



AI Risk Assessment: A Scenario-Based, Proportional Methodology for the AI Act

Claudio Novelli¹ · Federico Casolari¹ · Antonino Rotolo¹ · Mariarosaria Taddeo^{2,4} · Luciano Floridi^{1,3}

Received: 14 September 2023 / Accepted: 15 February 2024
© The Author(s) 2024

Abstract

The EU Artificial Intelligence Act (AIA) defines four risk categories for AI systems: unacceptable, high, limited, and minimal. However, it lacks a clear methodology for the assessment of these risks in concrete situations. Risks are broadly categorized based on the application areas of AI systems and ambiguous risk factors. This paper suggests a methodology for assessing AI risk magnitudes, focusing on the construction of real-world risk scenarios. To this scope, we propose to integrate the AIA with a framework developed by the Intergovernmental Panel on Climate Change (IPCC) reports and related literature. This approach enables a nuanced analysis of AI risk by exploring the interplay between (a) risk determinants, (b) individual drivers of determinants, and (c) multiple risk types. We further refine the proposed methodology by applying a proportionality test to balance the competing values involved in AI risk assessment. Finally, we present three uses of this approach under the AIA: to implement the Regulation, to assess the significance of risks, and to develop internal risk management systems for AI deployers.

Keywords Risk assessment · AI Act · IPCC · Proportionality · Artificial Intelligence

1 Introduction: From Broad Scopes to Risk Scenarios

The European Artificial Intelligence Act (AIA) introduces a risk-based regulatory framework for AI systems (AIs), categorising them into four levels of risk: unacceptable, high, limited, and minimal. The legislator allocates regulatory burdens to

✉ Claudio Novelli
claudio.novelli@unibo.it

¹ Department of Legal Studies, University of Bologna, Via Zamboni, 27/29, Bologna 40126, Italy

² Oxford Internet Institute, University of Oxford, 1 St Giles', Oxford, OX1 3JS, UK

³ Digital Ethics Center, Yale University, 85 Trumbull Street, New Haven, CT 06511, USA

⁴ Alan Turing Institute, British Library, 96 Euston Rd, London, NW1 2DB, UK

AIs' providers so that the greater the risk posed by AIs, the greater the legal safeguards to minimise it.

However, the AIA, in its capacity as a broad risk-based regulation, lacks a detailed risk assessment methodology for identifying risks, relying on a static view of AI risk. The four risk categories of AIs are mainly based on their technological features and the broad application areas. In short, AI is mostly seen as a product, akin to the EU product safety legislation.^{1,2} In doing so, the AIA does not consider the interaction among hazard sources, vulnerability profiles, and exposed values, but treats them as stand-alone technical standards. This is coupled by the lack of a proportionality judgement between the risk mitigation measures and the principles and rights involved. As a result, the impact that AIs may have on European fundamental values and interests seems predetermined.

The AIA may misestimate the magnitude of AI risks—i.e., the likelihood of detriment and severity of consequences on values like health, safety, privacy, and others—and make the overall legal framework ineffective, that is, with rules that are either too stringent or too soft for the actual applications of specific AIs. The root of this problem is that the AIA has not yet progressed to the standardization phase, which is crucial for developing detailed guidelines. To address this challenge, the paper suggests a risk assessment methodology aimed at improving the accuracy and relevance of the AIA's provisions. This methodology is not only applicable to enhancing the AIA but could also be beneficial for any other risk-based legislative framework governing AI.

We propose a risk assessment model that identifies and combines specific risk factors influencing real-world AI application scenarios.³ While some legal arguments have been presented, suggesting a reading of the AIA's risks approach considering tort law (Chamberlain, 2022), we draw from research and policy reports on climate change risk. In particular, we refer to the framework developed by the Intergovernmental Panel on Climate Change (IPCC) working groups and refined by the subsequent literature (Simpson et al., 2021). Accordingly, the risk of an event is assessed by the interplay between (1) determinants of risk (i.e., hazard, exposure, vulnerability, and responses), (2) individual drivers of determinants, and (3) other types of risk (i.e., extrinsic, and ancillary risks). This framework can provide a more accurate risk magnitude of AIs under a specific scenario. This is a measure defined based on hazard chains, the trade-off among impacted values, the aggregation of vulnerability profiles, and the contextualisation of AI risk with risks from other sectors.

This qualitative analysis is grounded on a quantitative assessment. In fact, the risk magnitude should be assessed by weighing the fundamental values (positively and negatively) affected by AIs against the intensity of the interference of AIA's risk

¹Cf. European Commission, Explanatory Memorandum to AIA, para 1.3.

²Still, it is compatible with a risk management standard employed in safety-critical industries in other legislations. An example we shall consider is the United Kingdom's 'As Low As Reasonably Practicable' (ALARP) principle, which we will use as a framework to introduce the AIA's risk classification.

³We shall use the expression 'risk factors' to refer in a general way to all variables potentially able to increase or decrease the risk of an event. We shall specify its meaning by referring to determinants and drivers.

containment measures on the same values. This type of judgment for interference between constitutional principles is the object of the proportionality test by Robert Alexy (Alexy, 2002). The outcome of the test would indicate whether a risk category is appropriate for an AI under a specific risk scenario or whether it introduces grossly disproportionate limitations and trade-offs for competing values.⁴

We consider three uses of this semi-quantitative risk framework under the AIA. First, for implementing the AIA with a different way of categorizing the risk. This means moving from a vertical approach, which categorizes AIs according to risk factors in isolation, to a horizontal approach, which categorizes AIs according to the interactions between various risk factors across different real-world scenarios. This nuanced assessment strategy is particularly suited for the delegated acts phase of the AIA, where implementation standards are established by the Commission with contributions from stakeholders and the scientific community. This phase offers an optimal setting for adopting the risk methodology recommended in this paper. Second, for enabling deployers of AIs initially classified as high-risk, based on the AIA's predefined criteria, to dispute this categorisation. According to Recital 32, they must demonstrate that considering the severity, intensity, likelihood, duration, and potential targets, the overall risk is not significant.⁵ Third, for setting up the internal risk management system for high-risk AI deployers as mandated by Article 9 of the AIA. This requirement involves identifying risks, considering AI uses and misuses, and evaluating new risks from post-market data. The second and third uses are seen as more feasible given the current legislative stage of the AIA.

The article is structured as follows. Section 2 presents the risk-based regulation of the AIA, bridging the risk model within the EU proposal and the ALARP principle. Section 3 discusses the strengths and weaknesses of the AIA risk-based regulation. Section 4 shows how to overcome the AIA's model gaps by using the IPCC framework for climate change risk assessment updated by the relevant literature. Section 5 offers a quantitative support to the model through a proportionality test. Section 6 discusses the three potential uses of our semi-quantitative proposal under the AIA. Section 7 outlines the advantages of modifying the AIA's strategy towards risk in its enforcement and regulation of General Purpose AI (GPAI). Section 8 concludes the article.

2 AIA's Risk-Based Regulation

Generally, risk-based regulations consist of (at least) three phases: assessment, categorisation, and management (Millstone et al., 2004). In this article, we shall focus more on the first two phases and less on the AIA's risk management system, that is, legal safeguards and requirements.

⁴Robert Alexy has been a pioneer in developing quantitative approaches to balancing principles. His work is widely respected for its clarity and rigor, making him a valuable reference for scholars and practitioners alike. Additionally, his Weight Formula offers a promising framework for implementing the fundamental rights impact assessment (RIA) recently introduced in the draft legislation.

⁵The risk significance evaluation has been introduced by the European Parliament compromise text on the AIA from 11 May 2023.

The AIA relies on the traditional conception that risk is the likelihood of converting a source of hazard into actual loss, injury, or damage.⁶ Sources of danger are those uses of AI that are most likely to compromise safety, health, and other values.⁷ Being the likelihood of damage, risk can be expressed through the ratio between hazard and safeguards so that, as the safeguards increase, the risk quotient decreases:

$$\text{risk} = \frac{\text{hazard}}{\text{safeguards}}$$

The risk may become untenable if safeguards do not offset severe hazards. The regulatory intervention should be proportionate to the hazards net of safeguards. Risk tolerance thresholds—in the AIA, the risk categories—indicate which risks are accepted without (strong) precautions and which instead require (further) mitigation practices.

In the AIA, the benchmark to calculate the risk of AIs is their potential adverse impact on health, safety, and EU fundamental rights. As a result, the AIA classifies AIs according to four risk categories: unacceptable, high, limited, and minimal (Kaplan & Garrick, 1981).⁸ Stricter requirements are prescribed for suppliers and users of riskier AIs. This is explicitly stated in Recital 14 of the draft:

In order to introduce a proportionate and effective set of binding rules for AI systems, a clearly defined risk-based approach should be followed. That approach should tailor the type and content of such rules to the intensity and scope of the risks that AI systems can generate.⁹

This is why the AIA modulates the legal requirements to make the risk of deploying AIs at least tolerable. The tolerance thresholds that constitute the AIA's risk categorisation is compatible with the ALARP principle. ALARP is a general principle in UK law for risk management systems in safety-critical industries (Abrahamsen et al., 2018; Jones-Lee & Aven, 2011), and in the UK health system.¹⁰

⁶This conception can easily be deduced from some sections of the draft, for instance, Recital 32 referring to high-risk systems: “[...] high-risk AI systems other than those that are safety components of products...it is appropriate to classify them as high-risk if, in the light of their intended purpose, they pose a high risk of harm to the health and safety or the fundamental rights of persons, taking into account both the severity of the possible harm and its probability of occurrence [...]”.

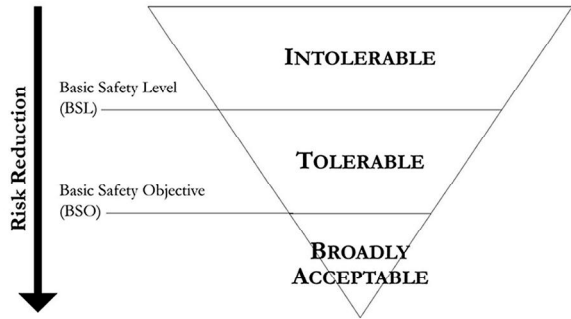
⁷The AI Act will with the OECD definition of AI: “An AI system is a machine-based system that, for explicit or implicit objectives, infers, from the input it receives, how to generate outputs such as predictions, content, recommendations, or decisions that [can] influence physical or virtual environments. Different AI systems vary in their levels of autonomy and adaptiveness after deployment”. Source: <https://www.euractiv.com/section/artificial-intelligence/news/ai-act-eu-policymakers-nail-down-rules-on-ai-models-butt-heads-on-law-enforcement/>.

⁸The text refers to three categories, but a fourth sub-category of high-risk systems can be derived from the presence of lighter obligations.

⁹Proposal for a Regulation of the European Parliament and of the Council Laying Down Harmonised rules on Artificial Intelligence (Artificial Intelligence Act) and Amending Certain Union Legislative Acts, Com/2021/206 final, Recital 14.

¹⁰UKHSE. Risk management: Expert guidance - ALARP at a glance. <https://www.hse.gov.uk/managing/theory/alarplance.htm>. However, the principle is also recognised in other legal systems, like the US, sometimes under the formula “as low as reasonably achievable” (ALARA)

Fig. 1 The tolerance ranges of the ALARP principle¹¹



ALARP-type approaches involve a proportionality review of risk reduction measures so that they are not exorbitant to the improvement gained (Bai & Jin, 2016). Typically, ALARP provides the following risk tolerance ranges (Hurst et al., 2019):

Although the transposition into EU law of ALARP is limited¹² and controversial¹³, AIA’s risk categories overlap with the tolerance ranges shown in Fig. 1.¹⁴ These risk categories can be summarised as follows.¹⁵

Unacceptable risk includes (AIA, Title II):

- AIs that may cause significant harm through (a) subliminal manipulation of individuals’ consciousness that distorts their behaviour¹⁶ or (b) exploitation of vulnerabilities—age, physical or mental disability—of a specific group of people that distorts the behaviour of its members.

¹¹This is a simplified version of the figure found in (Hurst et al., 2019). Note that the size of the three categories of the inverted pyramid is related to the severity and not to the numerousness of the relative risks.

¹²A specific reference may be found in the EU legislation on medical devices: cf. Annex I, Regulation (EU) 2017/745 of the European Parliament and of the Council of 5 April 2017 on medical devices, [2017] OJ L117/1.

¹³In particular, the debate developed when the European Commission argued that the use of ALARP (in the homologous version of SFAIRP) in the UK’s Health & Safety at Work Act was not consistent with the European “Framework Directive” for occupational safety and health (Directive 89/391/EEC), asking thus the Court of Justice to declare that the Member State failed to fulfil its obligation to correctly transpose the Directive. However, the European Court of Justice, without taking a specific position on the compatibility of ALARP with the Directive’s provision, dismissed the action brought by the European Commission, maintaining that the EU institution did not clearly identify the legal standard enshrined in the Directive that the UK failed to implement (Case C-127/05 *Commission v UK* EU: C:2007:338, para 58).

¹⁴Indeed, a textual reference to ALARP can be found in the section where the AIA describes the mandatory risk management system for high-risk systems: “In identifying the most appropriate risk management measures, the following shall be ensured: (a) elimination or reduction of risks as far as possible through adequate design and development; [...]”.

¹⁵This list is updated to the May 2023: Draft Compromise Amendments on the Draft Report (COM (2021)0206 - C9 0146/2021 - 2021/0106(COD)). However, the text is not yet conclusive.

¹⁶Literature indicates that the notion of subliminal techniques remains ambiguous under the AI Act (Neuwirth, 2022).

- AIs for social scoring that evaluate or classify natural persons or groups based on their social behaviour when social scoring leads to detrimental or unfavourable treatment (a) in social contexts that are unrelated to the contexts in which the data was originally generated or collected; (b) detrimental or unfavourable treatment are unjustified or disproportionate to the social behaviour of natural persons or groups.
- AIs for biometric categorisation that categorise natural persons according to sensitive or protected attributes or characteristics (e.g., gender, ethnicity, political orientation, religion, disability) or based on the inference of those attributes or characteristics.¹⁷
- AIs for risk assessments of natural persons or groups to assess the risk for offending or reoffending or for predicting the occurrence or reoccurrence of (an actual or potential) criminal or administrative offence based on assessing personality traits and characteristics, such as the person's location, past criminal behaviour of natural persons or groups of natural persons.
- AIs for inferring emotions of a natural person in the areas of law enforcement, border management, in workplace and education institutions.

High-risk includes (AIA, Title III):

- AIs used as safety components of products covered by the European New Legislative Framework (NLF) and other harmonised European regulations (Annex II, Sects. A and B). Regulated areas include, e.g., automotive, fossil fuels and medical devices.
- AIs deployed in (a) biometric identification (when this is not forbidden) (b) management and operation of critical infrastructure, (c) education and vocational training, (d) employment, worker management and access to self-employment, (e) access to and enjoyment of essential private services and public services and benefits (e.g., healthcare), (f) law enforcement (g) migration, asylum and border control management, (h) administration of justice and democratic processes (Annex III).

Limited risk includes (AIA, Title IV):

- AIs that interact with natural persons, e.g., chatbots, when this is not obvious from the circumstances and the context of use or is not permitted by law to detect, prevent and investigate criminal offences.
- AIs that generate or manipulate images, audio, or video to simulate people, objects, places or other existing entities or events (i.e., deep fakes).

Minimal risk includes (AIA, Title IX):

- Residual AIs, some examples are AIs for video games or spam filters.

AIs posing unacceptable risks fall into the ALARP 'Intolerable' risk range, i.e., situations whose risk cannot be justified except in extraordinary circumstances. Under the AIA, specific exempt circumstances, like terrorist attacks, allow the time-

¹⁷When the exonerating circumstances provided for in Articles 5(1)(d) and 5(2)(4) are not met.

limited use of AIs for remote biometric identification in publicly accessible spaces for law enforcement (Article 5(d)).

AIs posing high and limited risks fall into the ‘Tolerable’ risk range. That is where the ALARP principle comes fully into play: risk is tolerated only if all reasonably practicable mitigation measures are implemented. However, what counts as ‘reasonably practicable’ might be tricky to determine. A predominant interpretation is that: “Efforts to reduce risk should be continued until the incremental sacrifice is grossly disproportionate to the value of the incremental risk reduction achieved. Incremental sacrifice is defined in terms of cost, time, effort, or other expenditures of resources” (Baybutt, 2014).

This judgement should therefore consider the expected utility of risk containment. In the AIA, reasonable efforts consist of the legal requirements and guarantee mechanisms that providers (and deployers) must comply with to place high-risk AIs on the single market (Article 6 et seq.). We shall analyse the ALARP principle, seeking to improve its enforcement in the AIA, in greater detail in Sect. 5.

AIs posing minimal risks fall into the ALARP ‘Broadly Accepted’ risk range. In these cases, the risk is tolerable enough that no specific intervention is required, except to ensure compliance with good practices. This is also what the AIA prescribes by encouraging the adoption of voluntary codes of conduct either by individual providers of AIs or by their representative organisations (Article 69).

Much of the legal framework concerns high-risk AIs, prescribing conformity assessment procedures, technical documentation, and certification duties to place them on the market (e.g., Article 43). Sometimes these safeguards involve post-market monitoring (e.g., Article 61). The other three risk categories produce fewer and simpler regulatory burdens: AIs that pose unacceptable risks are prohibited (Article 5), those that pose limited risk trigger a general transparency obligation (Article 52), while for those that pose minimal risks the AIA fosters voluntary codes of conduct (Article 69).¹⁸ An exception to these rules is provided in the AIA insofar as it requires the Member States to introduce regulatory sandboxes: controlled environments in which AIs can be developed and tested for a limited time, before putting them on the market, prioritising small providers and start-ups (Article 53 seq.).

3 Strengths and Weaknesses of the AIA’s Risk Regulation

The supranational legislator expects the regulation of AI to increase legal certainty in this field and to promote a well-functioning internal market: reliable for consumers, attractive for investment, and technologically innovative.¹⁹ This might trigger the Brussels effect, ensuring a competitive advantage over other international policy-makers while shaping their regulatory standards (Bradford, 2020). Nevertheless, should the AIA prove to be unsustainable or ineffective, the EU may lose its attractiveness for the production and commercialisation of AI

¹⁸Codes of conduct can be created by individual providers or their representative organisations.

¹⁹These are explicitly stated objectives of the AIA draft (p. 3).

technologies. To prevent this, the AIA must introduce norms that promote safety while not disincentivising the production or deployment of AIs.²⁰ In this regard, the AIA's risk-based approach has its strengths and weaknesses. Let us start with the strengths.

First, risk-based regulations rationalise governance interventions by setting their priorities and objectives. Well-delineated priorities and objectives facilitate accountability mechanisms towards the policy-maker (Black, 2010b). In this respect, the AIA declares its priorities and objectives: the protection of the fundamental values and rights of the Union and the development of the AI market.

Second, risk-based regulations facilitate the fair distribution of resources (e.g., for supervision and certification) and costs. For example, costs are distributed according to the specific risks posed to a target community, and they are so transparently, as the criteria for distributing resources and costs are made evident in the regulation (Black, 2010a). As the compliance cost is proportional to the risk, AIA introduces a kind of Pigouvian tax on the negative externalities of high-risk AIs (Baumol, 1972). To be acceptable, the AIA should allocate costs and resources efficiently among market players. However, the AIA does not consistently distribute resources in the best possible way, as we shall see when discussing its weaknesses.

Third, risk-based regulations cope with the uncertainty of phenomena—i.e., “when there is a lack of knowledge in qualitative or quantitative terms”^{21,22}—for example, by qualifying predictions about the occurrence of specific hazards probabilistically (Rothstein et al., 2013). Moreover, risk-based regulations adapt to the political context or technological and market changes (Black & Baldwin, 2010).²³ In this regard, the AIA offers the possibility of updating its list of risky AIs at Articles 84–85. Unfortunately, the current version allows new AIs to be added only if they fall within the already established scopes. For this reason, some suggestions have been made to include reviewable risk categorisation criteria (Smuha et al., 2021).

By contrast, one of the main limitations of the AIA is the uncertainty about criteria for reviewing risk categorisation, which depends instead on the broad scopes of AIs. AI providers may be reluctant to invest in the EU's AI market due to the perceived rigidity of the AIA guidelines and the absence of a mechanism for revising or adapting limitations and prohibitions as technological advancements occur. These advancements could potentially make certain AI systems that are currently considered risky less so. The AIA may preclude adapting risk categorisation to the interplay of hazard sources, vulnerability profiles of the exposed

²⁰Of course, other factors will determine the success of the European AI strategy, like taxation and administrative efficiency. However, in this paper, we will only address the regulatory framework, namely the risk-categorisation of the AIA.

²¹van der Heijden J. Risk as an approach to regulatory governance: An evidence synthesis and research agenda. *SAGE Open* 2021;11(3):215.

²²Sometimes, the concepts of risk and uncertainty are kept separate, the former being considered calculable and the latter not. For this purpose, the distinction between *epistemic* and *aleatory* uncertainty may be relevant, with only the latter being effectively addressable through risk assessment. On this, see Renn O. Risk governance: coping with uncertainty in a complex world. London: Routledge 2011; 368.

²³At the same time, excessive uncertainty must be seen as a limitation of any risk model.

community, or values and interests at stake. No doubt, the model enshrined in the AIA heavily relies on a fundamental rights-based approach—as confirmed by the amendment introducing a fundamental rights impact assessment (AIA, Article 29a)—which characterizes the entire structure of the legislative proposal and, more broadly, the most recent pieces of legislation adopted at EU level in the digital context (Ufert, 2020). However, as legal compliance always comes at a cost (Khanna, 2021), if there is no possibility to ease regulatory burdens by a proportionality assessment, then the AIA might become unsustainable for AI providers or deployers. This would be a severe loss for the EU AI strategy, disincentivising innovation and losing the benefits AI technologies can bring to those values the AIA aims to protect. The May 2023 amendment significantly advanced the regulation by allowing revisions to high-risk system classifications based on an assessment of the risk’s significance, i.e., its probability, severity, intensity, and potential population impact (AIA, Recital 32). However, this revision process is currently without a defined methodology or metrics. Our goal is to furnish support and clarify this mechanism (among others) by introducing a semi-quantitative risk assessment approach.

4 Addressing the Model Flaw: The IPCC Framework for Risk Assessment

The model flaw results from an insufficiently granular risk assessment model: the relevant factors of AI risk are not accurately identified and/or combined.

As argued in Sect. 2, the AIA’s risk model is compatible with the ALARP principle and considers mainly two risk factors (a) the inherent risk of AI technology and (b) a value asset consisting of fundamental principles and rights of the Union. The EU legislator prescribes risk mitigation measures proportionate to the risk magnitude. As a result, risk management measures are allocated according to the four risk categories of the AIA.

Hence, the risk considered in the AIA is legal in nature, expressing the potential detriment that comes from the violation of a legal norm by an AI (i.e., principles and rules) (Mahler, 2007).²⁴ However, the AIA’s risk assessment model does not fulfil the distinctive nature of the legal risk as it does not evaluate comparatively and proportionately the specific weight of legal norms. Quite the opposite, risk assessment in the AIA seems modelled as a neutral tool that treats legal norms as technical standards which are either met or not (Smuha et al., 2021). Consequently, the risk is categorised through a list of AI scopes potentially detrimental to fundamental principles and rights. But risk assessment is not a neutral tool: it reflects the risk appetite of a specific community, weighing the costs and benefits of risk mitigation (Krebs, 2011),

²⁴This is at least one of the meanings that the concept of legal risk can take, and it is the one associated with the Basel Committee on Banking Supervision’s definition: “Legal risk includes, but is not limited to, exposure to fines, penalties, or punitive damages resulting from supervisory actions, as well as private settlements”.

balancing the interests and values of that community, and all this dynamically and diachronically; while promoting a legal value, it may be the case that the unexpected demotion occurs of other equally fundamental legal values. Accordingly, risk management measures should be modulated according to the outcome of such a balancing process. This adaptability is what the AIA needs to incorporate. In fact, despite claiming to be informed by the trade-off between economic development interest and the protection of fundamental rights,²⁵ the AIA seems to predetermine the proportionality judgment that settles the interference between values. Also significantly, not only the list of fundamental rights protected by the proposal is particularly rich, but it also includes inter-related rights,²⁶ making thus difficult a horizontal balance between competing fundamental rights.

The model flaw does not concern only the lack of granularity in the analysis of values and rights. The AIA also lacks an accurate representation of the hazards' sources of AIs, of what makes people vulnerable to these hazards, and of whether hazards and vulnerabilities are mitigated by mechanisms, including legal ones, that already exist (i.e., the net risk) (Black & Baldwin, 2012).

Against this background, the May 2023 compromise text's requirement for deployers of high-risk AI systems to conduct a fundamental rights impact assessment before market introduction is a progressive move. The methodology we propose in this section aims to enhance the accuracy of the proposed assessment outlined by the EU policymaker.

To improve the implementation of the AIA (Simpson et al., 2021), we propose a risk assessment methodology that includes multiple risk factors, and their interferences, and provides a proportionality judgement to review risk categories. This, however, without dismantling or multiplying the draft's tolerance ranges. On the contrary, we suggest applying the four risk categories horizontally to each of the AIs listed in the AIA, so that under varying conditions—e.g., a specific interference among fundamental rights involved—the same system can be treated as unacceptable, high-risk, limited-risk or minimal-risk. This implies that risk categories would not depend by default on AI scopes, but on the real-world risk scenarios associated with the application of AI systems due to the incidence and combination of multiple risk factors.

To build risk scenarios, the Intergovernmental Panel on Climate Change (IPCC) provides a multifaceted risk assessment model, which has then been refined the subsequent literature (Simpson et al., 2021) and which we can use to assess risks of AIs. The risk magnitudes associated with both climate change and AI are influenced by a range of interacting factors, resulting in context-dependent outcomes. Recognizing this, we look to the IPCC model, which offers a detailed and widely

²⁵This is clearly stated in the Explanatory Memorandum of the Proposal: “To achieve those objectives, this proposal presents a balanced and proportionate horizontal regulatory approach to AI that is limited to the minimum necessary requirements to address the risks and problems linked to AI, without unduly constraining or hindering technological development or otherwise disproportionately increasing the cost of placing AI solutions on the market”.

²⁶Cf. European Commission, Explanatory Memorandum to AIA, para 3.5.

recognized framework for assessing the trade-offs inherent in devising risk mitigation strategies.

The IPCC has often conceived the climate change risks—e.g., disaster risk—as the consequence of three determinants: hazard (H), exposure (E), and vulnerability (V).²⁷ Broadly speaking, hazard refers to the sources of potential adverse effects on exposed elements; exposure refers to the inventory of elements within the range of the hazard source; vulnerability refers to the set of attributes or circumstances that makes exposed elements susceptible to adverse effects when they impact the hazard source (Cardona et al., 2012; Liu et al., 2018).²⁸ The IPCC’s approach can be developed further, as in the framework for climate change risk assessment proposed by (Simpson et al., 2021), which evaluates risk at a lower level of abstraction by including the individual components of the risk determinants, i.e., the drivers. Simpson et al. expand the IPCC approach by incorporating a fourth risk determinant: the response (R), which refers to existing measures that counteract or mitigate risk. They also contextualise risk assessment by including multiple types of risk with their own determinants. Thus, according to their framework, the overall risk results from the interaction among (1) determinants, (2) drivers, and (3) risk types (Fig. 2). These three sets of relations occur at stages of increasing complexity. The AIA only considers the lowest complexity stage, where the relevant risk factors are the determinants taken statically, that is, overlooking interactions among their drivers (or with cross-sectorial risk types).

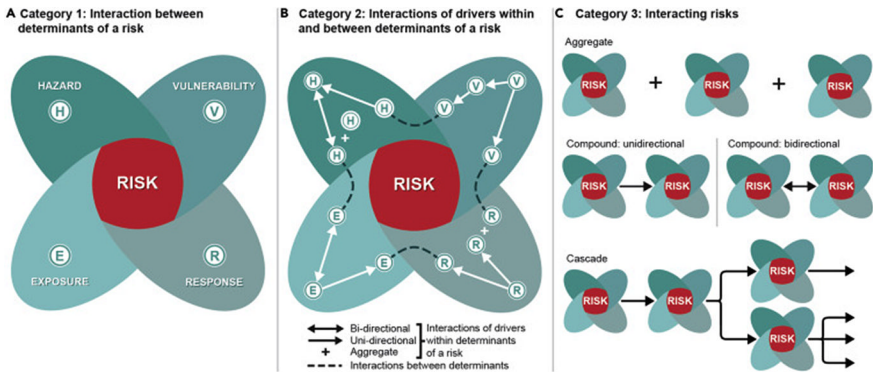


Fig. 2 Three categories of increasingly complex climate change risk by (Simpson et al., 2021)

²⁷This conceptual approach is clearly set out in (Cardona et al., 2012). This approach also emerges in special IPCC reports, e.g., Special Report on Climate Change and Land—IPCC site 2019 <https://www.ipcc.ch/srcccl/> and Special Report on the Ocean and Cryosphere in a Changing Climate 2018 <https://www.ipcc.ch/srocc/>.

²⁸That hazard, exposure, and vulnerability are relevant to risk assessment is also widely believed in the literature other than climate change, such as in (Renn, 2011). In studies on global environmental change and sustainability, the same four determinants were considered as parts of a risk sequence chain (Turner et al., 2003).

The weight of each determinant is given by the drivers and their interactions, both within and across determinants. Interactions among drivers may be (i) aggregate, if drivers emerge independently of each other but jointly influence the overall risk assessment; (ii) compounding, if drivers produce a specific effect on risk assessment when combined, unidirectionally or bi-directionally; (iii) cascading, when drivers trigger others which themselves may produce further drivers in a cascading process. The same applies to interactions between multiple risk types (Simpson et al., 2021).²⁹ Figure 2 below shows the three sets of interactions.

In climate change, the drivers of the hazard (H) can be natural or human-induced events. In AI, these drivers may be either purely technological or caused by human-machine interactions: e.g., the opacity of the model, data biases, interaction with other devices, and mistakes in coding or supervision. The last three hazard drivers interact in an aggregate way. Interactions are compounded when, e.g., low data representativeness compounds with overfitted machine learning models or biased data. The interaction between drivers is cascading when, e.g., model opacity triggers cascading hazards of unpredictability, unmanageability, or threats to security and privacy. An accurate reconstruction of these interactions can provide evidence about the simplicity or complexity of the causal chain between hazard and harm, as well as its likelihood and distribution (Black & Baldwin, 2012).

Drivers of exposure (E) in climate change risk may be people, infrastructure, and other social or economic assets. For AI risk, exposure drivers may be tangible assets, like goods or environment, or intangible assets, like values and rights. As already stressed, the exposed asset of the AIA mainly consists of fundamental rights and values, such as health, safety, employment, asylum, education, justice, and equality. Interactions between drivers of exposure may be aggregated if, e.g., an AI has adverse effects on the right to asylum and the privacy of asylum seekers. It is compounded when, e.g., an AI's adverse effect on the environment compounds with those on health. The interaction between drivers of exposure is cascading when, e.g., an AI's adverse effect threatens access to education, and thus equality and democratic legitimacy (and so on).

Vulnerability (V) drivers of climate change risk may concern the propensity to suffer adverse effects of communities—e.g., poverty—and infrastructure—e.g., lack of flood containment. Drivers of vulnerability in AI risks are multiple and overlapping, e.g., income, education, gender, ethnicity, health status, and age. The lack of appropriate control bodies, procedures, or policies should be included among the drivers of vulnerability for AI risk.³⁰ The AIA shows two conceptions of vulnerability: a generic one, whereby the mere entitlement to fundamental rights entails the propensity to suffer adverse effects of hazards; and a more specific one, whereby all those AIs that “[...] exploits any of the vulnerabilities of a specific group of persons due to their age, physical or mental disability”, (AIA, Article 5) should be banned. In the latter case, the list of vulnerability drivers is rather poor.

²⁹Simpson NP, et al. (n 6).

³⁰Further specifications can be made. For instance, it has been proposed to categorize exposure and vulnerability drivers into four types: ontological, passive, active, and intentional, as in (Liu et al., 2018).

The interaction between vulnerability drivers is aggregated when, e.g., an AIs is deployed in a vulnerable environment, and there are few surveillance or feedback mechanisms. The compounding interaction is perhaps the most interesting one, as an intersectional reading of vulnerabilities can also be advocated in AI risk: ethnicity, gender, health, age, education, economic status, and other characteristics are profiles of vulnerability that have to be considered in the way they intersect and influence each other. The vulnerability stems from various interconnected social processes that lead to multiple dimensions of marginalisation (Kuran et al., 2020). In this sense, the intersectional approach to vulnerability is a risk management principle that enables policy-makers to identify the most appropriate measures to counter hazards to individuals and groups. These interactions make vulnerability a multi-layered condition (Luna, 2019). The interaction between vulnerability drivers is cascading when, e.g., the absence of AIs liability rules triggers several other vulnerabilities for those under the adverse effect of AIs use.³¹

The analysis by Simpson et al. introduces a fourth determinant, i.e., response (R), which concerns existing measures that counteract or mitigate risk. The response indicates the environment's resilience to a specific risk and includes governance mechanisms. Regarding AI risk, the response drivers can be institutional safeguards on the development, design, and deployment of AIs or data quality rules. Consequently, risk assessment and categorisation within the AIA should consider already existing legal measures to avoid the adverse effects of AI technologies, e.g., those contained in the GDPR.³²

Adaptation and mitigation responses may increase or decrease the risk level of specific AIs. As a result, the response determinant can be used to discriminate intrinsic from net risk, the latter adjusted to risk management measures:

[...] where the potential harm is higher than for the intrinsically lower risks, but the probability and/or impact is reduced by risk management and other control measures, or by systems of resilience – such as capital requirements in financial institutions, or engineered safety controls in power stations, or by the possibility of remediation (Black & Baldwin, 2012, 5).

Simpson et al. also introduce a third stage of interaction, between climate change risk and other types of risk, which are extrinsic to it and have their own determinants. Risk types that interact with AI risk may be, e.g., market, liability, and infrastructure risks. Some of these risk types are created by the AI risk itself—i.e., cascading interactions—others are independent but may affect the overall assessment of AI risk—i.e., aggregate or compounded interactions. For instance, an

³¹In the proposed fundamental rights impact assessment, interest in vulnerability is emphasized: “This assessment should include [...] (f) specific risks of harm likely to impact marginalised persons or vulnerable groups” (Article 29a, (f)).

³²Consider, for example, art. 35 on data protection impact assessment: “1. Where a type of processing in particular using new technologies, and taking into account the nature, scope, context and purposes of the processing, is likely to result in a high risk to the rights and freedoms of natural persons, the controller shall, prior to the processing, carry out an assessment of the impact of the envisaged processing operations on the protection of personal data.”

aggregate interaction occurs between AI risk and policy risk, in the sense that adverse effects of ineffective policies or regulations—perhaps external to AI—cumulate with the adverse effects of AIs' deployment. In a healthcare setting, for example, an AI system used for skin cancer diagnoses might inaccurately diagnose patients due to algorithmic biases (Gupta et al., 2016). Concurrently, outdated tort liability laws may not adequately address AI's role in healthcare. This leads to compounded risks: the AI's misdiagnoses are exacerbated by unclear liabilities, heightening overall risk magnitude in patient care. AI risk can then compound with the risk of the digital infrastructure in which an AIs operates.

Finally, AI risk can cascade into multiple other types of risk, the risk to innovation, to digital sovereignty, to economic sustainability, to power concentration, and so forth.

This third stage of interaction should be linked to that of ancillary risks, i.e., risks posed or increased by the risk regulation itself. For example, banning AIs should be justified also against the loss of opportunity benefit of their use, the potential barriers to technological innovation that the ban raises, and the risk posed by the systems replacing the banned ones (Sunstein, 2004). The AIA's regulatory choices cannot be justified just by their positive impact on the intended scope—i.e. the protection of fundamental rights—but also by the (difference between) the marginal gains and harms they generate for other values at stake (Karliuk, 2022).

To sum up, the risk magnitude of each AIs listed in the AIA should be assessed in terms of the interactions among determinants, drivers, and other risk types. Although AIA considers some interactions among determinants—e.g., the scale and the likelihood of adverse effects on values—it does not account for the interaction among the individual drivers of those determinants, nor does it evaluate the risk of AIs in relation to other types of risk. Therefore, the AIA misestimates the AI risk magnitude and anchors risk categories to static, coarse-grained factors.

Once the determinants, drivers and external types of risk are identified, adaptation and mitigation become easier, i.e., to reduce the risk of AI by planning actions (including policies) that address the factors of hazards, exposure, and vulnerability (Simpson et al., 2021).

The granular risk assessment we propose has a higher degree of variability. The risk categories of the AIA become risk scenarios (Renn, 2011), which change depending on the interactions among risk factors. This leads to a more accurate representation of the risk magnitude—i.e., the likelihood of detriment and severity of consequences on values—with connections among risk factors being made explicit. Even if the EU legislator intends to keep the current framework—where risk categories are pre-determined based on the AI's scopes—this model can aid in the proposed additional assessments that could revise the risk categorisation (i.e., risk significance). However, what we have presented in this section, is just a general framework. While risk magnitudes may correspond in the abstract to risk categories, as a preliminary evaluation, this assignment also must pass the proportionality test that we shall describe in the next section.

5 A Quantitative Basis for the Model: The Proportionality Test

Though not directly mentioned in the AIA, an issue shared with the ALARP principle is setting risk management measures, without defining what qualifies as a “grossly disproportionate” containment measure.

A way to offer quantitative support for ALARP-based legislative choices is through the traditional cost-benefit analysis (CBA) (French et al., 2005). While the ALARP allows the costs of risk mitigation to exceed the benefits as long as they are not exorbitant, the CBA specifies that intervention is justified only if costs are less than or equal to the benefits. CBA does not account for uncertain costs and benefits (Jones-Lee & Aven, 2011). Despite this drawback, CBA can support ALARP as a preliminary informational input: as far as possible, CBA quantifies known costs and benefits so that this information can be combined with a qualitative assessment of what is “reasonably practicable” (Ale et al., 2015). The risk assessment model presented in the previous section helps us to combine CBA, the ALARP principle, and the AIA to account for the likelihood and distribution of adverse effects, the causal chain between hazards and harms, the effects of AI risk regulation (i.e., ancillary risks), and alternative measures for risk mitigation. However, CBA remains an imperfect tool for the AIA, as the former expresses the value of things with a single numerical parameter, usually market prices, while the latter concerns a legal risk, whose exposed asset consists of fundamental rights and values, which, respectively, are intended to represent “principles of [EU] law of a constitutional nature”³³ and the “very identity” of the EU legal order.³⁴

However, we suggest an alternative (semi-)quantitative approach to ascertain when sacrifices to mitigate risk are “grossly disproportionate” (within the scope of the AIA). This quantitative assessment should be seen as complementary—a second step—to the risk assessment model of the previous section: to assign the appropriate risk category for a specific scenario, we need to compare the impact each risk category has on the assets served by the intended scope of the AIs (P_x)—e.g., law enforcement—against those of the exposed asset (P_y)—e.g., safety, health, and equality. Thus, if applying the high-risk category to AIs for law enforcement under a specific risk scenario has a sub-optimal impact on the joint realisation of principles and rights, it is desirable to opt for an alternative risk category. Whereas, if the marginal gains to law enforcement outweigh the marginal harms to other rights, then the risk category is justified.

Robert Alexy proposed a well-known method in legal theory to quantify this type of choices (Alexy, 2002). According to this approach, a legal norm that interferes with fundamental values³⁵ is legitimate when it meets a proportionality test characterised by the following optimisation principles:

³³Joined cases C-402/05 P and C-415/05 P *Kadi and Al Barakat International Foundation* EU: C:2008:461, para 276.

³⁴Case C-156/21 *Hungary v European Parliament and Council of the European Union* EU:C:2022:97, para 127

³⁵In Alexy’s theory, these fundamental values are typically constitutional principles.

- *Suitability*, which “excludes the adoption of means obstructing the realisation of at least one principle without promoting any principle or goal for which they were adopted” (Alexy, 2003). In the AIA, the legislative choice of assigning a risk category R_1 to an AI that negatively impacts one principle P_2 is suitable if it impacts positively another principle P_1 .
- *Necessity*, which “requires that of two means promoting P_1 that are, broadly speaking, equally suitable, the one that interferes less intensively in P_2 ought to be chosen” (Alexy, 2003). In other words, R_1 with a negative impact on P_2 is necessary if it has a positive impact on P_1 and there is no alternative, R_2 , having a higher positive impact on P_2 and non-inferior on P_1 (Sartor, 2018). In the AIA, as in many other cases, Pareto-optimality equilibria are rather unstable: multiple values are involved, and a principle P_3 that is negatively interfered with by R_1 can easily occur. These unavoidable costs call, according to Alexy, for a third principle.
- *Proportionality in the narrow sense*, which states that “The greater the degree of non-satisfaction of, or detriment to, one principle, the greater the importance of satisfying the other” (Alexy, 2002, 102). This principle provides a basis for determining whether or not the importance of satisfying P_1 with R_1 justifies the impairment or failure to satisfy P_2 . When multiple values are involved, as in the AIA, we will say that R_1 with a negative impact on P_2 is balanced if there is no alternative R_2 having a lower negative impact on P_2 and a higher overall utility on $P_3, P_4 \dots P_n$. (Sartor, 2018).

Such a proportionality test, which is (by and large) in line with the proportionality test the EU Court of Justice applies while balancing competing rights and values (Alexy, 2003; Tridimas, 2018), may support legislative choices and trade-offs within the AIA, i.e., the exposed asset of AI risk.³⁶ We suggest that it may serve to justify trade-offs between fundamental values/rights that (should) inform the risk categorisation of AIs. The outcome of the test may warrant the ascription of a risk category R_1 (e.g., high-risk) to specific AIs or shifting an AI to a new category R_2 (e.g., minimal risk). For this purpose, proportionality in the narrow sense should be broken down into three evaluations:

- (1) the intensity of interference (I_x), the degree of non-satisfaction or detriment to a principle P_x to the benefit of a competing one P_y
- (2) the concrete importance (C_y) of satisfying P_y
- (3) the concrete weight of P_x ($W_{x,y}$), namely the ratio between I_x and C_y , which determines whether the importance of satisfying P_y justifies the non-satisfaction or detriment to P_x (Alexy, 2003).

Finally, the abstract weights of P_x (W_x) and P_y (W_y) also play a role in the overall balance.³⁷

³⁶Other formal models have been developed to balance the impact of actions and decisions on values in the legal domain. For instance, consider the model developed by (Maranhão et al., 2021).

³⁷Alexy also includes another variable in his weight formula, namely the epistemic reliability of the balancing premises. For simplicity of exposition, we will not consider them here.

In some cases, P_x will prevail over P_y , e.g., when I_x is severe, and C_y is weak. In other cases, P_y will prevail over P_x . There may also be cases where there is no prevalence between P_x and P_y , $I_x = C_y$, which creates deadlocks, increasing discretion in balancing. The outcome of the ratio between the intensity of the interference on a specific principle and the concrete importance of the competing one is expressed by the following, simplified version, of the weight formula (Alexy, 2003):

$$W_{x,y} = \frac{I_x \cdot W_x}{C_y \cdot W_y}$$

Applying the weight formula to the AIA, I_x would correspond to the degree of interference a risk category, with its containment measures, has on a (set of) value(s) served by the intended scope of AIs: e.g., the interference to public safety (P_x) as served by biometric categorisation systems. C_y would correspond to the concrete importance of satisfying a competing (set of) value(s) explicitly protected by the AIA, which is part of risk-exposed asset in biometric categorisation systems: e.g., the right to privacy (P_y). The concrete importance expressed in C_y depends on qualitative assessments in relation to the risk scenario, i.e., what are the hazard factors, vulnerability profiles and response mechanisms that determine the magnitude of risk in the concrete scenario (as described in the framework shown in the previous section). Therefore, whether the EU legislator is authorised to restrict the use of AIs for biometric categorisation will depend on whether the magnitude of the privacy risk posed by these systems (C_y) justify the impairment of public safety caused by the measures of the relevant risk category (I_x).

Although the weight formula relies on non-numerical premises—like judgments about the degree of interference of a risk category or the abstract weight of principles (Alexy, 2003)—numerical values can still be assigned to I_x and C_y . This can be done using a geometric sequence, like $2^0, 2^1, 2^2, 2^4$, to assign numerical ranges to the AIA's four risk categories according to the degree of interference, or non-satisfaction, they cause to the intended scope of an AIs (I_x): unacceptable risk = 16, high-risk = 4, limited risk = 2, minimal risk = 1. The same numerical ranges may be assigned to the importance of satisfying the competing principle—(C_y): major = 16, severe = 4, moderate = 2, light = 1³⁸—and to the abstract weights of principles (W_x and W_y). As shown below, where the asset served by the intended scope of an AIs prevails over the exposed asset, the concrete weight $W_{x,y}$ will be greater than 1. Conversely, $W_{x,y}$ will be less than 1.

- (A) $I_x \cdot W_x (16 \cdot 4)/C_y \cdot W_y (8 \cdot 2) = 4$
 (B) $I_x \cdot W_x (4 \cdot 4)/C_y \cdot W_y (8 \cdot 16) = 1/8$

This quotient describes the concrete weight of the asset served by the intended scope of an AIs given the interference of a risk category on it (I_x) and a competing asset protected for being partly exposed to the AIs (C_y). The inclusion of the

³⁸This requires assuming that the abstract weights have the same impact on the concrete weight as the intensity of interference.

vulnerability and response determinants' values in the ratio can make the proportionality test fully aligned with the risk assessment model outlined in Sect. 4.

To sum up, the quotient of the weight formula is a quantitative criterion to assess whether the risk control measures are “grossly disproportionate” in the AIA, given the balance of relevant values, and therefore whether a risk category is suitable for the risk scenario of an AIs or whether it should be changed. In particular, what is grossly disproportionate can be quantified over a range. In our example, according to the numerical parameters we employed, it is reasonable to argue that the quotient of the weight formula should not be less than 1 or greater than 4. If it falls outside this range, then the balancing between principles is disproportionate and it is advisable to alter the risk category. Indeed, out of the range, a specific risk category may be inadequate for the risk scenario, with measures too stringent or too soft to balance competing EU values, like privacy and technological innovation. In this way, the AIA fails to achieve one of its main objectives: a uniform protection of EU fundamental rights.

We are aware that compulsory numerical values of EU principles and fundamental rights cannot be pre-assigned. Also, attempts to establish a strict hierarchy among EU fundamental values and rights have so far failed.³⁹ While acknowledging the importance of these circumstances, we believe that a quantitative method for assessing risk containment measures could help relevant actors make policy decisions and avoid significant imbalances when implementing the AIA. Numerical values have been assigned to the coefficients in the proportionality test to enhance clarity, but these coefficients can also be compared through non-numerical preferences or magnitudes, such as the Paretian superiority illustrated in (Sartor, 2018).

On a different note, we cannot ignore the role that EU institutions—and, in particular, the role that the EU Court of Justice—shall play in preserving the constitutional framework of the Union and the untouchable core of the EU legal order, which include its fundamental values and rights. This is why in the next section, we shall discuss the allocation of competences and roles in scenario building and proportionality assessments.

6 Using the Semi-Quantitative Risk Approach: Three Applications Under the AIA

We illustrate three potential uses of this semi-quantitative risk framework under the AIA. These applications are not mutually exclusive; they can be implemented simultaneously. The first is intended for policymakers, while the second and third are aimed at deployers of high-risk AI systems.

(1) When implementing the AIA. This implies transitioning from a scope-oriented categorisation of risk to a scenario-based model that considers the interplay

³⁹This incapacity can stem from non-commensurability of values involved. While our model assumes commensurability, Amartya Sen has demonstrated that even in the presence of non-commensurability, rational decisions can still be made using a systematic approach. This approach involves acknowledging the limitations of our value judgments and making reasoned choices based on available information and analysis. Sen proposes the ‘weighted average principle’ as a method for making decisions when direct comparisons are not possible, considering the relative differences between options (Sen, 2004, 50).

of multiple factors in specific situations (as in Sect. 4). The four risk categories should then be applied horizontally to AIs so that, under varying risk scenarios, the same system can be estimated as unacceptable, high-risk, limited-risk or minimal-risk.

However, this application would pose some practical issues. In fact, although the categorisation of risk in the AIA is coarse-grained, its strategy of connecting risk measures to broad scopes of AIs makes it easier to approve and monitor them. Indeed, as previously highlighted in the context of regulation by design in the GDPR (Almada et al., 2023; Michelakaki & Vale, 2023), a finely-grained approach may fall short in offering sufficient guidance to regulated actors. In contrast, a legal framework with risk scenarios built on interacting factors and tested by proportionality-based balancing, as the semiquantitative model we are presenting, might complicate the procedures laid down in the AIA.

This issue is still manageable: under the existing AIA framework, national supervisory authorities (AIA, Title VI) could undertake the task of constructing risk scenarios. However, this approach would alter the governance structure of the AIA, which currently operates predominantly at a supranational level for regulating high-risk AIs. For this reason, it would be crucial to determine the competences, functions and interactions of supranational institutions and national bodies in the risk assessment of AIs. Considering the shared nature of the competences exercised by the EU legislator to adopt the AIA,⁴⁰ it remains undisputed that the EU legislator should retain a primary role in shaping the risk-assessment model at stake.⁴¹ Meanwhile the European Commission should keep its role of guardian of the AIA enforcement and the EU Court of Justice's authority in judging whether the risk assessment is consistent with the essential core of EU fundamental values (Lenaerts, 2019). This is particularly important considering the systematic backsliding on fundamental values and rights taking place in some EU countries.

More to the point, under our semi-quantitative model, the implementation acts of the EU Commission might establish (a) the key drivers of the four risk determinants, (b) the extrinsic types of risk to account for and (c) the (abstract) weight of the principles involved in the proportionality test. This task may be done by the AI Office, new body set up by the AI Act within the European Commission. This entity is designed to ensure a harmonized application of the AI Act through standardisation, provide guidance, and coordinate joint cross-border investigations. Anyway, these factors could be linked precisely to the scopes already identified in the AIA through the broader risk categorisation (e.g., Annex III).⁴² In this way, the broad scopes of AIs would still play a primary role in risk regulation—which means that the text of the AIA would not require substantial changes—and EU institutions would limit the discretion of Member States. National authorities would be

⁴⁰As it is well-known the AIA proposal is based in the first place on Article 114 TFEU, providing a EU shared competence in adopting measures to ensure the establishment and proper functioning of the internal market. In addition, the proposal is based on Article 16 TFEU, due to its connection to the processing of personal data.

⁴¹On the fundamental role the EU legislator should play in this respect, see (Fontanelli, 2016).

⁴²In the next section, we offer a case study that illustrates how key drivers of the four determinants and extrinsic risk types may be identified in connection to the scope of an AIs, i.e., justice (Sect. 7).

responsible for assessing risk in particular cases—thereby enhancing their powers over what is in the AIA—through scenarios and proportionality tests.⁴³

Detractors could claim that the proposed solution may lead to a partially diversified enforcement of the AIA within the EU, conflicting with the EU’s uniform values and rights. However, this position overlooks the necessity of context-sensitive risk assessment. Moreover, our method does not necessarily weaken the effectiveness of the EU’s fundamental values and rights as it is based on the idea of introducing a robust rational procedure, under the strict supervision of the European Commission and the ultimate control exercised by the EU Court of Justice.

To conclude, while some would prefer centralized risk assessment, our proposal to recalibrate the regulation during implementation seems more feasible. This means adhering to the AIA’s existing risk categories but applying them with a focus on the interplay of risk factors and proportionality. At the same time, among the three uses of the scenario-model we are discussing, this one appears to be the least practicable. This is due to the characteristics of the normative environment surrounding the AIA—specifically, the product safety framework (Almada & Petit, 2023)—and the legislative stage at which the AIA currently stands. The next two applications we shall explore are less constrained by path dependency and the AIA’s normative integrity. However, they also are finely-grained and this generates challenges that we shall address at the end of this section.

(2) When Evaluating the Significance of Risk (AIA, Recital 32a). The assessment of risk *significance* was initially introduced by the EP compromise text from 11 May 2023, and later confirmed by the consolidated version dated January 26, 2024, as per the COREPER text. According to the proposal, AI systems to be classified as high-risk must also pose what is called a ‘significant risk’, requiring evaluation of the risk’s severity, intensity, likelihood, duration, and potential targets, whether an individual, multiple people, or a specific group (e.g., AIA, art. 3 (1b)).⁴⁴ If a deployer can demonstrate that her system does not pose a significant risk, contrary to the initial categorisation based on broad AI scopes, she can then reclassify the risk level of their system. The compromise text mandates these evaluations to be conducted by deployers but with the obligation to inform national supervisory authorities, relevant stakeholders, and representative groups of individuals who may be impacted by the application of the (high-risk) AI system (e.g., Recital 32 and Article 29a). We believe this option could be extended to systems initially categorized as posing unacceptable risks.

However, the proposal does not offer a methodology for integrating these metrics. Our model addresses this, offering a way to calculate risk significance. This approach enables a flexible review procedure, allowing deployers to reduce regulatory burdens when they can provide evidence of lower risk.

(3) When Implementing Deployers’ Internal Risk Management (AIA, Art. 9). Deployers of high-risk AI systems, as mandated by Article 9 of the Artificial

⁴³Or any other type of proportionality-based balancing.

⁴⁴A second update from the compromise text mandates deployers of high-risk systems to conduct a fundamental rights impact assessment and develop a risk mitigation plan in coordination with the national supervisory authority and relevant stakeholders before market entry (e.g., AIA, Recital 58a and Article 29a).

Intelligence Act (AIA), must establish an internal risk management system. This system is aimed at ensuring compliance with the AIA and the reliability of high-risk AI systems. Article 9 requires deployers to identify both known and foreseeable risks associated with their AI systems, considering their intended uses and potential misuses. This includes the assessment of additional risks that may emerge, informed by data from post-market monitoring. Deployers are required to implement effective risk management measures, such as risk elimination or reduction, through thoughtful design and development processes. The risk management systems must be thoroughly documented and regularly updated throughout the AI system's lifecycle.

The AIA's outline of this risk management system, while not highly detailed, aligns with international standards (such as ISO)⁴⁵ and the historical evolution of the concept (Dionne, 2013). However, the AIA does not provide a shared procedure for these in-house risk management systems. Likely, further details will be provided in the Commission's delegated and implementing acts.

In this context, our scenario-based model offers a common framework, a detailed terminology, granular risk factors, and a dynamic methodology for integrating them in real-life situations. And the internal risk management system may benefit from it. The scenario-based model may help identifying "known risks" and assists in evaluating risks as more or less "foreseeable" (as per Article 9, letter (b)). This evaluation depends on factors like response mechanisms (e.g., prevention measures), hazard drivers and the way they combine (e.g., aggregating or cascading). Our model also includes extrinsic and ancillary risks and their impact on the main risk magnitude, as outlined in Sect. 4, which may help assessing "other possibly arising risks based on the analysis of data gathered from the post-market monitoring system" (Article 9, letter c). This approach enhances the identification of "suitable risk management measures" as required by Article 9, letter d, of the AIA. It clarifies the effectiveness of current response mechanisms in mitigating vulnerability drivers and explains how hazard drivers impact these vulnerabilities. This insight is crucial for understanding the overall risk magnitude, allowing for a more targeted and effective refinement of risk management strategies within the AI system. Furthermore, the scenario analysis may also be probabilistically qualified through the risk matrix that we shall illustrate in Sect. 7.

Finally, the proportional testing of rights and interests, discussed in Sect. 5, plays a significant role in various aspects of the internal risk management system. This includes identifying mitigating measures and testing procedures (art. 9 point 2 and point 6) that are both appropriate and proportional to the community's risk appetite (seen as the ratio of values and the interests involved).

Unrelated to risk management yet somewhat linked to the obligations of high-risk systems' deployers, proportional testing could also be relevant for the fundamental rights impact assessment (FRIA) recently introduced in the draft.⁴⁶

There is a common aspect to all three applications that needs consideration. Granular risk categorisation based on scenarios which combines multiple risk factors,

⁴⁵ISO 31000:2018 Risk Management—Guidelines, Clause 3.2: "coordinated activities to direct and control an organization with regard to risk"

⁴⁶The FRIA was introduced in the European Parliament proposal through Article 29a.

as we suggest, might be supererogatory and complicate the AIA procedures. To figure out the severity of this, we need to distinguish short-term from long-term aspects.

In the short run, a scenario-based risk assessments may indeed deter AI deployment and investment. To mitigate this, different strategies might be recommended to make our proposal more sustainable. Firstly, European legislation might indicate, in the AIA's implementing acts, the key risk drivers for each broad AI scopes already outlined in the regulation (e.g., in the Annex III). This would ease the task of deployers and minimizes arbitrariness in the AIA's enforcement. We shall illustrate this in the next section. Second, automating risk identification and management can streamline processes. Finally, a phased, iterative approach starting with a granular risk assessment only for a few deployers—maybe with lower risky systems and then with lower compliance costs—might enable procedural refinement and prepares others for a smoother implementation.

In the long run, these short-term costs will be offset by the benefits of decreased compliance costs as contextually tailored risk assessments yield less over-inclusive risk categories and more effective risk prevention or mitigation measures.

7 An Illustration Case for Showing the Contributions of the Semi-Quantitative Approach

Our analysis offers two contributions to the enforcement of the AIA and the regulation of general-purpose AI (GPAI).

First, the risk assessment model shown in Sect. 4, supported by the quantitative proportionality test shown in Sect. 5, improves the enforcement of the AIA. This improvement applies whether it is used for the horizontal implementation of risk categories, evaluating risk significance, or within deployers' internal risk management systems. It would provide risk management measures that are more appropriate to estimate and contain the dangers of AI, more specific for national regulators (and judges), more sustainable for AI providers, and ultimately more likely to achieve the AIA's goals of protecting all fundamental EU values involved. Ideally, such granular risk management measures can help avert, or more effectively handle, issues related to the under-inclusiveness or over-inclusiveness of risk categories (Hacker, 2023).

To show how risk scenario building might work in the AIA, let us consider a case study, that of AIs used to assess the recidivism rate of natural persons in criminal trials. The semiquantitative approach consists of two stages: risk-scenario building and the proportionality test. The risk drivers here identified can be easily inferred from the AIA. Of course, applying our proposed assessment model during the AIA implementation stage would necessitate enhanced legislative transparency in setting the drivers and extrinsic risks.

Starting from the risk-scenario building, the four determinants of AI risk, the interaction among their drivers and with other risk types may be thus combined:

- (a) Hazards. These drivers of an AI for recidivism rate assessment may be the inner opacity of the system, its reliability, the poor quality or misuse of the training data (e.g., outdated or incomplete data sets), and its validity. Validity is critical to ensure that the instrument measures what it is intended to measure

(Quattrococo, 2020). When these hazard drivers compound, they can lead to the AIs perpetrating discrimination biases. For instance, the AI might consistently and incorrectly predict higher recidivism rates for certain demographic groups due to poor training data, compounded by its opacity and unreliability. This risk magnitude differs from a scenario where these drivers do not compound, such as when the decision-making process is transparent, allowing stakeholders to identify and rectify biased data issues more readily. The greater these hazard drivers are, and the more likely they combine to produce such wrongdoing, the “heavier” the hazard determinant will be in the specific risk scenario. What is more, the hazards must be connected to the vulnerability drivers of a specific environment in which AIs are deployed, not least because these will be inclined to replicate the social discriminations of the environment.

- (b) **Exposure.** These driving factors include the directly impacted parties, such as defendants, and, indirectly, the fundamental values potentially affected by the use of an AI to assess the recidivism rate. This would involve the interest to an efficient trial together with some substantive legal principles—e.g., the principle of criminal culpability and of equality—or procedural ones—e.g., the principle of transparency and the right of/to defence (Garrett & Monahan, 2020). These drivers also interact with each other, and where they interfere, it is necessary to balance them to assess the overall weight of the exposure determinant. This also requires balancing those values that the use of AIs is intended to enhance (consistent with the proportionality test in Sect. 5), such as the principle of predictability, legal certainty, safety, and efficiency. An example of interference occurs when a recidivism assessment AI system is designed for procedural fairness and equality by only using unbiased data. This approach neglects defendants’ personal narratives and circumstances, clashing with the interest for an individualized justice. This focus on depersonalized data can overlook important personal factors like rehabilitation efforts or life challenges, crucial for a fair assessment.
- (c) **Vulnerability.** These drivers may be the attributes that make individuals or groups susceptible to the adverse effects of automatic recidivism rate assessment, such as ethnicity, economic conditions, and education. So, for instance, the lack of data literacy among judges, lawyers, and defendants may compromise their ability to grasp the workings of AI systems, their limitations, and the implications of their use in legal settings. This lack of data literacy and awareness can lead to over-reliance on or misinterpretation of AI-generated recidivism predictions. When these drivers interact with each other, perhaps compounding or cascading, vulnerability should be treated as a multi-layered condition (Luna, 2019): e.g., the compound of ethnicity and socio-economic conditions often leads to a heightened sensitivity to the biases of prediction systems. As mentioned above, drivers of vulnerability compound with hazard drivers: e.g., biases in the recidivism rate assessment will be greater where social discriminations are already in place.
- (d) **Response.** These are pre-existing strategies designed to mitigate the risks associated with automated recidivism rate assessments. They might be governance measures, like standards for data quality and data collection, transparency, bias examination, and human oversight. An example includes regulations mandating the exclusion of specific indicators that, while predicting some degree of social

dangerousness, are directly or indirectly linked to ethnic or social background, e.g., the postal code (van Dijk, 2022).

- (e) Extrinsic risks. The risk of AIs for recidivism rate assessment would finally interact with extrinsic risk types. Some extrinsic risks, in this case, would be compliance risk, liability risk, and economic risk. Indeed, AI risk may be influenced by the lack of effective rules for the allocation of liabilities for adverse effects. If the system's prediction leads to an unjustly prolonged detention, the developers or users of the AI could face lawsuits. This scenario could then lead to or exacerbate economic risks within the AI market, as potential legal repercussions might deter investment. The overall risk magnitude should also consider ancillary risks. For example, risks arising from inefficient regulation or poor risk management can impact innovation, lead to lost opportunities, and affect digital sovereignty. The introduction of regulatory burdens, or entry barriers, on AIs' providers may weaken technological innovation and, in the case of a radical ban, resulting in the loss of opportunity for the general social interest.

The interactions among these risk factors determine the two input variables of the overall risk magnitude of the specific scenario: (1) the likelihood of the event depend on the interaction between hazard drivers and response drivers (e.g., preventive measures); (2) likewise, the severity of the detriment can be higher or lower depending on the hazard sources, exposed asset, and vulnerability profiles.⁴⁷ As a result, risk magnitude is associated with the four risk categories of the AIA—i.e., unacceptable (U), high (H), limited (L) and minimal (M) risk—as illustrated in the risk matrix below⁴⁸ (Table 1):

The five levels of severity are described qualitatively, and those of likelihood in percentages in a range between 0 and 1 (where 0.20 – 1 is remote risk, while 0.80 – 1 the risk is almost certain). Under this matrix, the intersection of the input variables correlates with one of the four risk categories of the AIA.⁴⁹

Table 1 Risk matrix inspired by (Ni et al., 2010)

<i>Severity</i>	Major	L	H	U	U	U
	Serious	M	H	H	U	U
	Moderate	M	L	H	H	H
	Light	M	M	L	H	H
	Negligible	M	M	M	L	L
		0–0.20	0.20–0.40	0.40–0.60	0.60–0.80	0.80–1
	Likelihood (%)					

⁴⁷These are the same input variables of the conception of risk magnitude embraced by the AIA (e.g., Title III, art. 7).

⁴⁸The risk matrix approach is widespread in semi-quantitative risk assessments, such as the one we are suggesting (Ni et al., 2010).

⁴⁹For example, someone else might think it more correct that a moderate detriment with a probability between 0.20 and 0.40 percent should correspond to a high-risk category.

The second step is to evaluate the suitability of the resulting risk category in relation to the asset exposed to the use of an AIs, by means of the proportionality test. Let us assume that the risk magnitude for a specific recidivism rate assessment system matches its current categorisation in the AIA, i.e., unacceptable risk (U). One of the principles served by AIs for assessing recidivism rates is safety (P_x) and, according to the geometric sequence seen in Sect. 5, its abstract weight can be quantified with a score of 4 (W_x). The degree of interference (I_x) of the AIA's high-risk category on legal certainty is 16. In the denominator of the Weight Formula, the abstract weight of a competing principle, e.g., criminal culpability (P_y), might be 4 (W_y) as well as the concrete importance of satisfying it (C_y). Applying all these values to the ratio— $W_{x,y} = (I_x \cdot W_x)/(C_y \cdot W_y)$ —the outcome would be 4, which falls within the proportionality range we have assumed. As a result, we might conclude that the risk category is appropriate as it correctly balances the values involved. Of course, if the competing principle was deemed to be less significant, for instance, it held a light value such as 2, then the outcome of the equation might not fall within the range and the risk category should be revised.

The second contribution of the risk assessment model presented here concerns one of the regulatory issues that emerged from the debate on the AIA: the governance of general-purpose AI (GPAI), which is the label that includes the so-called Generative AI (e.g., large language models like ChatGPT). The issue was raised in an amendment proposing a definition of GPAI and classifying them as high-risk systems.⁵⁰ In the consolidated version, they have an autonomous three-tier risk classification. Anyway, GPAIs are systems that can be deployed in multiple fields and with different tasks, some of which were unintended by the developers). This definition would also include open-source AI models (e.g., open-source datasets).

Indeed, if the intended purposes are not foreseeable, due to their scalability and vast flexibility, neither are the fundamental values that AIs would affect and based on which their risk would be categorised. This implies that the application of the AIA would be even more static than for AIs with intended purposes. Therefore, the construction of risk scenarios based on determinants, drivers and types seem the only way to categorise and regulate GPAI in a granular manner and avoid treating them all the same. Given these AI technologies' success on the market, undifferentiated regulatory treatment might negatively impact AI industry innovation.

The semi-quantitative model outlined in this article would facilitate risk assessment and categorisation for all those situations that the AIA leaves uncovered, for example, where it recognizes the discretion of the European Commission in updating or modifying the list (and to remove use-cases) of high-risk AIs provided that:

[...] (b) the AI systems pose a significant risk of harm to health and safety, or an adverse impact on fundamental rights, to the environment, or to democracy

⁵⁰GPAIs were excluded from the previous draft of the AIA. However, they are given more room in the compromise text, as implementations of foundation models, and are no longer equate with high-risk systems. They must still adhere to certain documentation and transparency rules. For instance, generative foundation models must always disclose that the content was AI-generated (AIA, Recital 60g).

and the rule of law, and that risk is, in respect of its severity and probability of occurrence, equivalent to or greater than the risk of harm or of adverse impact posed by the high-risk AI systems already referred to in Annex III. (AIA, Article 7, point 1 (b))

However, to determine whether the risk associated with new AI systems matches or exceeds those already classified as high-risk, thus justifying their addition to the high-risk list, a robust and transparent risk assessment methodology is necessary. This is something the AIA does not currently offer, but an attempt to provide such a methodology has been made in this paper.

8 Conclusions

In this paper, we have offered a semi-quantitative approach to AI risk assessment, articulated in two stages: (1) the construction of risk scenarios and (2) a proportionality-based quantitative assessment. For scenario construction, we have referred to the IPCC's theoretical framework and the literature on climate change risk. Accordingly, risk results from the interaction among four determinants, among individual drivers of determinants, and among extrinsic types of risk. For the second stage, we have referred to the quantitative approach developed by Alexy for balancing legal principles. Such a quantitative assessment aims to check whether the risk category assigned following the scenario construction is proportionate to the values involved in employing AIs.

We have discussed three applications of a semi-quantitative risk framework under the AI Act (AIA). Firstly, the framework supports the horizontal implementation of the AIA by introducing a scenario-based risk model. It would categorize AI systems variably as unacceptable, high-risk, limited-risk, or minimal-risk, depending on specific situations. Secondly, it aids in the evaluation of risk significance, providing deployers with the opportunity to reassess the risk levels of their AI systems, which could lead to reduced regulatory burdens. Thirdly, it facilitates and harmonizes the implementation of the internal risk management system as outlined in Article 9 of the AIA, by offering a comprehensive framework to help deployers in identifying, evaluating, and managing risks. Our findings suggest that the latter two applications—evaluating risk significance and implementing internal risk management—are more practical and less challenging to realize, considering the AIA's framework and legislative process.

Finally, we have pointed out that a semi-quantitative approach can improve the enforcement of the AIA and help address issues uncovered by the EU regulation, e.g., risk assessment for the GPAI, without undermining the protection of EU fundamental values and rights. Future research should investigate further governance issues, including identifying which institutional bodies are called upon to apply the semi-quantitative risk analysis, with what specific faculties and with how much discretion in evaluating risk factors.

Author Contributions Not applicable.

Funding Open access funding provided by Alma Mater Studiorum – Università di Bologna within the CRUI-CARE Agreement.

Data Availability Not applicable.

Declarations

Informed Consent Not applicable.

Competing Interests Luciano Floridi is part of the Editorial Board of Digital Society.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- Abrahamsen, E. B., Abrahamsen, H. B., Milazzo, M. F., & Selvik, J. T. (2018). Using the ALARP principle for safety management in the energy production sector of chemical industry. *Reliability Engineering & System Safety*, 169(January), 160–165. <https://doi.org/10.1016/j.res.2017.08.014>
- Ale, B. J. M., Hartford, D. N. D., & Slater, D. (2015). ALARP and CBA All in the Same Game. *Safety Science*, 76(July), 90–100. <https://doi.org/10.1016/j.ssci.2015.02.012>
- Alexy, R. (2002). *A theory of constitutional rights*. Oxford University Press.
- Alexy, R. (2003). On balancing and subsumption. A structural comparison. *Ratio Juris*, 16(4), 433–449. <https://doi.org/10.1046/j.0952-1917.2003.00244.x>
- Almada, M., Maranhão, J., & Sartor, G. (2023). *Art. 25. Data protection by design and by default*. Nomos, Beck, and Hart Publishing. <https://cadmus.eui.eu/handle/1814/75913>
- Almada, M., & Petit, N. (2023). The EU AI act: A medley of product safety and fundamental rights? Working Paper. European University Institute. <https://cadmus.eui.eu/handle/1814/75982>.
- Bai, Y., & Jin, W.-L. (2016). Chapter 38 - Risk assessment methodology. In Y. Bai & W.-L. Jin (Eds.), *Marine Structural Design* (2nd ed., pp. 709–723). Butterworth-Heinemann. <https://doi.org/10.1016/B978-0-08-099997-5.00038-1>.
- Baumol, W. J. (1972). On taxation and the control of externalities. *The American Economic Review*, 62(3), 307–322.
- Baybutt, P. (2014). The ALARP principle in process safety. *Process Safety Progress*, 33(1), 36–40. <https://doi.org/10.1002/prs.11599>
- Black, J. (2010a). *Risk-based regulation: Choices, practices and lessons being learnt*. OECD. <https://doi.org/10.1787/9789264082939-11-en>
- Black, J. (2010b). *The role of risk in regulatory processes* (R. Baldwin, M. Cave, & M. Lodge, Eds.) (pp. 302–348). New York, USA: Oxford University Press. <http://ukcatalogue.oup.com/>
- Black, J., & Baldwin, R. (2010). Really responsive risk-based regulation. *Law & Policy*, 32(2), 181–213. <https://doi.org/10.1111/j.1467-9930.2010.00318.x>
- Black, J., & Baldwin, R. (2012). When risk-based regulation aims low: Approaches and challenges. *Regulation & Governance*, 6(1), 2–22. <https://doi.org/10.1111/j.1748-5991.2011.01124.x>
- Bradford, A. (2020, March). *The brussels effect: How the European union rules the world*. Faculty Books. <https://scholarship.law.columbia.edu/books/232>

- Cardona, O. D., Van Aalst, M. K., Birkmann, J., Fordham, M., Mc Gregor, G., Rosa, P., Pulwarty, R. S., et al. (2012, January). Determinants of risk: Exposure and vulnerability. *Managing the risks of extreme events and disasters to advance climate change adaptation: Special report of the intergovernmental panel on climate change*, pp. 65–108. <https://doi.org/10.1017/CBO9781139177245.005>
- Chamberlain, J. (2022, December). The risk-based approach of the European union's proposed artificial intelligence regulation: Some comments from a tort law perspective. *European Journal of Risk Regulation*, 1–13. <https://doi.org/10.1017/err.2022.38>
- Dijck, G. V. (2022). Predicting recidivism risk meets AI act. *European Journal on Criminal Policy and Research*, 28(3), 407–423. <https://doi.org/10.1007/s10610-022-09516-8>
- Dionne, G. (2013). Risk management: History, definition, and critique. *Risk Management and Insurance Review*, 16(2), 147–166. <https://doi.org/10.1111/rmir.12016>
- Fontanelli, F. (2016, January). The court of justice of the European union and the illusion of balancing in internet-related disputes. *The internet and constitutional law: The protection of fundamental rights and constitutional adjudication in Europe*, 94–118. <https://doi.org/10.4324/9781315684048>
- French, S., Bedford, T., & Atherton, E. (2005). Supporting ALARP decision making by cost benefit analysis and multiattribute utility theory. *Journal of Risk Research*, 8(3), 207–223. <https://doi.org/10.1080/1366987042000192408>
- Garrett, B., & Monahan, J. (2020). Judging risk. *California Law Review*, 108(2), 439–493.
- Gupta, A. K., Bharadwaj, M., & Mehrotra, R. (2016). Skin cancer concerns in people of color: Risk factors and prevention. *Asian Pacific Journal of Cancer Prevention: APJCP*, 17(12), 5257–5264. <https://doi.org/10.22034/APJCP.2016.17.12.5257>
- Hacker, P. (2023). The European AI liability directives—Critique of a half-hearted approach and lessons for the future. arXiv. <https://doi.org/10.48550/arXiv.2211.13960>
- Hurst, J., McIntyre, J., Tamauchi, Y., Kinuhata, H., & Kodama, T. (2019). A summary of the 'ALARP' principle and associated thinking. *Journal of Nuclear Science and Technology*, 56(2), 241–253. <https://doi.org/10.1080/00223131.2018.1551814>
- Jones-Lee, M., & Aven, T. (2011). ALARP—What does it really mean? *Reliability Engineering & System Safety*, 96(8), 877–882. <https://doi.org/10.1016/j.res.2011.02.006>
- Kaplan, S., & Garrick, B. J. (1981). On the quantitative definition of risk. *Risk Analysis*, 1(1), 11–27. <https://doi.org/10.1111/j.1539-6924.1981.tb01350.x>
- Karliuk, M. (2022, October). Proportionality principle for the ethics of artificial intelligence. *AI and Ethics*. <https://doi.org/10.1007/s43681-022-00220-1>
- Khanna, V. S. (2021). Compliance as costs and benefits. In B. van Rooij & D. D. Sokol (Eds.), *The Cambridge handbook of compliance* (pp. 13–26) Cambridge Law Handbooks. Cambridge University Press. <https://doi.org/10.1017/9781108759458.002>
- Krebs, J. R. (2011). Risk, uncertainty and regulation. *Philosophical Transactions of the Royal Society A: Mathematical, Physical and Engineering Sciences*, 369(1956), 4842–4852. <https://doi.org/10.1098/rsta.2011.0174>
- Kuran, C. H. A., Morsut, C., Kruke, B. I., Krüger, M., Segnestam, L., Orru, K., Nævestad, T. O., et al. (2020). Vulnerability and vulnerable groups from an intersectionality perspective. *International Journal of Disaster Risk Reduction*, 50(November), 101826. <https://doi.org/10.1016/j.ijdrr.2020.101826>
- Lenaerts, K. (2019). Limits on limitations: The essence of fundamental rights in the EU. *German Law Journal*, 20(6), 779–793. <https://doi.org/10.1017/glj.2019.62>
- Liu, H.-Y., Lauta, K. C., & Maas, M. M. (2018). Governing boring apocalypses: A new typology of existential vulnerabilities and exposures for existential risk research. *Futures, Futures of Research in Catastrophic and Existential Risk*, 102(September), 6–19. <https://doi.org/10.1016/j.futures.2018.04.009>
- Luna, F. (2019). Identifying and evaluating layers of vulnerability—A way forward. *Developing World Bioethics*, 19(2), 86–95. <https://doi.org/10.1111/dewb.12206>
- Mahler, T. (2007). Defining legal risk. SSRN Scholarly Paper. Rochester, NY. <https://papers.ssrn.com/abstract=1014364>
- Maranhão, J., de Souza, E. G., & Sartor, G. (2021). A dynamic model for balancing values. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law, ICAIL '21*, pp. 89–98. New York, NY, USA: Association for Computing Machinery. <https://doi.org/10.1145/3462757.3466143>

- Michelakaki, C., & Vale, S. B. (2023, May). Unlocking data protection by design & by default: Lessons from the enforcement of article 25 GDPR. <https://policycommons.net/artifacts/3838751/fpf-article-25-gdpr-a4-final-digital/4644643/>
- Millstone, E., van Zwanenberg, P., Marris, C., Levidow, L., & Torgersen, H. (2004). Science in trade disputes related to potential risk: Comparative case studies. Other. Seville, Spain: European Commission. <http://ftp.jrc.es/EURdoc/eur21301en.pdf>
- Neuwirth, R. J. (2022). The EU artificial intelligence act: Regulating subliminal AI systems. SSRN Scholarly Paper. Rochester, NY. <https://doi.org/10.2139/ssrn.4135848>
- Ni, H., Chen, A., & Chen, N. (2010). Some extensions on risk matrix approach. *Safety Science*, 48(10), 1269–1278. <https://doi.org/10.1016/j.ssci.2010.04.005>
- Quattrocchio, S. (2020). Artificial intelligence, computational modelling and criminal proceedings: A framework for a European legal discussion. Vol. 4. Legal studies in international, European and comparative criminal law. Cham: Springer International Publishing. <https://doi.org/10.1007/978-3-030-52470-8>.
- Renn, O. (2011). *Risk governance: Coping with uncertainty in a complex world*. Routledge. <https://doi.org/10.4324/9781849772440>
- Rothstein, H., Borraz, O., & Huber, M. (2013). Risk and the limits of governance: Exploring varied patterns of risk-based governance across Europe. *Regulation & Governance*, 7(2), 215–235. <https://doi.org/10.1111/j.1748-5991.2012.01153.x>
- Sartor, G. (2018). A quantitative approach to proportionality. In C. Aitken, A. Amaya, K. D. Ashley, C. Bagnoli, G. Bongiovanni, B. Brožek, C. Castelfranchi, et al. (Eds.), *Handbook of legal reasoning and argumentation* (pp. 613–636). Springer Verlag.
- Sen, A. (2004). Incompleteness and reasoned choice. *Synthese*, 140(1/2), 43–59.
- Simpson, N. P., Mach, K. J., Constable, A., Hess, J., Hogarth, R., Howden, M., Lawrence, J., et al. (2021). A framework for complex climate change risk assessment. *One Earth*, 4(4), 489–501. <https://doi.org/10.1016/j.oneear.2021.03.005>
- Smuha, N., Ahmed-Rengers, E., Harkens, A., Wenlong, L., Maclaren, J., Piselli, R., & Yeung, K. (2021, August). How the EU can achieve legally trustworthy AI: A response to the European commission's proposal for an artificial intelligence act. *Artificial Intelligence - Law, Policy, & Ethics eJournal*. https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3899991
- Sunstein, C. R. (2004). *Risk and reason*. Cambridge Books, Cambridge University Press. <https://ideas.repec.org/b/cup/cbooks/9780521016254.html>
- Tridimas, T. (2018). The Principle of Proportionality. In R. Schütze & T. Tridimas (Eds.), *Oxford principles of European union law: The European union legal order: Volume I*. Oxford University Press. <https://doi.org/10.1093/oso/9780199533770.003.0010>
- Turner, B. L., Kasperson, R. E., Matson, P. A., McCarthy, J. J., Corell, R. W., Christensen, L., Eckley, N., et al. (2003). A framework for vulnerability analysis in sustainability science. *Proceedings of the National Academy of Sciences*, 100(14), 8074–8079. <https://doi.org/10.1073/pnas.1231335100>
- Ufert, F. (2020). AI regulation through the lens of fundamental rights: How well does the GDPR address the challenges posed by AI? *European Papers - A Journal on Law and Integration*, 5(2), 1087–1097. <https://doi.org/10.15166/2499-8249/394>

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.