# Enterprise Liability for Human-AI-Decisions:

## A Multidisciplinary Approach for Identifying Principals' Duties of Care

## Ann-Kristin Mayrhofer\*

## **Contents**

I. Introduction
II. Fault-Based Liability and Principals' Duties of Care: The Notions of
"Possibility" and "Reasonableness" as Methodological Gateway to
Multidisciplinary Approaches
III. "Possible" Measures: Specific Risks of Human-AI-Decisions and Risk
Mitigation
A. Automation Complacency and Automation Bias: Over-Reliance on Al
Systems44
B. Algorithm Aversion: Under-Reliance on AI Systems
C. Conclusions: Which Measures are "Possible"?49
IV. "Reasonable" Measures: Limits of Risk Mitigation in Human-AI
Decisions
A. Principals' Duties of Care Under the Proposal for an AI Liability
Directive54
B. Summary and Perspectives
V. Bibliography58

\*Ann-Kristin Mayrhofer is Akademische Rätin a. Z. (Research Fellow) at the Chair of Civil Law, Civil Procedure, European Private Law and Procedure of Prof. Dr. Beate Gsell at Ludwig-Maximilians-Universität München.



#### I. Introduction

Artificial intelligence (AI) systems are already outperforming humans in a significant number of tasks, such as diagnosing diseases, classifying objects in images or predicting loan default risks,<sup>2</sup> and it is likely that this number will keep rising.<sup>3</sup> AI systems typically excel over humans when the task involves processing vast amounts of data and making general assumptions.<sup>4</sup> At the same time, AI systems are still limited in ways humans are not. They usually have difficulties when the task requires them to deviate from a general pattern and take into account the particularities of an individual case. The following example, concerning driving, illustrates the respective strengths and weaknesses of AI systems and humans: An AI driver confronted with

<sup>&</sup>lt;sup>5</sup> Cf. Kevin Bauer et al., 'Die KI braucht bei der Bankberatung immer noch menschliche Hilfe', available at https://www.boersen-zeitung.de/kapitalmarktforschung/in-derbankberatung-braucht-die-ki-menschliche-hilfe-90edbb42-86a4-11ed-a311-f90ecc32c8e4 accessed 25 April 2025); for a philosophical view on AI-decision-making, cf. Andreas Kaminski, 'Gründe geben. Maschinelles Lernen als Problem der Moralfähigkeit von Entscheidungen', in Klaus Wiegerling et al. (eds.), Datafizierung und Big Data (Wiesbaden: Springer VS, 2020) 151-74.



<sup>&</sup>lt;sup>1</sup> In this paper, the - controversial - term of "Artificial Intelligence" is understood in a broad sense: It is meant to include "simple" algorithms which are not based on machine learning (ML) techniques; cf. the definition of "AI system" in Art. 3(1) of the Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828 (AI Act): "a machine-based system that is designed to operate with varying levels of autonomy and that may exhibit adaptiveness after deployment, and that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments".

Cf. Yotam Liel and Lior Zalmanson, "Turning Off Your Better Judgment - Conformity to paper) Algorithmic Recommendations' (Working (2022),available https://www.researchgate.net/publication/366412145\_Turning\_Off\_Your\_Better\_Judgment\_-Conformity\_to\_Algorithmic\_Recommendations (last accessed 25 April 2025), pp. 5-6 (with further references; an abbreviated version of this working paper was published in Academy of Management Proceedings 2023(1); hereafter only the extended working paper will be referenced).

<sup>&</sup>lt;sup>3</sup> For an overview of use cases, cf. Deutsches Institut für Normung and Deutsche Kommission Elektrotechnik, Elektronik, Deutsche Normungsroadmap Künstliche Intelligenz, 2nd edn. (2022), available at https://www.din.de/de/forschung-und-innovation/themen/kuenstliche-intelligenz/fahrplanfestlegen (last accessed 25 April 2025); Katja Grace et al., 'Viewpoint: When Will AI Exceed Human Performance? Evidence from AI Experts' (2018) 62 Journal of Artificial Intelligence Research 729-54 conducted a survey among ML researchers in 2016 and found that, on average, they believed there was a 50 % chance AI would outperform humans in all tasks in 45 years' time.

<sup>&</sup>lt;sup>4</sup> Cf. Kathleen L. Mosier and Linda J. Skitka, 'Human Decision Makers and Automated Decision Aids: Made for Each Other?', in Raja Parasuraman and Mustapha Mouloua (eds.), Automation and human performance: Theory and application (Mahwah: Lawrence Erlbaum, 1996) 201-20, pp. 201 and 209.

the usual road signs may cause less accidents than a human driver. However, if a new sign appears, e.g., in a new country, which does not match the system's knowledge, the AI driver could cause an accident a human driver may have avoided by stopping and asking a pedestrian about the sign's meaning. Furthermore, small changes to traffic signs, that are invisible to the human eye and do not change the signs' meaning, could fool AI systems and lead to accidents humans would not have caused. Generally, there are some damage risks that are better avoided by AI systems and others that only humans could help to prevent.8 This suggests that "human-ML [machine learning] augmentation, where humans and technology work together to perform organisational tasks jointly" may be "the most promising path". In fact, it is expected that enterprises will increasingly integrate Human-AI-Decisions into their organisations. In principle, a Human-AI-Decision could replace any decision previously made by either a human or an AI system. It could, for example, concern medical diagnoses,<sup>11</sup> credit scoring<sup>12</sup> or the safety of products or services offered by the enterprise<sup>13</sup>. Combining humans and AI systems could enable the utilisation of both human and AI potential and therefore increase efficiency and safety. However, while many damage risks may be reduced by this combination, the Human-AI

<sup>&</sup>lt;sup>13</sup> For a concept of an autonomous warehouse, see e.g., Ahmet Börütecene and Jonas Löwgren, 'Designing Human-Automation Collaboration for Predictive Maintenance', in Companion Publication of the 2020 ACM Designing Interactive Systems Conference (New York: Association for Computing Machinery, 2020) 251-6.



<sup>&</sup>lt;sup>6</sup> Cf. the similar example by Erik J. Larson, The myth of artificial intelligence, p. 124.

<sup>&</sup>lt;sup>7</sup> Cf. Kevin Evkholt et al., 'Robust Physical-World Attacks on Deep Learning Visual Classification', in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (Salt Lake City: IEEE, 2018) 1625-34.

<sup>&</sup>lt;sup>8</sup> Cf. Philipp Hacker, 'Verhaltens- und Wissenszurechnung beim Einsatz von Künstlicher Intelligenz' (2018) 9(3) RW 243-88, p. 263 who distinguishes between three types of errors: errors which are only committed by AI systems, errors which are committed both by (reasonable) humans and AI systems and errors which are only committed by humans.

Cf. Mike H. M. Teodorescu et al., 'Failures of Fairness in Automation Require a Deeper Understanding of Human-ML Augmentation' (Minneapolis: University of Minnesota 2021) 45(3) MIS Quarterly 1483-500, p. 1484 on "achieving fairness".

<sup>&</sup>lt;sup>10</sup> Cf. Liel and Zalmanson, 'Turning Off Your Better Judgment', p. 2 (with further references).

<sup>&</sup>lt;sup>11</sup> Cf. Ekaterina Jussupow et al., 'Augmenting Medical Diagnosis Decisions? An Investigation into Physicians' Decision-Making Process with Artificial Intelligence' (2021) 32(3) Information System Research 713-35.

<sup>&</sup>lt;sup>12</sup> Cf. Rita Gsenger and Toma Strle, 'Trust, Automation Bias and Aversion: Algorithmic Decision-Making in the Context of Credit Scoring' (2021) 19(4) Interdisciplinary Description of Complex Systems 542-60.

cooperation may also give rise to new damage risks. <sup>14</sup> These risks trigger the question of liability. The most obvious issue seems to be the liability of the human who directly cooperates with the AI system. However, this paper addresses the liability of the enterprise that makes use of the Human-AI cooperation. More precisely, it is the liability of the *principal*, the corporation or the single entrepreneur that holds the enterprise, which is at stake. <sup>15</sup> Compared to the human agent's liability, enterprise liability may provide considerable advantages for victims: Usually, it is easier to identify the enterprise the Human-AI-Decision is integrated in than the individual human cooperating with the AI system. Furthermore, the principal is usually in a better financial situation.<sup>16</sup>

This paper focuses on the enterprise's non-contractual and fault-based liability.<sup>17</sup> Its aim is to use findings from other disciplines to identify some of the duties of care that principals must comply with. The paper will illustrate the need for a multidisciplinary approach in the context of technology particularly regarding Human-AI-Decisions<sup>18</sup> whose facets "are as complex as the environments in which they function". <sup>19</sup> Given the author's background, the analysis is based on a German perspective. However, it sems that the general ideas can also be applied to other legal systems. The paper first lays some principles of fault-based liability law and, at the same time, sets out the methodological framework of the applied multidisciplinary approach (II.). It will be shown that the notions of "possibility" and "reasonableness", commonly used to define duties of care, can serve as a methodological gateway to such approach.



<sup>&</sup>lt;sup>14</sup> Cf. the examples at Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids',

<sup>&</sup>lt;sup>15</sup> Non-human principals, namely corporations, generally act through human representatives (see, for example, § 31 of the German Civil Code (BGB)).

<sup>&</sup>lt;sup>16</sup> Cf. Helmut Koziol, 'Concluding Remarks', in Helmut Koziol (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015) 182-95, paras 5-27.

For a broader analysis of liability for autonomous systems (humans, animals, and AI systems), cf. Ann-Kristin Mayrhofer, Außervertragliche Haltung für fremde Autonomie (Tübingen: Mohr Siebeck, 2023).

<sup>&</sup>lt;sup>18</sup> Cf. Anna Beckers and Gunther Teubner, *Three Liability Regimes for Artificial Intelligence* (Oxford: Hart, 2022), p. 16. One may argue whether the following approach is "multidisciplinary" or already "interdisciplinary" (or "transdisciplinary"). The borders are fluid, and the terminology does not influence the following considerations; for definitions, see Benard C.K. Choi and Anita W.P. Pak, 'Multidisciplinarity, interdisciplinarity, and transdisciplinarity in health research, services, education and policy: 1. Definitions, objectives, and evidence of effectiveness' (2006) 29(6) Clinical and Investigative Medicine 351-64; Jacqueline Fawcett, 'Thoughts About Multidisciplinary, Interdisciplinary, and Transdisciplinary Research' (2013) 26(4) Nursing Science Quarterly 376–9; Eric Hilgendorf, 'Bedingungen gelingender Interdisziplinarität' (2010) 65(19) JZ 913-22, pp. 914-5.

<sup>&</sup>lt;sup>19</sup> Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', p. 202.

Second, the paper examines - as its main part - measures which have been identified by non-legal experts as "possible" mitigators of risks associated with Human-AI-Decisions (III.). Third, the paper explores - more briefly - how economic considerations may help to limit principals' duties to the "reasonable" measures (IV.).20 Finally, it looks at the European Commission's Proposal for an AI Liability Directive and its significance for Human-AI-Decisions and the principals' corresponding duties of care (A.). The AI Liability Directive was proposed in 2022 and abandonded in 2025.

#### II. Fault-Based Liability and Principals' Duties of Care: The Notions of "Possibility" and "Reasonableness" as Methodological Gateway to Multidisciplinary **Approaches**

Most legal systems recognise various forms of non-contractual liability.<sup>21</sup> This paper focuses on fault-based liability and more specifically on negligence liability. Faultbased liability "has been recognised for centuries as an element of liability" and, unlike e.g., strict liability, <sup>23</sup> applies in all kinds of situations. Comparative studies have shown that the term of "fault" could refer to a variety of different concepts. 24 Yet, there seems to be a certain agreement that negligence liability presupposes at least the violation of a *duty of care*. <sup>25</sup> Therefore, this paper seeks to identify some of the principals' duties of care in relation to Human-AI-Decisions integrated in their

<sup>&</sup>lt;sup>25</sup> Cf. Koziol, 'Comparative Conclusions', in Koziol (ed.), Basic Questions of Tort Law from a Comparative Perspective, paras 8/224 and 8/226. However, there is no complete consensus on this concept, cf. Gert Brüggemeier, Haftungsrecht: Struktur, Prinzipien, Schutzbereich (Berlin et al.: Springer, 2006), pp. 52-5 (with further references); De Bruyne et al., 'The European Commission's approach to extra-contractual liability and AI', p. 9 (with further references).



<sup>&</sup>lt;sup>20</sup> For an economic analysis of AI liability, see Gerhard Wagner, 'Roboter als Haftungssubjekte? Konturen eines Haftungsrechts für autonome Systeme', in Florian Faust and Hans-Bernd Schäfer (eds.), Zivilrechtliche und rechtsökonomische Probleme des Internet und der künstlichen Intelligenz (Tübingen: Mohr Siebeck, 2019) 1-39.

For an overview of tort law in different jurisdictions, cf. Helmut Koziol (ed.), Basic Questions of Tort Law from a Germanic Perspective (Vienna: Sramek, 2012); Helmut Koziol (ed.), Basic Questions of Tort Law from a Comparative Perspective (Vienna: Sramek, 2015) and Helmut Koziol (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015).

<sup>&</sup>lt;sup>22</sup> Helmut Koziol, 'Comparative Conclusions' in Helmut Koziol (ed.), *Basic Questions of Tort Law* from a Comparative Perspective (Vienna: Sramek, 2015) 685-838, para 8/218.

With strict liability "European legal systems show much more diversity than in other areas of tort law", ibid., para 8/35.

<sup>&</sup>lt;sup>24</sup> Cf. ibid., paras 8/219-224; Jan De Bruyne et al., 'The European Commission's approach to extracontractual liability and AI - An evaluation of the AI liability directive and the revised product liability directive' (2023) 51 Computer Law & Security Review 105894, pp. 8-9.

organisations. The practical significance of principals' duties of care for the enterprise's liability varies between jurisdictions: <sup>26</sup> Under German law, for example, principals are only liable for damages caused by their - human or artificial - assistants if the principals acted faultily themselves, e.g., by violating their supervision duties.<sup>27</sup> Other jurisdictions, for example the French system, contain stricter rules.<sup>28</sup> In these jurisdictions, the question of whether a duty of care has been violated is less pressing. Liability for Human-AI-Decisions under such specific and, to a significant extent, system dependent rules will not be discussed in this paper either.<sup>29</sup>

Principals who integrate humans into their organisations must carefully select, instruct, and supervise each person.<sup>30</sup> In principle, the same applies to technical devices, including AI systems. However, duties of care are not limited to the direct relationship between principals and - human or artificial - assistants. Rather, principals also need to coordinate the cooperation between the assistants. 32 Regarding German law, the German Federal Court of Justice (Bundesgerichtshof, BGH) holds that a person does not need to avoid every risk but must (only) do what is "possible" and "reasonable". 33 As the idea that "no one is obliged beyond what they are able to do" (ultra posse nemo obligator) and the standard of the "reasonable person" are

<sup>&</sup>lt;sup>33</sup> BGH, 6 February 2007, VI ZR 274/05 (2007) NIW 1683-5, para 14; cf., also, Gerhard Wagner, '§ 823 BGB', in Jürgen Säcker et al. (eds.), Münchener Kommentar zum Bürgerlichen Gesetzbuch, 9th edn., 13 vols. (Munich: C.H. Beck, 2024), vol. VII, para. 528 (with further references).



<sup>&</sup>lt;sup>26</sup> For a comparative overview cf. Koziol, 'Comparative Conclusions', in Koziol (ed.), *Basic Questions* of Tort Law from a Comparative Perspective, paras 8/250-5 and the contributions in Helmut Koziol et al., 'Liability for Agents and Agents' Liability', in Helmut Koziol (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015) 182-95, pp. 147-95.

<sup>&</sup>lt;sup>27</sup> Cf. §§ 831, 823 BGB. However, § 831 BGB contains a presumption of the violation of a duty of

<sup>&</sup>lt;sup>28</sup> Cf. Art. 1242 of the French Civil Code.

Therefore, the question of whether specific principal-agent-liability applies to AI systems will also remain open. At least, in more restrictive systems, such as the German system, an analogy would not lead to principals' strict or vicarious liability, for the discussion, cf. Hacker, 'Verhaltens- und Wissenszurechnung', pp. 265-7; Herbert Zech, 'Entscheidungen digitaler autonomer Systeme: Empfehlen sich Regelungen zu Verantwortung und Haftung?", in Ständige Deputation des Deutschen Juristentages (ed.), Verhandlungen des 73. Deutschen Juristentages (Munich: C.H. Beck, 2020) A 1-112, pp. A 76-81 and A 95-96; Beckers and Teubner, Three Liability Regimes for Artificial Intelligence, pp. 46-87 and passim.

<sup>&</sup>lt;sup>30</sup> Cf. § 831(1) sentence 2 BGB.

<sup>&</sup>lt;sup>31</sup> Cf. Susanne Horner and Markus Kaulartz, 'Haftung 4.0' (2016) 32(1) CR 7-14, p. 8.

Mayrhofer, Außervertragliche Haftung für fremde Autonomie, p. 111 and 343 with further references regarding so-called "organisational duties" ("Organisationspflichten").

largely recognised in most national laws, 34 it seems that the duties of care are defined similarly in other jurisdictions.

In simple cases the required safety standard may be obvious. For instance, it is certainly "possible" and "reasonable" to slow down before a sharp turn to avoid accidents. Difficulties arise in more complex cases where assessing the "possibility" and "reasonableness" of a measure requires a lot of non-legal knowledge. This is where multidisciplinary approaches can be useful:<sup>35</sup> Such approaches can help to concretise abstract concepts, <sup>36</sup> such as "possibility" and "reasonableness" <sup>37</sup>. These open-formulated requirements can serve as a methodological "gateway" to multidisciplinary approaches.<sup>38</sup> Cases of damages that involve Human-AI-Decisions tend to be complex: Humans and AI systems are complex entities, and the complexity increases when they work together.<sup>39</sup> Determining whether a particular measure provides a safety benefit and whether that benefit outweighs the effort that a principal must take to adopt the measure, is very difficult.<sup>40</sup>

Written standards, such as product safety legislation and technical standards, can provide some relief. However, they usually do not cover every situation exhaustively,

<sup>&</sup>lt;sup>40</sup> Cf. in the context of product liability, Mayrhofer, *Außervertragliche Haftung für fremde Autonomie*, p. 269.



<sup>&</sup>lt;sup>34</sup> Both the idea that "no one is obliged beyond what they are able to do" (*ultra posse nemo obligator*) and the standard of the "reasonable person" seem to be largely recognised in most national laws, cf. on the latter Art. 4:102 (1) Principles of European Tort Law (PETL), available at http://www.egtl.org/PETLEnglish.html (last accessed 25 April 2025); Koziol, 'Comparative Conclusions', in Koziol (ed.), Basic Questions of Tort Law from a Comparative Perspective, paras 8/229-36.

<sup>35</sup> Cf. Franz Hofmann, 'Disziplinarität, Interdisziplinarität und Interdisziplinarität am Beispiel der Grundsätze "mittelbarer Verantwortlichkeit" (2018) 73(15-16) JZ 746-54, pp. 749-50 and 753-4.

<sup>&</sup>lt;sup>36</sup> Cf. Hilgendorf, 'Bedingungen gelingender Interdisziplinarität', p. 920.

<sup>&</sup>lt;sup>37</sup> Cf. Hofmann, 'Disziplinarität, Intradisziplinarität und Interdisziplinarität', pp. 750 and 753-754 who uses interdisciplinary approaches to identify duties of care in the form of so called "Verkehrspflichten".

<sup>38</sup> Cf. ibid., p. 750 ("Einfallstor"); Rolf Stürner, 'Die Zivilrechtswissenschaft und ihre Methodik – zu rechtsanwendungsbezogen und zu wenig grundlagenorientiert?' (2014) 214(1-2) AcP 7-54, pp. 31-32 ("Einbruchstellen").

<sup>&</sup>lt;sup>39</sup> Cf. Expert Group on Liability and New Technologies, Liability for Artificial Intelligence and Other Technologies' (2019),https://www.europarl.europa.eu/meetdocs/2014 2019/plmrep/COMMITTEES/JURI/DV/2020/01-09/AI-report EN.pdf (last accessed 25 April 2025), pp. 32-3.

especially when it comes to rapidly evolving technologies such as AI. 41 The AI Act, 42 for example, will certainly play an important role in the assessment of fault-based liability. It focuses mainly on "high-risk" AI systems and "general-purpose AI models". 43 Yet, many of its requirements leave room for interpretation, 44 such as the need for high-risk AI systems to be "effectively overseen by natural persons". <sup>45</sup> Thus, the written standards will most likely need to be supplemented by courts which will need to equip themselves with the non-legal knowledge necessary to adequately assess the specific risks of the Human-AI-Decision. In theory, the need for multidisciplinary approaches seems to be recognised by the courts. The BGH, for example, refers to the "state of scientific and technical knowledge" when determining whether a safety measure was "possible". 46 According to the BGH, enterprises are able to avoid a risk "if, according to the assured expert knowledge of the relevant specialist groups, solutions are available that can be used in practice". Furthermore, when determining, what is "reasonable", the BGH considers "all the circumstances of the individual case", in particular "the magnitude of the risk", but also "the economic effects of the safeguarding measure". 48 This wording suggests that the safety standards are indeed set using knowledge from other disciplines, including technical, sociological and economic aspects. In practice, however, it can be difficult for judges to make use of non-legal knowledge.49 Therefore, the paper will now show how studies on Human-AI-Decisions can effectively be made fruitful for the legal analysis of the principals' duties of care. It will present examples of non-legal studies and draw some conclusions on the legal assessment of "possibility" and "reasonableness". While the first step - presentation of the non-legal knowledge - is rather descriptive,



<sup>&</sup>lt;sup>41</sup> Cf. in the context of product liability, ibid., pp. 275-85.

<sup>42</sup> Cf. note 1 above.

<sup>&</sup>lt;sup>43</sup> Principals that integrate Human-AI-Decisions into their organisations will frequently qualify as "deployer" under the AI Act (Art, 3(4)). However, if they develop the AI system or general-purpose AI model themselves or have the AI system developed, they can also be classified as "providers" (Art. 3(3)).

<sup>44</sup> Cf. David Bomhard and Marieke Merkle, 'Europäische KI-Verordnung - Der aktuelle Kommissionsentwurf und praktische Auswirkungen' (2021) 1(6) RDi 276-83, p. 283.

<sup>&</sup>lt;sup>45</sup> Cf. Art. 14 AI Act.

<sup>&</sup>lt;sup>46</sup> BGH, 16 September 2009, VI ZR 107/08, BGHZ 181, 253-68, para 16 (juris). The decisions mainly concern manufacturers' liability. However, there is no reason for limiting this approach to such enterprises so that it may also apply to e.g., service providers.

<sup>47</sup> Ibid.

<sup>48</sup> Ibid., para 18.

<sup>&</sup>lt;sup>49</sup> Cf. Hilgendorf, 'Bedingungen gelingender Interdisziplinarität', p. 917.

the second step - legal conclusions - aims to provide some new guidelines for the assessment of duties of care.

# III. "Possible" Measures: Specific Risks of Human-AI-Decisions and Risk Mitigation

Human-AI cooperation can take various forms. Three categories are frequently employed: First, AI systems can *provide information* to the human. Second, they can make recommendations. Third, AI systems can take the decision on behalf of the human. This paper focuses on AI systems that generate an output that the human can accept or reject before it causes damage. The human is not completely replaced by the AI system but stays "in the loop". Mostly, these systems give recommendations (second category). However, the output could also consist of information the human can disregard (first category) or a decision the human can override (third category). 52

The specific risks of Human-AI-Decisions explored here are risks associated with precisely this option to accept or reject an AI output. The aim of Human-AI cooperation - utilising both AI and human strengths - is only achieved if humans correctly "exercise their own judgment [...] to minimise risks of bad or biased decisions". 53 Otherwise, correct AI output may be rejected or wrong AI output may be accepted, and damages may occur, also to third parties. Particularly, just as when working with other humans, humans need to have the right level of trust and mistrust

<sup>&</sup>lt;sup>53</sup> Cf. Liel and Zalmanson, 'Turning Off Your Better Judgment', p. 2 (with further references).



<sup>&</sup>lt;sup>50</sup> Cf. the definition of an "AI system" in Art. 3(1) AI Act: "generate outputs such as predictions, content, recommendations, or decisions that can influence physical or virtual environments"; Datenethikkommission der Bundesregierung, Gutachten (2019). https://www.bmi.bund.de/SharedDocs/downloads/DE/publikationen/themen/itdigitalpolitik/gutachten-datenethikkommission.html (last accessed 25 April 2025), pp. 161-2: "algorithm-based", "algorithm-driven" and "algorithm-determined and therefore completely automated" decisions.

Cf. High-Level Expert Group on Artificial Intelligence, Ethics Guidelines for Trustworthy Al' (2019), available at https://digital-strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai (last accessed 25 April 2025), p. 16, where a distinction is made between three approaches of human oversight: "human-in-the-loop (HITL), human-on-the-loop (HOTL), or human-in-command (HIC)" and where HITL is defined as "the capability for human intervention in every decision cycle of the system".

<sup>&</sup>lt;sup>52</sup> Cf. the "Levels of Automation of Decision and Action Selection" in Matthew Ball and Vic Callaghan, 'Explorations of Autonomy', in 2012 Proceedings of the 8th International Conference on Intelligent Environments (Guantajo: IEE Xplore, 2012) 114-21.

towards their artificial "colleague". However, three major phenomena appear to hinder humans from adopting the appropriate attitude towards AI systems, which can be categorized into two groups:55 (A.) Automation Complacency and Automation Bias on the one hand and (B.) Algorithmic Aversion on the other. Non-legal experts have conducted research on these problems and have suggested mitigating factors that can help to determine the "state of scientific and technical knowledge" regarding Human-AI-Decisions and to identify "possible" risk prevention measures.

#### A. Automation Complacency and Automation Bias: Over-Reliance on AI Systems

In the context of this paper, Automation Complacency and Automation Bias can be addressed together. They "represent different manifestations of overlapping automation-induced phenomenon" <sup>56</sup> as they both lead to "inappropriate overreliance on automation". 57 While Automation Complacency is generally linked to insufficient monitoring of a system's output,<sup>58</sup> Automation Bias describes humans' tendency to accept the output "as a heuristic replacement for vigilant information seeking and processing". 59 Research on these phenomena has started in the aviation sector: Human pilots frequently did not monitor the autopilots properly or accepted suboptimal flight plans. 60 Subsequently, many studies have been conducted in different contexts to find out more about causes and possible mitigators of human over-

<sup>&</sup>lt;sup>60</sup> Cf. Gsenger and Strle, 'Trust, Automation Bias and Aversion', p. 546; Parasuraman and Manzey, 'Complacency and Bias in Human Use of Automation', pp. 381-2; Mary L. Cummings, 'Automation Bias in Intelligent Time Critical Decision Support Systems' (2004) AIAA 1st Intelligent Systems Technical Conference, 20-22 September 2004, Chicago, IL, American Institute of Aeronautics and Astronautics 1-6 (with reference to further studies).



<sup>&</sup>lt;sup>54</sup> Cf. Gsenger and Strle, 'Trust, Automation Bias and Aversion', pp. 545-6.

<sup>&</sup>lt;sup>55</sup> There seem to be other phenomena, e.g., "selective adherence", cf. Saar Alon-Barkat and Madalina Busuioc, 'Human-AI Interactions in Public Sector Decision-Making: "Automation Bias" and "Selective Adherence" to Algorithmic Advice' (2023) 33(1) Journal of Public Administration Research and Theory 153-69, p. 154.

<sup>&</sup>lt;sup>56</sup> Raja Parasuraman and Dietrich H. Manzey, 'Complacency and Bias in Human Use of Automation: An Attentional Integration' (2010) 53(3) Human Factors 381-410, p. 405.

<sup>&</sup>lt;sup>57</sup> Ibid., p. 398.

<sup>&</sup>lt;sup>58</sup> Ibid., p. 382; the authors point out that "there is no consensus on the definition of complacency". Also, cf. Kate Goddard et al., 'Automation bias: a systematic review of frequency, effect mediators, and mitigators' (2012) 19(1) Journal of the American Medical Informatics Association 121-7, pp. 121-

<sup>&</sup>lt;sup>59</sup> Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', p. 205; cf. Goddard et al., 'Automation bias', p. 121.

reliance on algorithms. <sup>61</sup> Not all of them can be examined here. Rather, this paper focuses on two examples dealing with two distinct kinds of decisions.

#### 1. "Difficult" Decisions: Medical Diagnostic

Jussupow et al. conducted an experiment with physicians who had to make medical diagnoses based on radiological data and advice provided by an AI system. 62 In the control group - physicians without AI advice - the accuracy rate, manifesting the decision-making results, was about 77 %. The authors found that the accuracy rate was significantly lower among the physicians who received incorrect AI advice (about 55 % accuracy rate). To investigate the decision-making *process*, they also analysed think-aloud protocols, questionnaires and interviews. <sup>64</sup> According to *Jussupow* et al., three "decision pathways" could lead to erroneous conformity to AI advice: First, physicians tend to rely on AI advice for their own initial assessment. Second, when AI advice disconfirms their initial position, novice physicians are likely to doubt their own capabilities. They lean heavily on their *belief* in the *system* without further validation of the specific assessments ("belief conflict" vs. "validation conflict"). Experienced physicians more often ignore the AI advice or engage in such an evaluation of the advice. However, in case of such validation, physicians frequently do not reconsider both their own and the AI assessment but only collect data to reject or confirm *one* of the positions. This leads to the third pathway: If physicians opt for only checking the AI assessment, they often fail to find or even search for data that contradicts the AI advice. 65

Regarding mitigators against erroneous conformity to AI advice, the authors suggest that the system *design* could have an influence. Physicians seem to ignore wrong advice more often when it is presented after they have already made an initial assessment. 66 This idea is in line with findings of behavioural economists who suggest improving the decisions of (human) groups by having participants form their opinions



<sup>&</sup>lt;sup>61</sup> Cf. III.A.1. and III.A.2. below, and the examples provided by Liel and Zalmanson, 'Turning Off Your Better Judgment', pp. 3-8.

<sup>&</sup>lt;sup>62</sup> Jussupow et al., 'Augmenting Medical Diagnosis Decisions?'.

<sup>&</sup>lt;sup>63</sup> Ibid., pp. 721-2. In contrast, the accuracy rate was only marginally higher than the accuracy rate of the control group among those participants who received correct advice (about 90 % accuracy rate).

<sup>&</sup>lt;sup>64</sup> Ibid., p. 718.

<sup>65</sup> Ibid., pp. 729-30 and passim; cf. already Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', p. 202-205 where humans' tendency to interpret other information as being consistent with the automated decision aid is highlighted.

<sup>&</sup>lt;sup>66</sup> Jussupow et al., 'Augmenting Medical Diagnosis Decisions?', p. 730.

before the discussion. <sup>67</sup> Jussupow et al. also highlight the importance of training: Physicians need to learn how to validate the AI system's and their own position, so that first, they do not get stuck in the "belief conflict" and second, they are able to solve the "validation conflict" by reconsidering both their own and the AI's assessment ("system monitoring" and "self-monitoring").69

## 2. "Easy" Decisions: Simple Image Classification

Liel and Zalmanson cite Jussupow et al. as an example of a context "in which the purpose of the algorithm was to assist in complex calculations or judgment under uncertainty". According to *Liel* and *Zalmanson*, in such situations, erroneous conformity to AI advice could be "justifiable". Their own experiment, by contrast, dealt with "objective and straightforward tasks, with very little uncertainty".72 Participants were shown a set of three items alongside a fourth one and were asked to identify the item in the set which was equal in length to the fourth item. The participants were split into three groups: A control group which received no advice, a second group which received wrong advice framed as being generated by an AI, and a third group which received wrong advice framed as being generated by humans. Members of the control group were accurate in around 98 % of the tasks which demonstrates the tasks' simplicity. Members of the second group followed the erroneous AI advice in around 27 % of the tasks, whereas members of the third group followed the erroneous human advice in about 19 % of the tasks.<sup>74</sup>

According to *Liel* and *Zalmanson*, these results show that humans do not only overrely on AI advice in cases of uncertainty. Rather, they also exhibit an "'irrational'



<sup>67</sup> Cf. Daniel Kahneman et al., 'NOISE' (2016) 94(10) Harvard Business Review 38-46, p. 46.

<sup>&</sup>lt;sup>68</sup> Cf. also Andreas Fügener et al., 'Cognitive Challenges in Human-Artificial Intelligence Collaboration: Investigating the Path Toward Productive Delegation' (2022) 33(2) Information Systems Research 678-96: The authors conducted experiments where humans did not perform well in delegating tasks to AI systems. They explain this result mainly by humans' lack of knowledge about their own capacities ("metaknowledge").

<sup>&</sup>lt;sup>69</sup> Jussupow et al., 'Augmenting Medical Diagnosis Decisions?', pp. 731-2.

<sup>&</sup>lt;sup>70</sup> Liel and Zalmanson, 'Turning Off Your Better Judgment', p. 7.

<sup>&</sup>lt;sup>71</sup> Ibid., p. 8.

<sup>&</sup>lt;sup>72</sup> Ibid., p. 8.

<sup>&</sup>lt;sup>73</sup> Ibid., pp. 15-8.

<sup>&</sup>lt;sup>74</sup> Ibid., pp. 18-20.

tendency to adopt advice that contradicts their better judgment". The comparison with human advice indicates that at least in certain situations "algorithmic conformity" can even exceed "social conformity". <sup>76</sup> One could argue that, in real life, such simple tasks are unlikely to be carried out by Human-AI cooperation. However, it seems that the findings could be transferred to AI advice concerning more difficult tasks, e.g., driving. AI advice could also be obviously wrong, e.g., if the AI driver clearly misinterprets a new sign. Adherence to such AI advice can very well have serious real-life impact. This hypothesis is supported by anecdotes about accidents involving clearly erroneous navigation system.<sup>78</sup>

Liel and Zalmanson also investigated factors that could mitigate conformity. They conducted more experiments and found that the tendency to conform to incorrect AI advice was significantly lower when a second AI system provided *correct* advice.<sup>79</sup> In addition, participants seemed to perform better when they were told that their task was of "high significance", e.g., that it was to improve the safety of autonomous vehicles. 80 Consequently, they suggest including two or more AI systems in a Human-AI-Decision and encouraging users to perceive their tasks as important.<sup>81</sup>

### B. Algorithm Aversion: Under-Reliance on AI Systems

The aforementioned studies highlight that AI systems may give wrong advice in some cases. However, as AI systems frequently outperform humans, in other cases, humans should better trust the output of AI systems. In these situations, damages occur if humans do not use AI advice but rely on their own initial judgment or the



<sup>&</sup>lt;sup>75</sup> Ibid., p. 33.

The authors explicitly refer to experiments on social conformity. Their experimental design is based on studies conducted by Solomon E. Asch in the 1950s, cf. ibid., pp. 16-8.

<sup>&</sup>lt;sup>77</sup> Ibid., p. 8.

<sup>&</sup>lt;sup>78</sup> Cf. e.g., Greg Milner, 'Death by GPS: Are Satnavs changing our brains?', *The Guardian*, available at https://www.theguardian.com/technology/2016/jun/25/gps-horror-stories-driving-satnay-greg-milner (last accessed 25 April 2025). Reference to such "anecdotal evidence" of automation bias is also made by Alon-Barkat and Busuioc, 'Human-AI Interactions in Public Sector Decision-Making', p. 155.

<sup>&</sup>lt;sup>79</sup> Liel and Zalmanson, 'Turning Off Your Better Judgment', pp. 20-4.

One group was told the work served to improve the safety in autonomous vehicles, one group was told that it served to improve performance in smart warehouses, and one group did not receive any context, cf. ibid., pp. 24-7.

<sup>&</sup>lt;sup>81</sup> Ibid., pp. 34-5.

advice of other humans. This phenomenon is called Algorithm Aversion. 82 To explore the risks associated with this third problem and possible mitigators, again, two studies will be presented:

The first study was conducted by *Dietvorst* et al. Participants were provided with data about MBA students or U.S. states. They were then asked to predict the students' performance in the MBA program or the rank of U.S. states in terms of the number of airline passengers departing from that state.83 Participants had to choose between two advisors: a human or a statistical model. Prior to their predictions, participants either saw the AI system's performance on this task, the human's performance, both or neither.<sup>84</sup> The results suggest that Algorithm Aversion is particularly prevalent when humans have seen AI systems performing and making mistakes: Even after seeing the AI system outperforming the human, participants preferred human advisors.85 Therefore, it seems that "resistance at least partially arises from greater intolerance for error from algorithms than from humans". 86 At the same time, this finding suggests that Algorithm Aversion may be prevented by hiding the AI system's past errors from the human or by hiding the AI system itself.87

The second study, conducted by Yeomans et al., dealt with selecting jokes. 88 The authors found that AI systems generally outperformed humans in this task. They were usually able to select jokes which participants found funnier. 89 Nonetheless, participants seemed to prefer receiving recommendations from humans: They rated the advisor better when they believed it was human. 90 Yeomans et al. also found that advice that was framed as coming from a machine was perceived as less scrutable than advice framed as coming from a human.91 To the authors, this indicates the



<sup>82</sup> Cf. Berkeley J. Dietvorst et al., 'Algorithm aversion: People erroneously avoid algorithms after seeing them err' (2015) 144(1) Journal of Experimental Psychology: General 114-26, p. 114.

<sup>83</sup> Ibid., pp. 115 and 118.

<sup>&</sup>lt;sup>84</sup> Ibid., p. 115.

<sup>&</sup>lt;sup>85</sup> Ibid., p. 119.

<sup>86</sup> Ibid., p. 124.

<sup>&</sup>lt;sup>87</sup> Cf. Ibid., p. 124.

<sup>&</sup>lt;sup>88</sup> Michael Yeomans et al., 'Making sense of recommendations' (2019) 32(4) Journal of Behavioral Decision Making 403-14.

<sup>&</sup>lt;sup>89</sup> Ibid., pp. 404-8.

<sup>&</sup>lt;sup>90</sup> Ibid., pp. 409-10.

<sup>&</sup>lt;sup>91</sup> Ibid., pp. 410-1.

"subjective understanding" of the recommenders influence the assessment. 92 This idea was supported in another experiment, where one group received only a sparse explanation of the system's decision-making process whereas the other one got details about the functioning.<sup>93</sup> It suggests that "algorithmic sensemaking" could mitigate Algorithm Aversion.<sup>94</sup> Such sensemaking could, for example, be achieved by improving the AI system's explainability.95 Yeomans et al. themselves propose that more experience with the AI system could increase humans' subjective understanding. However, more experience with the system may be accompanied by seeing the system making more mistakes which, according to the aforementioned studies of *Dietvorst* et al., could foster Algorithm Aversion. According to *Yeomans* et al, another way of increasing subjective understanding could be creating AI systems with a more human-like design. 97 Yet, it seems that strong human resemblance could equally have negative effects on humans' attitude toward the AI system. <sup>98</sup> Besides, such a design could make humans belief that the AI system takes its decisions the same way a human does. This is usually wrong and could equally lead to a false reevaluation of the advice. 99

#### C. Conclusions: Which Measures are "Possible"?

The research just presented might assist in identifying "possible" measures to avoid the risks of Human-AI-Decisions. Naturally, the extent to which Automation Complacency, Automation Bias and Algorithm Aversion influence Human-AI-Decision-making depends on the specific task. 100 For example, a study by *Lee* which dealt with managerial decisions suggests that when a task requires mostly



<sup>&</sup>lt;sup>92</sup> Ibid., pp. 410-1.

<sup>&</sup>lt;sup>93</sup> Ibid., p. 411.

<sup>&</sup>lt;sup>94</sup> Ibid., p. 412.

<sup>95</sup> Cf. on transparency measures, Fraunhofer IAIS, Leitfaden zur Gestaltung vertrauenswürdiger Intelligenz (2021).https://www.iais.fraunhofer.de/de/publikationen/studien/2021/ki-pruefkatalog.html (last accessed 25 April 2025), pp. 63-85.

<sup>&</sup>lt;sup>96</sup> Michael Yeomans et al., 'Making sense of recommendations', p. 412.

<sup>&</sup>lt;sup>97</sup> Ibid., p. 412.

 $<sup>^{98}</sup>$  Cf. the idea of an "uncanny valley" developed by the Japanese professor *Mashiro Mori* in the 1970s, Mashiro Mori, 'The Uncanny Valley' (2012) 19(2) IEEE Robotics & Automation Magazine 98-100 (translated by Karl F. MacDorman and Norri Kageki).

<sup>&</sup>lt;sup>99</sup> Cf. already Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', pp. 210.

<sup>&</sup>lt;sup>100</sup> Cf. Gsenger and Strle, 'Trust, Automation Bias and Aversion', p. 548.

"mechanical" skills, e.g., work scheduling, AI and human advice appear equally trustworthy, whereas in more "human" tasks, e.g., work evaluation, AI advice is perceived less trustworthy. 101 Alon-Barkat and Busuioc conducted a study on school board decisions concerning the employment of teachers and found no evidence for Automation Bias. 102 According to the authors one reason might be a "relative skepticism about the performative capacity of AI algorithms" which distinguished this task from "areas well-accustomed to such devices (aviation, healthcare), characterised by routine use of reliable automation, resulting in high levels of trust in their performance". 103 Liel and Zalmanson point out that the situation they experimented with - image classification - "may be characterised by particularly high levels of trust in technology and algorithms", which could certainly encourage conformity. However, despite these differences and difficulties, the studies and their interpretation by their authors allow drawing some general conclusions regarding "possible" risk prevention measures.

Some of these measures concern the design of the Human-AI cooperation: 105 Principals can ensure humans make an initial assessment before receiving the AI advice, e.g., by making them write down their own opinion first (cf. III.A.1.). Besides, they can provide human agents with advice from multiple AI systems (cf. III.A.2.). Additionally, training could serve as a mitigator: Humans can be taught how both AI systems and humans make their decisions so that in case of a disagreement they are able to re-evaluate both the AI's advice and their own initial assessment (cf. III.A.1.). Furthermore, regarding specific task, *information* plays a significant role: Principals can inform the human agents about the general performance of specific AI systems and, if a comparison is possible, about the performance of competing humans. They can clarify to the human the *importance* of the task (cf. III.A.2.) and the organisation's expectations. 106 Besides, principals can ensure explainability of the AI system's decision-making process (cf. III.B.). However, this last measure presupposes

<sup>&</sup>lt;sup>106</sup> Ibid., p. 205 where the "extent to which organisations encourage the use of automated systems" is listed as one determining factor.



Min Kyung Lee, 'Understanding perception of algorithmic decisions: Fairness, trust, and emotion in response to algorithmic management' (2018) 5(1) Big Data & Society 1-16.

Alon-Barkat and Busuioc, 'Human-AI Interactions in Public Sector Decision-Making', pp. 157-64. <sup>103</sup> Ibid., p. 165.

Liel and Zalmanson, 'Turning Off Your Better Judgment', p. 33.

<sup>&</sup>lt;sup>105</sup> Cf. already Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', pp. 203-

first, that the AI system is not a complete "black box" 107 and second, that the explanation can be done understandably not only for experts, e.g., the AI's developers, but also for human agents who may not have specific technical knowledge 108. Effective information and explainability may require some additional training of the human agent.

When analysing these non-legal studies, however, it is essential to keep in mind the legal context in which the analysis is conducted. For instance, the findings of the studies may suggest that, in certain scenarios, erroneous Human-AI-Decisions could potentially be avoided by deliberately misleading humans: To mitigate Algorithm Aversion, e.g., in the context of "human" decisions, it may appear helpful to tell humans that the AI advice stems from another human. If the system is a "black box", humans may get fake explanations, which could foster their trust. Furthermore, to change humans' attitude towards an AI system in a more subtle way, previous mistakes could be concealed from them, or the AI system could be set up in a way that makes it appear more "human" (cf. III.B.). However, principals are generally not required or even allowed to take such measures: The "possibility" of a measure also needs to be assessed in light of the legal framework. For instance, to the extent that the algorithmic nature of advice must be made transparent, 109 hiding the source of advice does not constitute a *legally* possible risk mitigator.

#### IV. "Reasonable" Measures: Limits of Risk Mitigation in Human-AI-Decisions

Principals are not obliged to take all "possible" measures. Their duties of care only require them to do what is "reasonable". As mentioned above, to set this standard, the BGH, for example, considers *inter alia* "economic effects". <sup>110</sup> The precise role of economic aspects in liability law is a matter of debate beyond the scope of this



Cf. European Commission, 'White Paper on Artificial Intelligence - A European approach to excellence and trust', 19 February 2020, COM(2020) 65 final, p. 12, where "opacity ('black boxeffect')" is cited as one of the "specific characteristics of many AI technologies".

<sup>&</sup>lt;sup>108</sup> Cf. Mosier and Skitka, 'Human Decision Makers and Automated Decision Aids', p. 208; cf., also, Fraunhofer IAIS, Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher Intelligenz, p. 67, where the necessity to distinguish between different groups of users is highlighted.

Cf. e.g., Art. 13(2)(f) of the Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation); Art. 50 AI Act.

<sup>110</sup> Cf. notes 46 and 48 above.

paper. 111 All in all, it seems convincing to take them into account and even use them as a starting point, while bearing in mind they may need to be complemented, e.g., by social aspects.<sup>112</sup>

From an economic point of view, liability law should generally "reduce the sum of the costs of accidents and the cost of avoiding accidents". <sup>113</sup> Therefore, in principle, a safety measure should be taken if its costs are lower than the potential damage costs it prevents. 114 According to the so-called "Learned Hand formula", proposed in the U.S. by Judge Learned Hand, a person who does not take a risk prevention measure, acts negligently if " $B < (P \times L)$ ", where "B" refers to the burden of taking the measure, "P" to the probability of damage if the measure is not taken and "L" to the magnitude of such damage. 115 The usefulness of this formula is controversial. 116 However, it appears that courts, while refraining from explicitly citing "Learned Hand", at least tend to base their decisions implicitly on such cost-utility tests. 117

Regarding Human-AI-Decisions, economic considerations allow to draw at least the following conclusions which could provide additional guidelines to define the

<sup>&</sup>lt;sup>117</sup> Cf. the analysis in Schäfer and Ott, Lehrbuch der ökonomischen Analyse des Zivilrechts, 204-6 and Gerhard Wagner, '§ 823 BGB', para 531. The formula may be supplemented by other (economic) considerations, especially when the victim can equally take measures to avoid accidents, cf. Schäfer and Ott, Lehrbuch der ökonomischen Analyse des Zivilrechts, pp. 279-84 with further references, where the formula is complemented by the idea of the "cheapest cost avoider".



<sup>&</sup>lt;sup>111</sup> For an overview on the discussion on economic analysis of tort law in Germany, cf. Jochen Taupitz, 'Ökonomische Analyse und Haftungsrecht - Eine Zwischenbilanz' (1996) 196(1/2) AcP 114-67; Gerhard Wagner, 'Vor § 823 BGB', in Jürgen Säcker et al. (eds.), Münchener Kommentar zum Bürgerlichen Gesetzbuch, 9th edn., 13 vols. (Munich: C.H. Beck, 2024), vol. VII, paras 66-76; Hans-Bernd Schäfer and Claus Ott, Lehrbuch der ökonomischen Analyse des Zivilrechts, 6th edn. (Berlin et al.: Springer Gabler, 2020), pp. 165-71.

 $<sup>^{112}</sup>$  The need to complement the economic analysis is highlighted by Beckers and Teubner,  $\mathit{Three}$ Liability Regimes for Artificial Intelligence, pp. 16-7; Thomas M. J. Möllers, Juristische Methodenlehre, 5th edn. (Munich: C.H. Beck, 2023), p. 215-6; cf. equally the idea of a "socioeconomic analysis of law" mentioned by Josef Essers and Eike Schmidt, Schuldrecht, vol. I/1, 8th edn. (Heidelberg: C.F. Müller, 1995), p. 39; Taupitz, 'Ökonomische Analyse und Haftungsrecht', p. 126; cf., also, Martin Sommer, Haftung für autonome Systeme (Baden-Baden: Nomos, 2020), p. 270-2 (and passim) who convincingly suggests a "normativised" "risk-utility test".

Gudio Calabresi, *The Costs of Accidents* (Cumberland: Yale University Press, 1970), p. 26; cf. Schäfer and Ott, Lehrbuch der ökonomischen Analyse des Zivilrechts, p. 170; Wagner, 'Vor § 823 BGB', para 65.

<sup>114</sup> Cf. Wagner, 'Vor § 823 BGB', para 66.

<sup>&</sup>lt;sup>115</sup> United States v Carroll Towing Co., 159 F.2d 169 (1947). The formula is named after the author of the opinion, Judge Learned Hand, cf. Schäfer and Ott, Lehrbuch der ökonomischen Analyse des Zivilrechts, pp. 202-3.

<sup>&</sup>lt;sup>116</sup> Cf. Schäfer and Ott, *Lehrbuch der ökonomischen Analyse des Zivilrechts*, pp. 204-10.

boundaries of duties of care. First, a distinction must be made between the *enterprises*. The size of "B", the burden of the measure, depends not only on the amount of money to be paid but also on the paying organisation. For large and experienced enterprises in good financial standing, a safety measure is usually less of a burden than for small and medium-sized enterprises (SME) or start-ups with limited financial funds. When determining negligence, *individual* deficiencies are not considered. However, regarding the different "economic effects", one should distinguish between different *groups* of enterprises and e.g., make higher demands on wealthy market leaders.

Second, the *role of human intervention in the specific Human-AI cooperation* needs to be considered, as it determines "P", the probability of damage if the measure is not taken. In case of a Human-AI-Decision, the probability of damage depends on the skills of both the AI system and the human, and the extent to which the *principal's* measures can reduce this probability depends on the interplay between these skills: When the AI system and the human make their individual assessments in *distinctive* ways, e.g., because they consider different aspects of the situation, a lack of human control (Automation Complacency, Automation Bias) or an excessive human control (Algorithm Aversion) is more likely to cause damage. In contrast, if they make their decisions *similarly*, "P" will not be reduced considerably by assuring that the human correctly exercises his or her own judgment, e.g., by providing more explainability.

Third, one must assess the *specific impact* of the Human-AI-Decision. The AI system could e.g., suggest a translation the human needs to accept or modify. "L", the magnitude of the damage is obviously higher when the translation concerns the use of a medical drug as in case of a Christmas card for a business partner. <sup>122</sup>

<sup>&</sup>lt;sup>122</sup> Cf. the similar example by Fraunhofer IAIS, *Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher Intelligenz*, p. 12.



The distinction is in line with the idea of reducing "secondary costs": A loss is generally less burdensome if the loss is spread among lots of entities (Loss Spreading Method) or if it concerns an entity which does not suffer a lot from the loss because it is in good financial standing (Deep Pocket Method), cf. Calabresi, *The Costs of Accidents*, pp. 39-45; Taupitz, 'Ökonomische Analyse und Haftungsrecht', pp. 140-1.

<sup>&</sup>lt;sup>119</sup> Cf. BGH, 16 June 2009, VI ZR 107/08, BGHZ 181, 253-68, para 28 (juris); *Nettleship v Weston* 2 Q.B. 691 (1971) regarding English law; Schäfer and Ott, *Lehrbuch der ökonomischen Analyse des Zivilrechts*, p. 208; Wagner, '§ 823 BGB', paras 38-9.

<sup>&</sup>lt;sup>120</sup> "Verkehrskreise", cf. BGH, 2 October 2012, VI ZR 311/11, BGHZ 195, 30-42, para 7 (juris); Wagner, '§ 823 BGB', para 40.

<sup>&</sup>lt;sup>121</sup> Cf. Sommer, *Haftung für autonome Systeme*, pp. 238-42.

Finally, when weighing up the costs and benefits of a measure, it is important to remember that a measure that prevents one risk may at the same time create another risk. One example is training which increases *experience*: Experience may on the one hand lead to a better "subjective understanding" and therefore improve humans' capacity to re-evaluate the AI advice (cf. III.A.1.). On the other hand, if the system works well, experience may also lead to "routine use" and encourage Automation Complacency and Automation Bias (cf. III.C.). If the human sees the system err which comes with more experience - this may foster Algorithm Aversion. At the same time, seeing it making mistakes in specific situations could lead to a better understanding and therefore reduce such Aversion (cf. III.B.). Resolving such tradeoffs is challenging for both enterprises and courts. <sup>125</sup> The multidisciplinary approach will make it easier to define the duties of care, but not completely eliminate the difficulties.

## A. Principals' Duties of Care Under the Proposal for an AI Liability Directive

In 2022, the European Commission presented a Proposal for a Directive on adapting non-contractual civil liability rules to artificial intelligence (AILD Proposal). 126 However, in February 2025, the Commission announced its withdrawal, stating that an agreement was not foreseeable. 127 The Commission will now "assess whether another proposal should be tabled or another type of approach should be chosen". 128 As will be shown in the following, the application to Human-AI-Decisions is one point where the approach of the AILD Proposal should indeed be reconsidered.



<sup>&</sup>lt;sup>123</sup> Cf. note 96 above.

 $<sup>^{124}</sup>$  Cf. notes 102 and 103 above.

<sup>&</sup>lt;sup>125</sup> Cf. Schäfer and Ott, *Lehrbuch der ökonomischen Analyse des Zivilrechts*, p. 419 regarding general product liability law.

<sup>&</sup>lt;sup>126</sup> European Commission, 'Proposal for a Directive of the European Parliament and of the Council on adapting non-contractual civil liability rules to artificial intelligence (AI Liability Directive)', 28 September 2022, COM(2022) 496 final (AILD Proposal). On the same day, the European Commission presented a 'Proposal for a Directive of the European Parliament and of the Council on liability for defective products', COM(2022) 495 final (PLD Proposal). Unlike the proposed AILD, the new PLD was adopted in 2024, cf. Directive (EU) 2024/2853 of the European Parliament and of the Council of 23 October 2024 on liability for defective products and repealing Council Directive

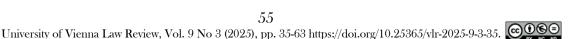
European Commission, 'Annexes to the Communication from the Commission to the European parliament, the Council, the European Economic and Social Committee and the Committee of the Regions, Commission work programme 2025', 11 February 2025, COM(2025) 45 final, Annex IV, p. 26.

<sup>&</sup>lt;sup>128</sup> Ibid, p. 26.

The Proposal does not contain a new form of AI liability but seeks to harmonise *fault-based liability* by laying down common rules on the disclosure of evidence and the burden of proof.<sup>129</sup> According to the Recitals and the explanations, the AILD Proposal does not want to touch on the definition of "fundamental concepts" like "fault".<sup>130</sup> However, Art. 4(1)(a) AILD Proposal presupposes that in national tort law "fault" could consist "in the non-compliance with a duty of care laid down in Union or national law directly intended to protect against the damage that occurred".<sup>131</sup> The AILD Proposal does not lay out duties of care itself.<sup>132</sup> It frequently refers to the definitions and obligations set out in the AI Act (cf. Art. 2, Art. 4(2) and (3) AILD Proposal).

Under the AILD Proposal, non-compliance with a duty of care can be both an *effect* and a *prerequisite* of a (rebuttable) presumption in favour of the victim: First, where a defendant fails to comply with a court order to disclose or preserve evidence, the court shall presume the defendant's breach of a duty of care (Art. 3(5) AILD Proposal - effect). Second, when the fault of the defendant is established and consists in the non-compliance with a duty of care laid down in Union or national law directly intended to protect against the damage that occurred, courts shall - under certain other conditions - presume a causal link between the fault and the output of the AI system (Art. 4(1) AILD Proposal - prerequisite). However, it is questionable whether the AILD Proposal applies to the Human-AI-Decisions analysed in this paper: Recital 15 suggests that it does not cover "liability claims when the damage is caused by a human assessment followed by a human act or omission, while the AI system only provided information or advice which was taken into account by the relevant human actor". This restriction has been criticised and the critic's arguments are convincing: As seen above, the impact of the human "in the loop" can be illusory. 134 The AILD Proposal's assumption that in such cases "it is possible to trace back the

<sup>&</sup>lt;sup>134</sup> Ibid., p. 13.



<sup>&</sup>lt;sup>129</sup> Cf. Art. 1 AILD Proposal. However, according to Art. 5(1) and (2), five years after the end of the transition period, the Commission must present a report that should in particular "evaluate the appropriateness of no-fault liability rules for claims against the operators of certain AI systems".

<sup>&</sup>lt;sup>130</sup> Cf. Explanatory Memorandum of the AILD Proposal, p. 11; Recital 10 AILD Proposal.

<sup>&</sup>lt;sup>131</sup> Cf. De Bruyne et al., "The European Commission's approach to extra-contractual liability and AI', p. 8: "the reliance on this concept [duty of care] in EU legislation is rather surprising".

<sup>&</sup>lt;sup>132</sup> Art. 2(9) AILD Proposal defines a "duty of care" as "a required standard of conduct, set by national or Union law, in order to avoid damage to legal interests recognised at national or Union law level, including life, physical integrity, property and the protection of fundamental rights".

<sup>&</sup>lt;sup>133</sup> Philipp Hacker, 'The European AI liability directives - Critique of a half-hearted approach and lessons for the future' (2023) 51 Computer Law & Security Review 105871, pp. 13-4.

damage to a human act or omission, as the AI system output is not interposed between the human act or omission and the damage" may be true but misses the actual problem of the victim. Proving the "factual" causation 185 between a human act or omission - which could also consist in the simple activation of the AI system - and a damage is usually possible. 136 The real difficulty consists in establishing whether the human act or omission breached a duty of care and whether the damage would not have occurred if this duty had been respected<sup>137</sup>. This difficulty also exists for victims of Human-AI-Decisions: As shown above, there are many factors to be considered when cooperating with AI systems and when organising such Human-AI cooperation. The victim's typical lack of insight into first, the AI system, second the Human-AI cooperation and third, the specific organisation this cooperation is integrated into, presents significant obstacles. 138 The wording of Art. 3 AILD Proposal would allow to extent the disclosing obligations and the presumption of non-compliance in case of a violation to these victims: It includes any "high-risk AI system that is suspected of having caused damage" (Art. 3(1) AILD Proposal) and any "claim for damages" (Art. 3(2) AILD Proposal). There is no requirement that the claim must be based on a faulty human behaviour *preceding* the AI output. In contrast, the *presumption* of causality in the case of fault contained in Art. 4 AILD Proposal seems to be tailored to fully automated AI systems: It only allows to presume the "causal link between the fault of the defendant and the output produced by the AI system". Taken literally, the presumption does not provide any relief if the faulty human behaviour succeeds the AI output. 40 Consequently, it would not cover claims against the human agent directly cooperating with the AI system. However, the wording of Art. 4 AILD Proposal would already allow applying the presumption to claims against the



<sup>&</sup>lt;sup>135</sup> Cf. Brüggemeier, *Haftungsrecht: Struktur, Prinzipien, Schutzbereich*, p. 28.

<sup>136</sup> Cf. Tianvu Yuan, 'Lernende Roboter und Fahrlässigkeitsdelikt' (2018) 9(4) RW 477-504, p. 493.

This latter aspect is part of the "legal" causation, cf. Brüggemeier, Haftungsrecht: Struktur. Prinzipien, Schutzbereich, p. 28.

<sup>138</sup> Cf. Hacker, 'The European AI liability directives', pp. 13-4: "it must be questioned whether it is really equally difficult to prove causality and non-compliance in cases of human intervention with and without AI involvement"; cf. equally De Bruyne et al., 'The European Commission's approach to extra-contractual liability and AI', p. 18 highlighting that - with regards to the Art. 8 PLD Proposal -"the complexity at stake is rather of an organisational nature and/or relates to the information asymmetry that the consumer endures concerning the apportion of responsibilities between the various actors at stake".

Cf. Hacker, 'The European AI liability directives', p. 14: "Hence, the presumption of noncompliance contained in Article 3(5) AILD Proposal should apply equally if a human agent took, or failed to take, the final decision leading to the damage caused by the AI output".

<sup>&</sup>lt;sup>140</sup> Cf. ibid., p. 23.

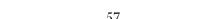
*principals* of these human agents: *Their* faulty behaviour, e.g., a violation of a duty to inform the human, usually *precedes* the AI output. <sup>141</sup> Nevertheless, a potential new proposal or other type of approach to AI Liabilityshould at least reconsider this limitation. <sup>142</sup>

#### B. Summary and Perspectives

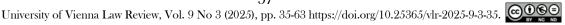
Enterprises increasingly make use of Human-AI-Decisions. When such decisions result in damage the question of the enterprise's liability arises. This paper has focused on enterprises' non-contractual and *fault-based* liability which generally requires a violation of a *duty of care* by the principal. When it comes to Human-AI-Decisions, it is challenging for courts to define these duties of care. The paper has shown how findings from other disciplines can help judges to determine "possible" and "reasonable" measures principals must take to prevent specific risks of Human-AI-Decision-making. There are three main sources of such risks: Automation Complacency and Automation Bias on the one hand and Algorithm Aversion on the other hand. Scientists have conducted numerous studies on these phenomena and have identified mitigating factors. Their findings are part of the "state of scientific and technical knowledge" and need to be considered when setting legal safety standards. Nonetheless, principals do not need to take every "possible" measure. To determine if a measure can be reasonably expected, knowledge from other disciplines, in particular from economics, is again of considerable value.

At the same time, this paper has shown that some challenges remain, especially when it comes to resolving trade-offs associated with a particular measure. The now abandoned AILD Proposal would have provided little relief in this regard. Thus, the question as to whether liability for AI systems should not be tightened *de lege ferenda* remains pressing.<sup>143</sup>

<sup>&</sup>lt;sup>143</sup> Cf. on this issue e.g., Gerhard Wagner, 'Verantwortlichkeit im Zeichen digitaler Techniken' (2020) 71 VersR 717-41, pp. 734-41; Zech, 'Entscheidungen digitaler autonomer Systeme', pp. A 87-110; Beckers and Teubner, *Three Liability Regimes for Artificial Intelligence*, pp. 45-166; Benedetta Cappiello, *AI-systems and non-contractual liability* (Torino, Giappichelli: 2022), pp. 45-96; Mayrhofer, *Außervertragliche Haftung für fremde Autonomie*, pp. 370-443 and 'Product liability in the age of AI – Proposal for a "two track" solution' (2024) 33(1) *Revista Electrónica de Direito* 105-127; cf., also, Philipp Hacker, *Proposal for a directive on adapting non-contractual civil liability rules to artificial intelligence, Complementary impact assessment* (2025), available at







Whether the Recitals can limit the scope of application of a directive cannot be discussed here; for an analysis of the significance of recitals in the methodology of EU law, cf. Tobias Gumpp, 'Stellenwert der Erwägungsgründe in der Methodenlehre des Unionsrechts' (2022) 8(4) ZfPW 446-76.

<sup>&</sup>lt;sup>142</sup> Cf. Hacker, "The European AI liability directives", p. 23: "What would be needed is a presumption that the act/omission of the user [..] caused the damage".

### V. Bibliography

Alon-Barkat, Saar and Busuioc, Madalina, 'Human-AI Interactions in Public Sector Decision-Making: "Automation Bias" and "Selective Adherence" to Algorithmic Advice' (2023) 33(1) Journal of Public Administration Research and Theory 153-69

Ball, Matthew and Callaghan, Vic, 'Explorations of Autonomy', in 2012 Proceedings of the 8th International Conference on Intelligent Environments (Guantajo: IEE Xplore, 2012) 114-21

Bauer, Kevin, et al., 'Die KI braucht bei der Bankberatung immer noch menschliche Hilfe', Börsen-Zeitung, available at https://www.boersenzeitung.de/kapitalmarktforschung/in-der-bankberatung-braucht-die-ki-menschlichehilfe-90edbb42-86a4-11ed-a311-f90ecc32c8e4 (last accessed 25 April 2025)

Beckers, Anna and Teubner, Gunther, Three Liability Regimes for Artificial Intelligence (Oxford: Hart, 2022)

Bomhard, David and Merkle, Marieke, 'Europäische KI-Verordnung - Der aktuelle Kommissionsentwurf und praktische Auswirkungen (2021) 1(6) Recht Digital (RDi) 276-282

Börütecene, Ahmet and Löwgren, Jonas, 'Designing Human-Automation Collaboration for Predictive Maintenance', in Companion Publication of the 2020 ACM Designing Interactive Systems Conference (New York: Association for Computing Machinery, 2020) 251-6

Brüggemeier, Gert Haftungsrecht: Struktur, Prinzipien, Schutzbereich (Berlin et al.: Springer, 2006)

Calabresi, Gudio, *The Costs of Accidents* (Cumberland: Yale University Press, 1970)

Cappiello, Benedetta, AI-systems and non-contractual liability (Torino, Giappichelli: 2022)

Choi, Benard C.K. and Pak, Anita W.P., 'Multidisciplinarity, interdisciplinarity, and transdisciplinarity in health research, services, education and policy: 1. Definitions, objectives, and evidence of effectiveness' (2006) 29(6) Clinical and Investigative *Medicine* 351-64

Cummings, Mary L., 'Automation Bias in Intelligent Time Critical Decision Support Systems' (2004) AIAA 1st Intelligent Systems Technical Conference, 20-22

https://www.europarl.europa.eu/RegData/etudes/STUD/2024/762861/EPRS\_STU(2024)762861\_E N.pdf (last accessed 25 April 2025).





September 2004, Chicago, IL, American Institute of Aeronautics and Astronautics 1-6

Datenethikkommission der Bundesregierung, Gutachten (2019), available at https://www.bmi.bund.de/SharedDocs/downloads/DE/publikationen/themen/itdigitalpolitik/gutachten-datenethikkommission.html (last accessed 25 April 2025)

De Bruyne, Jan, et al., 'The European Commission's approach to extra-contractual liability and AI - An evaluation of the AI liability directive and the revised product liability directive' (2023) 51 Computer Law & Security Review 105894

Deutsches Institut für Normung and Deutsche Kommission Elektrotechnik, Elektronik, Deutsche Normungsroadmap Künstliche Intelligenz, 2nd edn. (2022), available at https://www.din.de/de/forschung-und-innovation/themen/kuenstlicheintelligenz/fahrplan-festlegen (last accessed 25 April 2025)

Dietvorst, Berkeley J., et al., 'Algorithm aversion: People erroneously avoid algorithms after seeing them err' (2015) 144(1) Journal of Experimental Psychology: General 114-26

Essers, Josef and Schmidt, Eike, Schuldrecht, vol. I/1, 8th edn. (Heidelberg: C.F. Müller, 1995)

Expert Group on Liability and New Technologies, 'Liability for Artificial Intelligence and Other **Emerging** Digital Technologies' (2019).available https://op.europa.eu/en/publication-detail/-/publication/1c5e30be-1197-11ea-8c1f-01aa75ed71a1/language-en (last accessed 25 April 2025)

Eykholt, Kevin et al., 'Robust Physical-World Attacks on Deep Learning Visual Classification', in 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (Salt Lake City: IEEE, 2018) 1625-34

Fawcett, Jacqueline, 'Thoughts About Multidisciplinary, Interdisciplinary, and Transdisciplinary Research' (2013) 26(4) Nursing Science Quarterly 376-9

Fraunhofer IAIS, Leitfaden zur Gestaltung vertrauenswürdiger Künstlicher Intelligenz (2021).available https://www.iais.fraunhofer.de/de/publikationen/studien/2021/ki-pruefkatalog.html (last accessed 25 April 2025)

Fügener, Andreas et al., 'Cognitive Challenges in Human-Artificial Intelligence Collaboration: Investigating the Path Toward Productive Delegation' (2022) 33(2) Information Systems Research 678-96



Goddard, Kate, et al., 'Automation bias: a systematic review of frequency, effect mediators, and mitigators' (2012) 19(1) Journal of the American Medical Informatics Association 121-7

Grace, Katja, et al., 'Viewpoint: When Will AI Exceed Human Performance? Evidence from AI Experts' (2018) 62 Journal of Artificial Intelligence Research 729-54

Gsenger, Rita and Strle, Tomam, 'Trust, Automation Bias and Aversion: Algorithmic Decision-Making in the Context of Credit Scoring' (2021) 19(4) Interdisciplinary Description of Complex Systems 542-60

Gumpp, Tobias, 'Stellenwert der Erwägungsgründe in der Methodenlehre des Unionsrechts' (2022) 8(4) Zeitschrift für die gesamte Privatrechtswissenschaft (ZfPW) 446-76

Hacker, Philipp, Proposal for a directive on adapting non-contractual civil liability rules to artificial intelligence, Complementary impact assessment (2025), available at https://www.europarl.europa.eu/RegData/etudes/STUD/2024/762861/EPRS\_STU( 2024)762861\_EN.pdf (last accessed 25 April 2025)

Hacker, Philipp, 'The European AI liability directives - Critique of a half-hearted approach and lessons for the future' (2023) 51 Computer Law & Security Review 105871

Hacker, Philipp, 'Verhaltens- und Wissenszurechnung beim Einsatz von Künstlicher Intelligenz' (2018) 9(3) Rechtswissenschaft (RW) 243-88

High-Level Expert Group on Artificial Intelligence, 'Ethics Guidelines for Trustworthy Aľ (2019),available https://digitalat strategy.ec.europa.eu/en/library/ethics-guidelines-trustworthy-ai (last accessed 25 April 2025)

Hilgendorf, Eric, 'Bedingungen gelingender Interdisziplinarität' (2010) 65(19) Juristenzeitung (JZ) 913-22

Hofmann, Franz, 'Disziplinarität, Intradisziplinarität und Interdisziplinarität am Beispiel der Grundsätze "mittelbarer Verantwortlichkeit" (2018) 73(15-6)JuristenZeitung (JZ) 746-54

Jussupow, Ekaterina et al., 'Augmenting Medical Diagnosis Decisions? An Investigation into Physicians' Decision-Making Process with Artificial Intelligence' (2021) 32(3) *Information System Research* 713-35

Kahneman, Daniel et al., 'NOISE' (2016) 40 Harvard Business Review 38-46



Kaminski, Andreas, 'Gründe geben. Maschinelles Lernen als Problem der Moralfähigkeit von Entscheidungen', in Klaus Wiegerling et al. (eds.), *Datafizierung* und Big Data (Wiesbaden: Springer VS, 2020) 151-74

Koziol, Helmut (ed.), Basic Questions of Tort Law from a Comparative Perspective (Vienna: Sramek, 2015)

Koziol, Helmut (ed.), Basic Questions of Tort Law from a Germanic Perspective (Vienna: Sramek, 2012)

Koziol, Helmut (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015)

Koziol, Helmut, 'Comparative Conclusions' in Helmut Koziol (ed.), Basic Questions of Tort Law from a Comparative Perspective (Vienna: Sramek, 2015) 685-838

Koziol, Helmut, 'Concluding Remarks', in Helmut Koziol (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015) 182-95

Koziol, Helmut, et al., 'Liability for Agents and Agents' Liability', in Helmut Koziol (ed.), Comparative Stimulations for Developing Tort Law (Vienna: Sramek, 2015) 182-95

Larson, Erik J., The myth of artificial intelligence (Cambridge et al.: Harvard University Press, 2021)

Lee, Min Kyung, 'Understanding perception of algorithmic decisions: Fairness, trust, and emotion in response to algorithmic management' (2018) 5(1) Big Data & Society 1-16

Liel, Yotam and Zalmanson, Lior, 'Turning Off Your Better Judgment - Conformity to Algorithmic Recommendations' (Working paper) (2022), available https://www.researchgate.net/publication/366412145\_Turning\_Off\_Your\_Better\_Ju dgment\_-Conformity\_to\_Algorithmic\_Recommendations (last accessed 25 April 2025)

Mayrhofer, Ann-Kristin, 'Product liability in the age of AI – Proposal for a "two track" solution' (2024) 33(1) Revista Electrónica de Direito 105-127

Mayrhofer, Ann-Kristin, Außervertragliche Haftung für fremde Autonomie (Tübingen: Mohr Siebeck, 2023)

Milner, Greg, 'Death by GPS: Are Satnavs changing our brains?', The Guardian, available at https://www.theguardian.com/technology/2016/jun/25/gps-horror-storiesdriving-satnay-greg-milner (last accessed 25 April 2025)



Möllers, Thomas M. J., Juristische Methodenlehre, 5th edn. (Munich: C.H. Beck, 2023)

Mori, Mashiro, 'The Uncanny Valley' (2012) 19(2) IEEE Robotics & Automation *Magazine* 98-100

Mosier, Kathleen L. and Skitka, Linda J., 'Human Decision Makers and Automated Decision Aids: Made for Each Other?', in Raja Parasuraman and Mustapha Mouloua (eds.), Automation and human performance: Theory and application (Mahwah: Lawrence Erlbaum, 1996) 201-20

Parasuraman, Raja and Manzey, Dietrich H., 'Complacency and Bias in Human Use of Automation: An Attentional Integration' (2010) 53(3) Human Factors 381-410

Schäfer, Hans-Bernd and Ott, Claus, Lehrbuch der ökonomischen Analyse des Zivilrechts, 6th edn. (Berlin et al.: Springer Gabler, 2020)

Sommer, Martin, *Haftung für autonome Systeme* (Baden-Baden: Nomos, 2020)

'Die Zivilrechtswissenschaft und ihre Methodik - zu rechtsanwendungsbezogen und zu wenig grundlagenorientiert?' (2014) 214(1-2) Archiv für die civilistische Praxis (AcP) 7-54

Susanne Horner and Markus Kaulartz, 'Haftung 4.0' (2016) 32(1) Computer und Recht (CR) 7-14

Taupitz, Jochen, 'Ökonomische Analyse und Haftungsrecht – Eine Zwischenbilanz' (1996) 196(1/2) Archiv für die civilistische Praxis (AcP) 114-67

Teodorescu, Mike H. M., et al., 'Failures of Fairness in Automation Require a Deeper Understanding of Human-ML Augmentation' (Minneapolis: University of Minnesota 2021) 45(3) MIS Quarterly 1483-500

Wagner, Gerhard, 'Roboter als Haftungssubjekte? Konturen eines Haftungsrechts für autonome Systeme', in Florian Faust and Hans-Bernd Schäfer (eds.), Zivilrechtliche und rechtsökonomische Probleme des Internet und der künstlichen Intelligenz (Tübingen: Mohr Siebeck, 2019) 1-39

Wagner, Gerhard, 'Verantwortlichkeit im Zeichen digitaler Techniken' (2020) 71 Versicherungsrecht (VersR) 717-41

Wagner, Gerhard, 'Vor § 823 BGB' and '§ 823 BGB', in Jürgen Säcker et al. (eds.), Münchener Kommentar zum Bürgerlichen Gesetzbuch, 9th edn., 13 vols. (Munich: C.H. Beck, 2024), vol. VII

Yeomans, Michael, et al., 'Making sense of recommendations' (2019) 32(4) Journal of Behavioral Decision Making 403-14



Tianyu, 'Lernende Roboter und Fahrlässigkeitsdelikt'  $(2018) \quad 9(4)$ Rechtswissenschaft (RW) 477-504

Zech, Herbert, 'Entscheidungen digitaler autonomer Systeme: Empfehlen sich Regelungen zu Verantwortung und Haftung?", in Ständige Deputation des Deutschen Juristentages (ed.), Verhandlungen des 73. Deutschen Juristentages (Munich: C.H. Beck, 2020) A 1-112