease of use, with their complex underlying procedures (intentionally) obfuscated. This leads to a lack of legibility, misconceptions about how AI works, even AI illiteracy. Using a Research through Design (RtD) approach, the authors address the challenge of AI legibility through iconography to enhance agency and understanding of AI and data-infused interactions.

A Process of Obfuscation: the Challenge

AI and data collection have become core activities within the cloud-based services empowering smart thermostats, streaming services, and AI assistants, such as Alexa, and are becoming increasingly ubiquitous in our daily activities. However, the underlying operations relating to AI and data collection and processing in and by these networked products are predominantly illegible such as the user's voice data being recorded to train AI assistants. The illegibility of AI operations diminishes user agency and their ability to negotiate the transactional nature of various human-machine activities, as obtaining functionality is often conditional on agreeing to provide personal data [1]. Designs that intentionally exploit the epistemologically indeterminate terrain of human and algorithmic logics, obfuscate AI operations for a variety of reasons: under the aim of being human-centred they simplify the

interaction to avoid overloading the user with auxiliary information [2]; for institutional protection; to conceal intellectual property; and, in some cases, to deceptively [3] collect data without explicit consent [4]. While the 2018 Cambridge Analytica scandal (which revealed that the personal data of over 87 million Facebook users had, from 2010 onwards, been collected and used for political advertising purposes), made many Facebook users more aware that data produced by their interactions with Facebook is archived and traded, it is less obvious *how* AI is processing the data and to what purpose, whether by Facebook, or the third parties Facebook supplies the data to. This lack of legibility presents many challenges, resulting in users' perception of AI being more influenced by totalitarian political narratives or science-fiction renderings of AI, than the mundane reality of narrow AI using machine learning embedded in services such as Netflix to predict viewing preferences. This dichotomy, usually referred to as the “definitional dualism of AI” [5], highlights part of the challenge of making the functional elements of AI more legible to users. While the notion of legibility is promoted in many frameworks to encourage ‘better’ strategies for AI implementation [6], there are no examples of how this is achieved in practice. This article discusses several concrete, practical suggestions for achieving (greater) legibility in human-AI interaction.

Researching AI Legibility through Design

To address the challenge of how AI legibility might be achieved in practice, we adopted a Research through Design (RtD) approach and developed a set of AI icons to communicate AI operations, to diffuse the complexity and signification indeterminacy, and raise user awareness of AI functions. As a practice-based approach, RtD provides a generative aptitude to “explore, speculate, particularise, and diversify” the challenge in hand towards manifesting findings into rich research artefacts [7]. To this end, our research interweaved methods from diverse disciplinary perspectives: human-computer interaction for theories of richer inter-relationships between users and computers [8]; semiotics for modelling the constructs of signifiers [9]; amalgamating design and AI research in tasks such as creating trust in AI products – given the bias issues prevalent in training data [10; 11]; and promoting legibility and explainability of AI functions over AI “transparency”, which is far more related to making AI products transparent to and auditable by specialists, than to everyday users [12].

Instead of focussing on text-based descriptions of AI, which tend to use technical jargon and/or require expert knowledge, we assessed the current levels of AI legibility by surveying AI iconography [13] as the visualisation of information is often used to increase accessibility. Our research showed that, although some AI imagery attempts to represent the underlying system such as a neural network (Fig. 1a) or highlighted its use, such as in face detection (Fig. 1b), the vast majority play into AI's definitional dualism by showing human-like machines (Fig. 1c & d), thus exacerbating misconstrued perceptions of human intelligence existing in AI. The survey also highlighted that current AI imagery rarely communicates the intricacies of how an AI functions and in what context, emphasising the need to develop a new visual approach to enhance AI legibility.

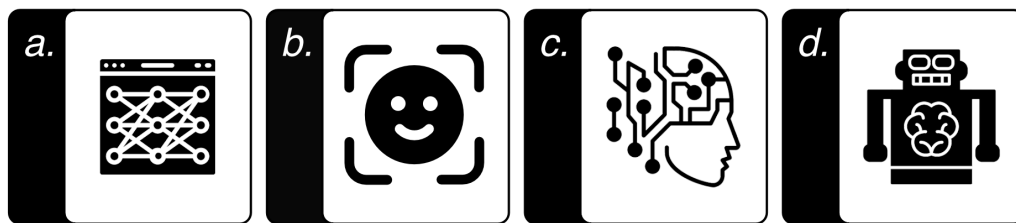


Figure 1. Examples of current AI iconography (a. © Noun Project. Vector image: Ben Davis) (b. © Noun Project. Vector image: Adrien Coquet) (c. © Adobe Stock. Vector image: Priyanka) (d. © Noun Project. Vector image: Med Marki).

Our initial design response resulted in seventeen graphical icons (Fig. 2), which individually and in their grouping detail a particular operational function or feature we identified as crucial rudimentary components to communicate. These include: *Learning Scope*, which communicates how the AI adapts, and is a fundamental factor for human-AI interaction [14]; *Data Provenance*, which declares the source of the training data, as data reflects the AI and its trustworthiness [15]; *Processing Location*, which defines where processing is occurring, and impacts users' perception of accountability [16]; and *Training Data Type*, which is a granular account of the type of data used to train the AI to reduce opacity, bias and increase trust. The final Icon is *Intrinsic Labour*, a critical reflection of the monetisation of data through the commodification of users and their interactions. These icons illuminate an AI's operational activity without engaging with the specifics of the AI's implementation, such as using neural nets.

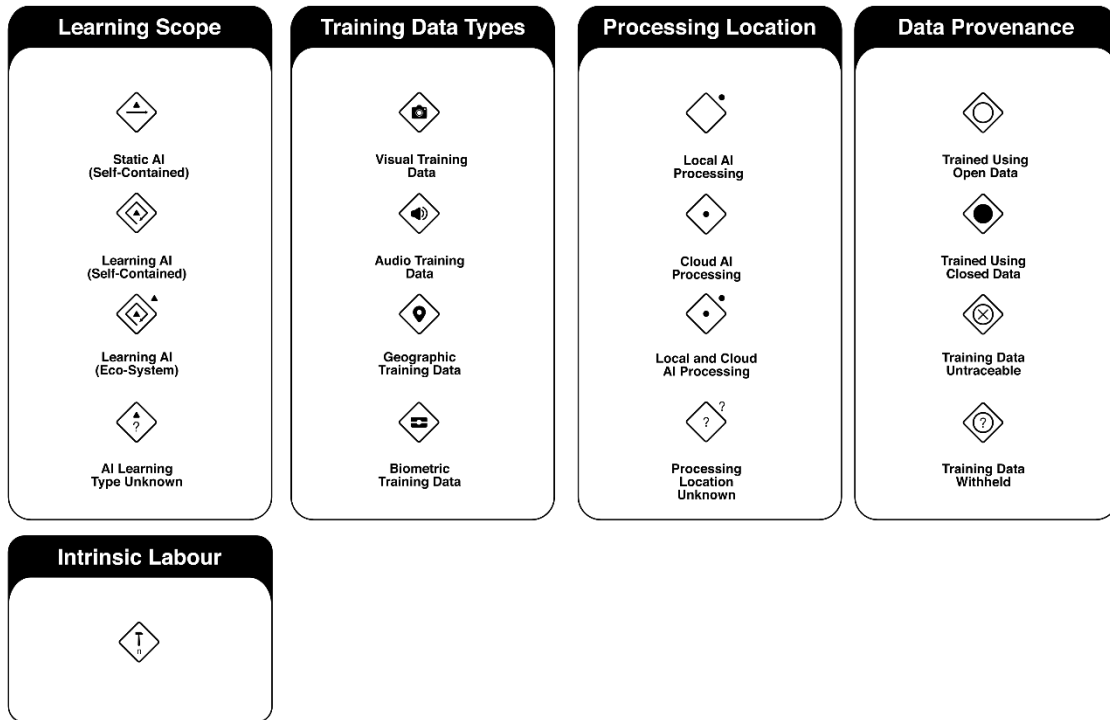


Figure 2. The first iteration of icons (© Pilling, Akmal, Lindley, Coulton).

These icons can then be used in contextualized combinations to map and communicate an AI’s “ontological constituent” and articulate an AI’s implementation to the user [17]. We designed three different icon styles during our ideation: pictorial, textual, and abstract variants. Our initial evaluation settled on the abstract icons as they best avoided the pitfalls previously discussed. Further, these abstract icons conformed to the principles of semiotics, where icons often hybridize symbolic, indexical, and iconic categories to communicate the intended concept. These icons were designed with a diamond shape to retain uniformity between icons and facilitate their combinations in different configurations. The abstract AI icons were also inspired by laundry care labels initially created in 1971; while we may not always take notice of these abstract icons, they provide a means of making legible how we can most easily maintain a working relationship with our clothes.

An icon works if the user can match the interpretant to the intended object—or, in the instance of a digital thing, which is arguably more challenging to capture having no “conventional representation” [18]—its concept or the implication [19]. However, through semiotics, a plausible representation can emerge [20] and can be embedded within an icon to become “iconic” [21] and be utilized by users. To ascertain how well the icons performed, we

engaged in iterative and interactive evaluation workshops with potential stakeholders, throughout 2020, including end-users, academics and industry practitioners who deploy AI.

Testing Intuitiveness and Icon Development

Due to Covid, the workshops were conducted online using a bespoke set of playful activities produced with the game **engine Godot**, which is an open-source and cross platform games engine that facilitates designers to code and build 2D and 3D games using a simplified variant of Python (a programming language).

A playful approach allowed us to ease participants of all knowledge levels towards discussing potentially complex ideas outside their field of expertise or even experience. The initial workshops evaluated the icons with forty-seven participants. In the following paragraphs, we note our observations and how these observations influenced the creation of the second set of icons. The first exercise of the workshop was matching the icons to their descriptors. It quickly became apparent that the *Training Data Types* and *Processing Location* icons were the most intuitive, with the highest number of correct matches. The success of the *Training Data* icons was due to the fact that they contained well-known “iconic” signifiers, such as an audio speaker for audio training data, and “symbolic” signifiers [22] that have become embedded in our society, such as the geographic pin for geographic data.

As all seventeen icons were present in one setting, participants developed non-verbal reasoning tactics to match icons and their textual descriptions and then use the resemblance to one another to collate into the corresponding groupings. For the *Processing Location* icons, participants commented on the fact that they were able to decipher these individually and group them as the icons employed the symbolic element of a circle to represent processing and were specifically positioned either inside or out of the aforementioned ‘AI diamond’ to represent internal or cloud processing. From the results of the workshop discussions, we identified the problematic aspects of some icons. For example, to communicate an AI trained once offline, the icon *Static AI* was presented with a triangle used to symbolize learning, beneath a directional arrow pointing to the right. Many participants observed that the arrow suggested movement rather than stasis. This arrow was changed to a triangle enclosed by a diamond shape that sat inside the icon’s AI diamond for the second icon iteration. This configuration better conformed to the group’s symbology where an open arrow path in a

diamond shape symbolized continuous learning; hence a closed diamond accentuated static (Fig. 3).

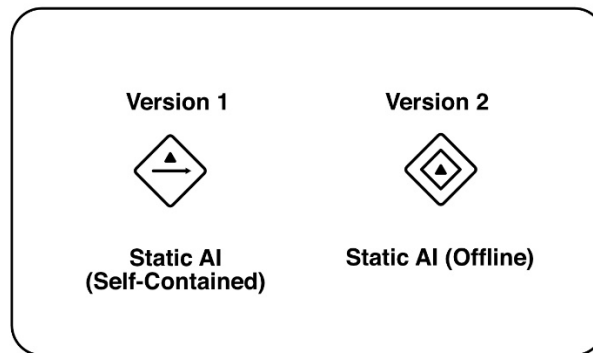


Figure 3. Comparison of Static AI icons (© Pilling, Akmal, Lindley, Coulton).

Framing AI Concepts

The workshops also tested how well the AI's relationship to data was being communicated and how we framed these concepts. We attempted to design the icons to be expansive enough yet not ambiguous. However, this proved problematic in the way we sought to frame the concepts drawn from discussions amongst those working in AI research. Consequently, the *Data Provenance* category was redefined as *Training Data Origin* as this proved more understandable during workshop discussions. Additionally, the training definitions we initially framed were found to be vague and beyond the scope of knowledge for everyday users because of the specialist terminology we used. For instance, the concept *Trained Using Open Data* where participants often asked what open data meant, to which facilitators would answer "data to be audited by an external body, to determine whether the data is representative of the activity it is being applied to" (workshop facilitator). Therefore, in the second iteration, we reframed what specific icons were communicating to be more accessible for general users. In particular, 'open data' changed to *Trained Using Auditable Data*, and most importantly, we created an icon for *Trained using User Data* which participants were often most concerned about.

The workshops and the icons as tools brought about supplementary observations regarding how users developed a better understanding of AI operations, as many participants remarked that they had more knowledge and critical awareness of AI technology and data gathering techniques after completing the workshop. This comprehension led to a handful of

participants claiming they would immediately disable the location permissions on their devices. We also identified which icons and their concepts needed to be self-explanatory. For example, the concept *Learning AI (Eco-system)* was confusing and required an understanding of what the ecosystem metaphor meant in the context of AI learning; thus, we framed the concept *AI to AI Learning*.

Developing New AI concepts

During the first workshops, it became apparent that the AI's overall application was not communicated, with participants seeking this information. Consequently, the category of *AI-Assisted Decisions* was designed for the second iteration, expanding the number of icons to Twenty-One (Fig 4.). This icon set proved problematic to design an abstract pattern due to existing understandings associated with particular terms; for example, using a crystal ball to signify prediction would play into the saturated discourse on technology and magic [23], which we wanted to avoid. The icons we settled with the use of letters, although these are not ideal for translation into other languages.

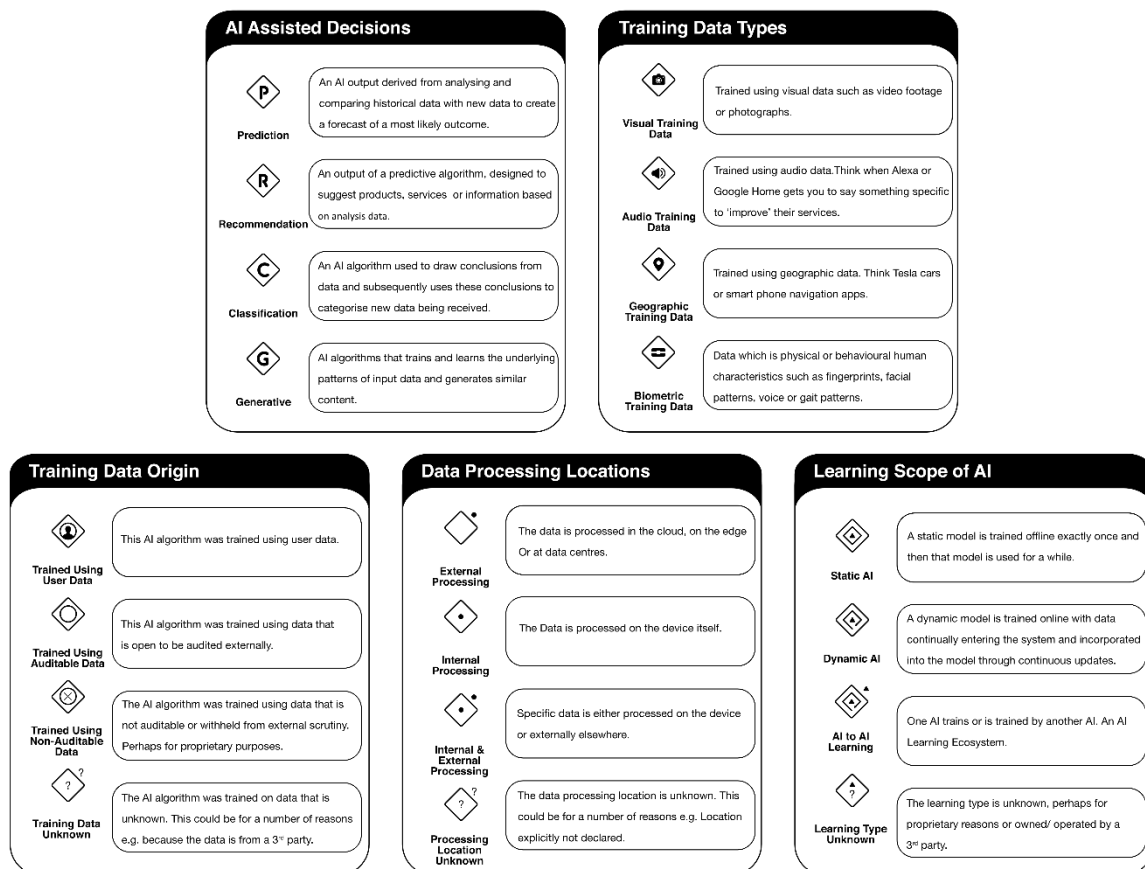


Figure 4. Version two of AI icons (minus Intrinsic Labour) (© Pilling, Akmal, Lindley, Coulton).

Designing a User Priority Arrangement

Discussing with participants the information they felt was most relevant to them resulted in speculating a hierarchical order for the icons. First, noting the ‘presence of AI’ with a subsequent breakdown of the AI in question and be in an order that conformed to users’ ‘priority of information’ (Fig .5). In this way, users could quickly access the most relevant information, and for users wanting a technical understanding, this would be detailed further down the hierarchy. Thus, a user can decide how much information they would need to make a conscious decision on how to interact with an AI device. Together, the icons abate the indeterminacy of interaction, and the hierarchy uniformly organises information in a way that further diminishes ambiguity, however at a comfortable depth of detail for the user. In the second set of the workshops, we asked participants to rate the icons based on information they wanted to know about their device to establish an order of importance. Interestingly, the *Training Data Types* as a category ranked the most important, with participants particularly wanting to understand what type of personal data an AI was recording and learning from. Correspondingly, *Trained Using User Data* was the common highest choice, whereas the remaining *Training Data Origin* category and *AI Learning Type* were considered less important to know and be placed further down the hierarchy.

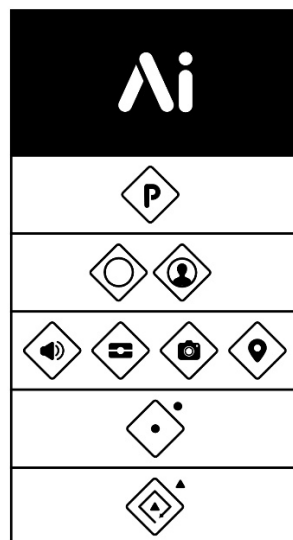


Figure 5. An example of the AI Hierarchy System (© Pilling, Akmal, Lindley, Coulton).

Second Iteration Performance

The second iteration of icons had a high success rate; it was seen as intuitively matched to the textual descriptors, with several participants matching all correctly. Granted, the participants had both the concept and the icon to work with, rather than arbitrarily guessing what the icon could mean. However, the empirical testing we chose to conduct reflects the reality of launching icons out into the ‘wild,’ where the semiotic design of the icons establishes a visual vocabulary for an immaterial concept. Over time these icons would become concretized in meaning and familiarized due to their tested and considered design. The icons that proved to be less intuitive (based on forty-five participants results) were the more abstract icons of *AI Learning Type* and *Training Data Origin*. Nevertheless, participants who did match these icons commented that crucial parts of the icons could be identified and considered a visual representation. For instance, the rotational arrows in *AI Learning types* resembled “a round of training or an iteration of learning” (anonymous workshop participant). It must be recognized that while the icons can facilitate a greater understanding and legibility of AI operations, for many systems, it would be dubious to claim that a full, in-depth explanation of how a particular decision was reached, is possible. AI computing is essentially un-interpretable. Although AI processes attempt to determine real-world indeterminacies, indeterminacy is also intrinsic and internal to the computing process itself, its iterative and recursive processes [24]. It is only what is already known that can be made legible and communicated.

Unaccounted Concepts

With the experience of the above-mentioned workshops and the rapid advancement of AI research, we have begun designing accompanying icons that mirror current ethical and responsible considerations (Fig 6). The icon *Intrinsic Labour* offered the opportunity to critically discuss and question the larger impacts of using AI technology, which is too often far removed from the immediate interactions and outputs a user experiences when using AI-assisted devices. Icons can address the notion of responsibility, enabling us to make better choices and empower responsible technology.

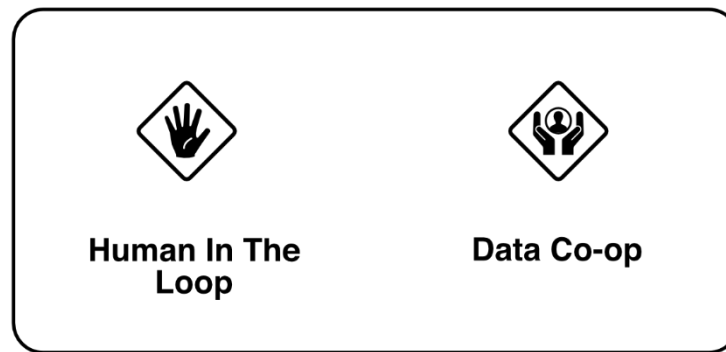


Figure 6. Test examples for responsible AI icons (© Pilling, Akmal, Lindley, Coulton).

The first developmental icon is *Human in-the-Loop*, signifying that some form of human labour occurred for the AI to operate. This icon could acknowledge the oft misconception that an AI operation is purely performed through computation. The reality is that a human's hand, expertise and sensitivity is still an integral part of the automation revolution with activities such as data labelling [25]. 'Human Out-of-the-Loop' systems can evolve beyond human intelligibility and perform tasks in a way that is literal, although ultimately incorrect by human standards. Frank Lantz's 2017 *Universal Paperclips* shines a light on this. In this game, the user plays the role of an AI programmed to produce paperclips. They first click on a box to create a single paperclip at a time. This is followed by options to sell paperclips to finance machines that produce huge quantities of paperclips automatically, without human intervention. The game is based on exponential growth; the user is continually required to invest money, material and immaterial resources, and computer cycles to invent ever-new paperclip-producing ideas in order to move to the next phase of growth. The game ends when the AI succeeds in converting the entire universe into paperclips, and destroying the world [26]. An icon of this prominence would aid in developing more critical awareness and scope to fortify workers' protection within these technological sectors through policy. It would also educate users about AIs' accurate foundations, guiding users to opt for more coherently programmed or 'rational' AIs.

While the second concept is still in its preliminary stages of development, we are keen to create an icon, or set, that promotes the ethical use of intimate data that has potential benefits for the common good, such as patient data to prevent or predict disease. Ethically approved AI applications or services could be certified with these *Cooperative Icons* to communicate to users that their data is used for good on neutral terms, which would embody trust in interacting with these services.

Conclusion

The research presented here is not expected to solve the evolving challenge of AI legibility but to demonstrate the potential of design-led responses to the problem of legibility and illustrate our RtD methodology of interweaving interdisciplinary perspectives to design graphical AI icons to communicate AI's intangible, complex functions for more informed use. We outlined the developing iterations of the icons to demonstrate the multifaceted challenges of legible AI. This is achieved by both acknowledging and steering away from AI's essentially indeterminate processes, which rely on machine-machine interaction, regular updates and additional information sourcing, as well as on iteration, recursion, and thus also change, inherent in computational processes. Rarely can we say for certain why an AI has reached a particular decision, even with the aid of expert knowledge. At the same time, it is important to limit the indeterminacy of information, modes of communication, and semantics, in order to pave the way for general AI literacy as well as make human-AI interaction design accountable. The design challenge for legible AI is to create an accessible and uniformly constructed AI lexicon not only to demystify AI, and the effect this mystification may have on users (confusion, uncertainty and/or erroneous use), but also to make it possible for the general user to understand and assimilate future AI developments while remaining aware and, where needed, critical of intentional or unintentional obfuscation of AI processes.

Acknowledgements

This research has been supported through the Ucanny AI project funded through the PETRAS National Centre for Excellence in Cybersecurity EPSRC project EP/S035362/1.

References

1. Richard Mortier et al , "Human-Data Interaction: The Human Face of the Data-Driven Society," *ARxiv* Abs/1412.6159 (2014).
2. Donald Norman. *The Invisible Computer: Why Good Products Can Fail, the Personal Computer is So Complex, and Information Appliances are the Solution* (MIT Press; 1998).

3. Jenna Burrell, “How the machine ‘thinks’: Understanding opacity in machine learning algorithms,” *Big Data & Society* 3. No. 1 (2016). DOI: 10.1177/2053951715622512
4. Shoshana Zuboff. *The Age of Surveillance Capitalism* (London. Profile Books Ltd; 2019).
5. Joseph Lindley et al “Researching AI Legibility through Design,” *CHI '20: Proceedings of the 2020 CHI Conference on Human Factors in Computing Systems* (2020) pp.1-13. DOI: 10.1145/3313831.3376792
6. Jessica Fjeld et al “Principled Artificial Intelligence: Mapping Consensus in Ethical and Rights-Based Approaches to Principles for AI,” *Berkman Klein Center Research Publication*, No. 1 (2020). DOI: 10.2139/ssrn.3518482
7. William Gaver, “What should we expect from research through design?,” *CHI '12: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems* (2012) pp. 937-946. DOI: 10.1145/2207676.2208538
8. John Bowers, Tom Rodden, “Exploding the interface: experiences of a CSCW network” *CHI '93: Proceedings of the CHI '93 Conference on Human Factors in Computing Systems* (1993) pp.255-262. DOI: 10.1145/169059.169205
9. Charles Peirce, ed., *On a New List of Categories* (University of North Carolina Press; 1991) pp. 23–33.
10. Matthew Arnold et al, “FactSheets: Increasing Trust in AI Services through Supplier’s Declarations of Conformity” *IBM Journal* 63 (2019) pp. 1-13.
11. Julia Angwin et al “Machine Bias There’s software used across the country to predict future criminals. And it’s biased against blacks” *ProPublica*, 23 May 2016, <www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing> (accessed 9 July 2020).
12. Joseph Lindley, Paul Coulton “AHRC Challenges of the Future: AI & Data,” *AHRC* (2020). DOI: 10.13140/RG.2.2.29569.48481
13. Lindley et al. [5]
14. Saleema Amershi et al, “Guidelines for Human-AI Interaction” *CHI '19: Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems* (2019) pp.1-13. DOI: 10.1145/3290605.3300233
15. Arnold et al. [11]

16. Emilee Rader et al, “Explanations as Mechanisms for Supporting Algorithmic Transparency”. *CHI '18: Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems* (2018) pp. 1-13. DOI: 10.1145/3173574.3173677
17. Franziska Pilling et al, “Design (Non) Fiction: Deconstructing/Reconstructing The Definitional Dualism of AI” *International Journal Film Media Arts* 6, No. 1 (2021) pp. 6–32.
18. Jennifer Ferreira et al, “A case for Iconic Icons” (2006). DOI: 10.1145/1151758.1151771
19. Barr P et al, “Icon R Icons: User interface icons, metaphor and metonymy,” *Victoria University of Wellington School of Mathematical and Computing Sciences* (2002).
20. Ibid.
21. Ferreira et al [18].
22. Ibid.
23. Erik Davis. *Techgnosis: Myth, Magic, Mysticism in the Age of Information* (New York: Tree Rivers Press 1998).
24. Beatrice Fazi. *Contingent Computation: Abstraction, Experience, and Indeterminacy in Computation Aesthetics* (London: Rowman & Littlefield 2018).
25. Sarayu Natarajan et al, “Just and Equitable Data Labelling towards a responsible AI supply chain,” *aapti institute* (2021).
26. Adam Rogers, “The Way the World Ends: Not with a Bang But a Paperclip,” *Wired*. 21 October 2017 <www.wired.com/story/the-way-the-world-ends-not-with-a-bang-but-a-paperclip/> (accessed 8 July 2021).