

Foreseeing the Impact of the Proposed AI Act on the Sustainability and Safety of Critical Infrastructures

The AI Act has been recently proposed by the European Commission to regulate the use of AI in the EU, especially on high-risk applications, i.e. systems intended to be used as safety components in the management and operation of road traffic and the supply of water, gas, heating and electricity. On the other hand, IEC 61508, one of the most adopted international standards for safety-critical electronic components, seem to mostly forbid the use of AI in such systems. Given this conflict between IEC 61508 and the proposed AI Act, also stressed by the fact that IEC 61508 is not an harmonised European standard, with the present paper we study and analyse what is going to happen to industry after the entry into force of the AI Act. More in detail we focus on how the proposed AI Act might positively impact on the sustainability of critical infrastructures by allowing the use of AI on an industry where it was previously forbidden. To do so, we provide several examples of AI-based solutions falling under the umbrella of IEC 61508 that might have a positive impact on sustainability in alignment with the current long-term goals of the EU and the Sustainable Development Goals of the United Nations, i.e. affordable and clean energy, sustainable cities and communities.

Additional Key Words and Phrases: AI Act, IEC 61508, safety standards, sustainability

ACM Reference Format:

. 2022. Foreseeing the Impact of the Proposed AI Act on the Sustainability and Safety of Critical Infrastructures. 1, 1 (December 2022), 10 pages. <https://doi.org/10.1145/nnnnnnnn.nnnnnnnn>

1 INTRODUCTION

Harnessing the full potential of AI could lead our society to more efficient and thus sustainable energy production, storage and transportation, in accordance with many of the Sustainable Development Goals of the United Nations [34], including: affordable and clean energy (goal 7), industry, innovation and infrastructure (goal 9), sustainable cities and communities (goal 11), responsible consumption and production (goal 12), climate action (goal 13). Indeed, AI can be used to optimally recognise, predict, detect, identify, determine, control, generate, and classify [33] in a wide range of tasks, sometimes also achieving or exceeding human performance in problems such as strategy games [17, 43], image and object recognition [19, 38], etc.. Nonetheless, the adoption of AI-based technological solutions for more sustainable energy (e.g. for decreasing the carbon emissions of coal-fired thermal power plants [44]) has been held back in the last decades by conservative international standards, i.e. IEC 61508 [45], a standard that regulates safety-critical electronic components and that practically forbid AI in many critical infrastructures.

Despite this, in 2021, the European Commission published a proposal of AI Act¹ [9] that is expected to become a legally binding regulation to all the Member States of the EU by 2024. Importantly, the objective of the AI Act is to set a common regulatory and legal framework for AI that applies to all sectors (except for military), and to all types of artificial intelligence, including (high-risk) AI for the *management and operation of critical infrastructure*.

¹EUR-Lex - 52021PC0206 - EN - EUR-Lex

Author's address:

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2022 Association for Computing Machinery.

Manuscript submitted to ACM

Manuscript submitted to ACM

53 Considering that one of the goals of the proposed AI Act is to regulate the use of AI also on those systems covered
54 by IEC 61508, i.e. ‘systems intended to be used as safety components in the management and operation of road traffic
55 and the supply of water, gas, heating and electricity’ (see Annex III, point 2.a), the research questions we are trying to
56 answer with the present letter are the following:
57

- 58 • Will the AI Act be disruptive with respect to IEC 61508?
- 59 • What will happen to those industries currently regulated by IEC 61508?

61 In fact, we believe that answering these has the potential to help both industry and academia to quickly seize the
62 opportunities offered by the new European policies enshrined in the AI Act.

63 This paper is structured as follows. Section 2 discusses the adopted methodology. In Section 3 we give enough
64 background to understand IEC 61508 and its implications for industry. While, in Section 4 we analyse the position
65 of the AI Act on IEC 61508 and other standards, providing in Section 5 our understanding of how AI could improve
66 sustainability and energy efficiency whilst maintaining safety. Finally, in Section 6 we try to give a conclusive answer
67 to each research question, discussing the consequences of our findings as well as some possible issues.
68
69

70 2 METHODOLOGY

71 In order to answer these questions, we analyse the main differences between IEC 61508 and the proposed AI Act.
72 Then, we identify significant and concrete examples of technical solutions that could increase sustainability and energy
73 efficiency but which, at the moment, are not feasible due to IEC 61508. Furthermore, we also study how such new
74 technological solutions might impact industry, trying to understand how disruptive the AI Act could be by loosening
75 the tight laces of IEC 61508 on AI. Hence, we try to align our findings to the medium- and long-term objectives of the
76 EU on sustainability and support for innovation.
77
78
79

80 3 BACKGROUND

81 With this section we provide a minimal amount of information about safety, AI, the IEC 61508 standard and the proposed
82 AI Act.
83
84

85 3.1 The safety standard IEC 61508

86 IEC 61508 [23] is an international standard describing how to design, deploy and maintain an Electrical/Electronic/Programmable
87 Electronic safety-related system. Examples of safety-related systems to which IEC 61508 can be applied are: emergency
88 shut-down systems, remote monitoring, operation or programming of a network-enabled process plant, information-
89 based decision support tool where erroneous results may affect safety. More in details, programmable electronic
90 safety-related systems typically incorporate programmable controllers, programmable logic controllers, micropro-
91 cessors, application specific integrated circuits, or other programmable devices (e.g. ‘smart’ devices such as sen-
92 sors/transmitters/actuators). The focus is in particular on safety functions and on the relative level of risk reduction
93 that they provide. Those levels are grouped in four Safety Integrity Levels (SILs), the higher the Safety Integrity Level
94 the greater the risk of failure.
95
96
97
98

99 Notice that it is expected² that IEC 61508 can be published as EN 61508, an European standard, but it does not have
100 the status of a *harmonized* European standard in relation to any EC product directive and it is not therefore listed in the
101 EC Official Journal. However, this does not prevent compliance with relevant parts of EN 61508 being used to support
102

103 ²<https://www.iec.ch/functional-safety/faq>

105 a declaration of conformity with an EC product directive, if that is appropriate. In any cases, IEC 61508 is followed
106 worldwide³.
107

108 3.2 Safety vs AI 109

110 Quoting [46]: ‘ There is no such thing as zero risk. This is because no physical item has zero failure rate, no human
111 being makes zero errors, and no piece of software design can foresee every operational possibility ’. Thus, perfect *safety*,
112 i.e., the absence of catastrophic consequences on the user(s) and the environment [3], is out of reach. During the last
113 decades, several standards on how to develop hardware and software artifacts in safety critical contexts have been
114 defined. These standards crystallize lessons learned, common practices and scientific research into concrete guidelines.
115 Each industry sector has its own standard, but the idea behind all of them is the same: a risk-based approach that
116 characterize the entire product life-cycle.
117

118
119 With respect to Artificial Intelligence (AI), at row 5 in Tables A.2 and C.2, Part 3, of IEC 61508 [23], it is clearly stated
120 that AI is not recommended for Safety Integrity Level 2 or above because it may complicate the achievement of one
121 or more of the following properties: correctness with respect to software safety requirements specification, freedom
122 from intrinsic design faults, simplicity and understandability, related with the observability-in-depth principle, aimed at
123 avoiding as much as possible a false sense of safety due to lack of information, predictability of behaviour, verifiable
124 and testable design.
125

126 IEC 61508 has influenced other standards [46], here called ‘second tier standards’, that are as rigid as IEC 61508 Part
127 3 with respect to AI. Among those, examples are software for train EN 50128 [11], process industry [24] and machinery
128 IEC 62061 [25]. Parallel to the family of standards originated from IEC 61508, other really important examples where AI
129 is banned, for high Safety Integrity Level, from computer-based systems employed in nuclear power plants, IEC 60880
130 [22], and avionic, DO-178 C [37].
131

132 To the best of the authors’ knowledge, the only safety standard that allows the employment of AI (because it does
133 not mention it, and then it is not ‘not recommended’) is ISO 26262 for the automotive industry sector [14, 20, 42].
134
135

136 3.3 The Proposed AI Act 137

138 The AI Act [9] is a proposed European law on AI. Differently from other domains, this act is specific to AI systems
139 and requires an *ad hoc* discussion rather than the framing of these systems in the discussion of other legal domains.
140 This is because AI technologies are not placed within an existing legal framework (e.g. banking), but the whole legal
141 framework (i.e. the proposed AI Act) is built around AI technologies.
142

143 The proposed AI Act assigns applications of AI to three risk categories. First, applications and systems that create
144 an unacceptable risk, such as government-run social scoring of the type used in China, are banned. Second, high-risk
145 applications, such as a CV-scanning tool that ranks job applicants, are subject to specific legal requirements. Lastly,
146 applications not explicitly banned or listed as high-risk are largely left unregulated.
147

148 Examples of high-risk AI are given by the proposed AI Act in Annex III, as they broadly include applications for:
149 biometric identification and categorisation of natural persons, management and operation of critical infrastructures,
150 etc. More in detail, for all those applications defined as ‘high-risk’, the AI Act provides several limitations and safety
151 assurance procedures including: a risk management system (art. 9), appropriate data governance and management
152 practices (art. 10), detailed technical documentation (art. 11).
153

154
155 ³<https://www.iec.ch/national-committees>
156

157 Finally, the AI Act defines in Annex I what are the AI techniques and approaches referred by the proposal. Among
158 them we have: machine learning approaches (e.g. neural networks), logic- and knowledge-based approaches (e.g.
159 inductive logic programming), and statistical approaches (e.g. Bayesian estimation).
160

161 4 ANALYSIS OF THE POSITION OF THE AI ACT ON IEC 61508 AND OTHER STANDARDS

162 It is crystal clear from the European Green Deal⁴ and the EU's commitment to global climate action under the Paris
163 Agreement that one of the big goals of the EU is to be climate-neutral by 2050. Citing the words of the European
164 Commission: 'The EU can lead the way [to climate-neutrality] by investing into realistic technological solutions,
165 empowering citizens and aligning action in key areas such as industrial policy, finance and research, while ensuring
166 social fairness for a just transition.' The reason why we are citing these statements is that we are going to use them
167 as interpretative key for the proposed AI Act, especially with respect to the importance of article 54.1.a, stating that
168 'innovative AI systems shall be developed for safeguarding substantial public interest in [...] a high level of protection
169 and improvement of the quality of the environment'.
170

171 In fact, the AI Act is (as mentioned in Section 3) regulating a vast range of AI applications, with due focus on those
172 listed as high-risk in Annex III. More in detail, it covers, among others, the AI applications for the 'management and
173 operation of critical infrastructure', i.e. the 'AI systems intended to be used as safety components in the management
174 and operation of road traffic and the supply of water, gas, heating and electricity.' But, considering that many critical
175 infrastructures are currently following a non-harmonized IEC 61508 standard that is de facto excluding the involvement
176 of any AI, we can see a non-alignment of it to the proposed AI Act.
177

178 Indeed, according to article 40 only the 'harmonised standards or parts thereof the references of which have been
179 published in the Official Journal of the European Union' are considered to be in conformity with the requirements set
180 out in the proposed AI Act for high-risk AI systems (see Chapter 2 of Title III). In other words, article 40, together
181 with the *Explanatory Memorandum*, article 54.1.a and the fact that IEC 61508 is a non-harmonised standard, make us
182 understand that the intent of the proposed AI Act is to promote innovative AI systems also for the 'management and
183 operation of road traffic and the supply of water, gas, heating and electricity'. As consequence, we envisage that the
184 proposed AI Act, without further modifications, can have a disruptive effect in the industry of critical infrastructures.
185 This effect can be disruptive in a positive way, by opening to new technological solutions that have the potential to
186 improve even further our quality of life, reducing costs and increasing efficiency. Nonetheless, it can be disruptive
187 also in a negative way, by ceding the control of critical infrastructures to automatic decision makers that are possibly
188 opaque, greedy, unfair and non-transparent in a way that would not allow to understand where the responsibility lies.
189

190 Although, despite the fact that IEC 61508 is a non-harmonized standard, thus not covered by article 40, we can see
191 that the proposed AI Act shares several and important similarities with it, suggesting that it is not the intent of the EU
192 Commission to fully upset existing standards.
193

194 Overall, we see that the intent of the proposed AI Act is to modernize existing critical infrastructures, to make them
195 more sustainable. To do so, the AI Act does not ignore or try to eliminate the currently adopted standards, although
196 it wants them harmonised with the EU's policies. This is why the CEN-CENELEC has established a joint technical
197 committee on AI⁵ and defined a road map for AI standardization [39] that includes the harmonization of IEC 61508
198 and other standards. In fact, according to article 2(1)(c) of Regulation (EU) No 1025/2012, the CEN-CENELEC is the
199

206 ⁴<https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=COM:2019:640:FIN>

207 ⁵<https://www.cenelec.eu/areas-of-work/cen-cenelec-topics/artificial-intelligence>

Table 1. **AI Act vs IEC 61508**: AI Act is centred on transparency while IEC 61508 on safety. This table shows how the proposed AI Act and IEC 61508 address the same process for risk-assessment, analysis, development and document production in different ways.

| | Differences | |
|--|--|---|
| | IEC 61508 | Proposed AI Act |
| Risk-based approach, in particular to establish the belongings to a pre-defined category | Quantitative (hazard analysis, risk assessment and identify the Safety Integrity Level) | Qualitative (one of the alternatives: no risk, application listed in Annex III, AI not applicable) |
| Normalised life cycle, with focus on accountability | V-shape development (focus on modularity and decomposability) | Clear definition of datasets (focus on data management, in particular for training the AI and how to use the product) |
| Ex-ante and ex-post analysis | Statistical methods (hardware), study of qualitative techniques (hardware and software) and structured testing campaigns | Declarative (identify high level characteristics, provide general description of components behaviour) |
| Document production | Assessment performed by an external institution | Fill a form in EU database, part of the information is of public domain (focus on transparency) |

European Union (EU) authority for standards. Nonetheless, we can see also that European countries start producing guidelines and roadmaps [47] on this subject.

So, given the very clear position of the proposed AI Act with respect to the possibility of using AI systems in particular critical infrastructures, we believe that the CEN-CENELEC, together with IEC will adapt IEC 61508, eventually opening to a safe use of AI systems also in critical infrastructures. Importantly, the CEN-CENELEC [39] has already identified article 41 as a possible source of uncertainty in industry, given that it would explicitly cut out any non-harmonized IEC standard, i.e. IEC 61508. In fact, article clearly 41.1 says that: ‘Where harmonised standards referred to in Article 40 do not exist or where the Commission considers that the relevant harmonised standards are insufficient or that there is a need to address specific safety or fundamental right concerns, the Commission may, by means of implementing acts, adopt common specifications in respect of the requirements set out in Chapter 2 of [Title III]. Those implementing acts shall be adopted in accordance with the examination procedure referred to in Article 74(2).’

Consequently, given all the aforementioned facts, we see an harmonization of IEC 61508 or its replacement by 2024, and this will open to at least one two scenarios. In the first scenario, we will have an opening to the use of AI systems in the context of critical infrastructures, whereas they can improve sustainability whilst guaranteeing safety. While in a second scenario a very strict policy against AI systems in critical infrastructures will be maintained.

Again, as consequence of the analysis presented in this Section, we believe that this very first scenario is the most likely. If that is correct, we envisage that a new stream of research on AI for critical infrastructures will be opening by the end of 2024, paving the way for AI systems to improve the sustainability of our society. Nonetheless, it is important to stress that the use of AI does not come free of problems related to safety, fairness, transparency and sometimes even sustainability. For this reason, in the following section we will discuss and classify existing AI techniques, to analyze their impact on sustainability and safety and to understand which AI-based solutions are likely to be allowed by a future harmonization of IEC 61508.

Table 2. **AI on Safety-Critical Environments:** This table shows examples of possible applications of AI on some safety-critical contexts. For each context we identify its Safety Integrity Level (SIL) and possible tasks where AI can be deployed to improve sustainability.

| SIL | Context | Use of AI |
|-------|------------------------------------|--|
| 4 | Nuclear power plant [18] | <ul style="list-style-type: none"> Anomaly detection [5, 6] In-core full management [35] |
| 3 | Railway, station management [36?] | <ul style="list-style-type: none"> Turning on/off switch heaters [8] Fault detection of sensitive components [2] |
| 2 & 1 | Chemical industry [1, 12] | <ul style="list-style-type: none"> Predicting chattering alarms [48] Plant health diagnosis [50] |

5 CLASSIFICATION AND DISCUSSION OF THE IMPACT OF AI ON SUSTAINABILITY AND SAFETY

As mentioned in Section 3.3, the techniques and approaches covered by the proposed AI Act include both symbolic (e.g. logic-based) and non-symbolic (e.g. neural networks, statistics) techniques. Nonetheless, each different type of AI may have its own characteristics, impacting on safety differently from others. Indeed, as suggested by Mohseni et al. in their taxonomy of machine learning safety [32], the decisions of state-of-the-art machine learning techniques can be, in some cases, completely unexplainable, non-transparent, biased and non-robust. On the other hand, the automatic decisions of fully symbolic approaches can be explainable by design but not as good as those of a state-of-the-art neural network [21]. Therefore, given such a trade-off between explainability and performance, being able to foresee and analyse the impact of AI on safety is not trivial, forcing us to analyse it differently for different types of applications.

This is why in the present letter we will study the impact of AI, on the safety and sustainability of critical infrastructures, by using as reference point the 4 Safety Integrity Level defined by IEC 61508. In fact, for each safety level we will show concrete examples of technological solutions based on AI that have the potential for significantly improving sustainability, analyzing what is the trade-off between sustainability and safety and how important that is.

5.1 Examples of "Forbidden" AI-based Solutions for Improving Sustainability in Safety-Critical Systems

Safety-critical subsystems of cyber-physical systems⁶ compliant with IEC 61508 are required to have the properties listed in Section 3.1 and these normally do not include AI. Nonetheless, few examples can illustrate how impactful can be AI on safety-critical systems, considering that in scientific and technical literature are available several studies that directly address the issue or propose promising approaches that well fit the kind of data relevant for safety-critical functions. In Table 2 we show these examples aligned to the Safety Integrity Level of IEC 61508.

Safety Integrity Level 4 systems compliant to IEC 61508 are quite rare. Nevertheless, the nuclear power plants industry offers examples of such systems [29]. Here, AI is envisioned to have a great impact in the relatively close future, in particular for safeguard and surveillance (filter and identify signatures of nuclear materials), monitoring and diagnosis of severe accidents or nuclear power plant transients [51]. All these actions are crucial to avoid environmental consequences of accidents and lives lost [18, 35].

For Safety Integrity Level 3 consider the railway industry, and in particular those systems for which energy efficiency is crucial, with focus on the heating system for rail-road switches [8]. This is a critical subsystem, responsible for keeping the switches free from snow and ice, necessary to guarantee the correct operation of the switches and so the correct train routing (always turned on increases safety). Depending on the climate conditions of the place where the railway

⁶A cyber-physical system comprises physical mechanisms that are monitored and/or controlled by Information and Communication Technologies.

313 system operates, the energy consumed by this heating system can be very relevant, i.e. the heating always turned on
314 implies greater ambient impact and cost. To provide concrete examples, in [36] it is reported that the cost for heating
315 the 6800 switches and crosses in Sweden can amount to 10 – 15 Million Euros per year. In Germany, Deutsche Bahn
316 (DB) alone has 64000 switches heated with electrical resistance and gas heaters, a combined power of 900 MW which
317 consume up to 230 GWh/year [?]. AI is expected to empower this subsystem with snow or ice prediction/detection and
318 by making the turning on/off algorithm more responsive.
319

320 Regarding Safety Integrity Level 2 and 1, [4] shown concrete examples in the process industry, where the second tier
321 standard IEC 61511 [24] apply, that can be generalized to the chemical industry. The chemical industry is one of the
322 most energy-intensive manufacturing industries and a major source of greenhouse gas emissions. Besides that, chemical
323 production often involves hazardous materials and high-pressure/high-temperature conditions, which may lead to fire,
324 explosion, and other types of chemical accidents. Those chemical accidents could cause casualties, financial and social
325 losses [30]. According to a survey conducted by Accenture [1], 94% of the executives of companies in the chemical
326 and advanced materials industry expect an industry-wide digitisation, and AI plays an essential role in enabling the
327 digital revolution [12]. In particular, fault detection and diagnosis is crucial to both safety and sustainability. As an
328 example, consider fault detection for a Tennessee Eastman process (chapter 8 of [41]) with few modes, where unit
329 operations include a reactor, a condenser, a recycle compressor, a vapour-liquid separator and a stripper [50]. Notice
330 that the adoption of AI is not ‘not recommended’ for Safety Integrity Level 1 in IEC 61508.
331

332 Indeed, AI has the potential to cope with high dimensional data, being able to generalise, handling novel inputs and
333 incomplete knowledge [28]. These features are expected [49] to greatly impact the way goals and targets in the 2030
334 Agenda for Sustainable Development are addressed.
335

336 Overall, we can say that AI may be critical to anomaly detection, for taking timely countermeasures, being able to
337 find patterns in data that do not conform to expected behaviour [7].
338

341 5.2 Discussion

342 Even though several metrics for AI performance and robustness appeared in literature and have been tested in several
343 contexts [52], only preliminary ones have been defined specifically to address safety or sustainability (e.g., [13, 15]), and
344 are yet to be tested extensively before some AI can become amenable for safety critical applications (where quantification
345 has a central role). Thus, it is expected that those AI for which will be available reliable metrics will be the first to be
346 employed in safety functions or safety critical systems.
347

348 It is desirable that interpretable or explainable-by-design AI [31] are the first to be employed, in particular for handling
349 tabular data [40]. This is indeed expected to cover, at least in part, simplicity, understandability and observability-in-depth
350 (Section 3.2).
351

352 The heart of the problem is that AI is difficult to be framed in safety standards because of the way it fails. Deter-
353 ministic software fails systematically, whereas hardware fails randomly [46]. Safety standards recommend to address
354 hardware failures through statistical methods and mitigate/tolerate deterministic software failures employing qualitative
355 techniques. In some standards, statistical methods for quantifying software failures are allowed (e.g., suggestions are
356 provided in Part 7 of IEC 61508 [23]) in others (e.g., DO-178 C [37]) are not recommended. After about forty years of
357 discussions, in industry and academia, no consensus has been reached, and strong opinions continue to emerge [10].
358

359 Among those listed as AI in Annex I of the proposed AI Act, some (e.g., statistical models or neural networks) are
360 intrinsically non-deterministic [27], and then does not fit current safety standards framework. Seen from a different
361 perspective, though, this removes many of the assumptions that prevent the use of statistical methods, opening up new
362
363
364

ways to address AI failures. Indeed, a positive byproduct of the discussions on statistical methods for deterministic software is the huge body of knowledge that is available but not enough explored for addressing non-deterministic software.

6 CONCLUSIONS

First of all, with this paper we performed an analysis of how the proposed AI Act might impact on the sustainability and safety of critical systems (e.g. power plants). We did it by looking at the differences, incompatibilities and similarities of the AI Act with IEC 61508, one of the most important non-harmonised standards for safety-critical infrastructures. Importantly, among the main differences, we show the incompatibility of IEC 61508 with the use of any AI in systems requiring a Safety Integrity Level greater than 1, pointing to the disruptive effect that the proposed AI Act might have on that part of industry aligned with IEC 61508. Then, we identified examples of AI-based solutions falling under the umbrella of IEC 61508 with a Safety Integrity Level greater than 1 that might have a positive impact on sustainability in alignment with the current long-term goals of the EU and the proposed AI Act.

Eventually, we collected enough material to answer our initial research questions and foresee a future where critical infrastructures may harness the full potential of AI to improve both sustainability and safety in accordance with the following Sustainable Development Goals of the United Nations [34]: affordable and clean energy (goal 7), industry, innovation and infrastructure (goal 9), sustainable cities and communities (goal 11), responsible consumption and production (goal 12), climate action (13). To be more precise, in accordance with the analysis we carried out in this paper, we believe that the AI Act will eventually soften the position of IEC 61508 with respect to AI, leading to a new generation of critical infrastructures harmonised with the European vision embodied by the proposed AI Act. This would clearly open to new research and technological solutions on this topic by the end of 2024.

Overall, with this paper, our focus was exclusively on those safety-critical contexts where AI is expected to enhance economic/environmental aspects of sustainability but is not employed yet because considered not enough mature or potentially in conflict with safety or technical aspects of sustainability, as per IEC 61508. Nonetheless, despite the promises made by state-of-the-art AI we can sceptically argue that using AI in safety-critical systems does definitely come with a risk. This risk is posed by the fact that ceding control to machines might lead to new unregulated unethical and immoral behaviours as well as a dangerous lack of transparency and accountability. Importantly, with respect to this specific issue, there are several flourishing discussions in literature and among policy makers, also taking into account that similar issues are addressed in other contexts as well [16]. This gives us hope that the technology of the future will be able to cope with such urgent problems to give us solutions based on AI capable of addressing the sustainability goals that have been set for the future. For this reason, we argue that any forthcoming harmonised version of IEC 61508 is unlikely to completely close to application of AI in safety-critical systems with a Safety Integrity Level greater than 1. This is why we are all waiting for the CEN-CENELEC and its technical commission to give us a final answer to our research questions in the form of new harmonised standards.

REFERENCES

- [1] Accenture. 2014. global digital chemicals survey.
- [2] J. A. Agirre, L. Etxeberria, R. Barbosa, S. Basagiannis, G. Giantamidis, T. Bauer, E. Ferrari, M. Labayen Esnaola, V. Orani, J. Öberg, D. Pereira, J. Proença, R. Schlick, A. Smrčka, W. Tiberti, S. Tonetta, M. Bozzano, A. Yazici, and B. Sangchoolie. 2021. The VALU3S ECSEL project: Verification and validation of automated systems safety and security. *Microprocessors and Microsystems* 87 (2021), 104349.
- [3] A. Avizienis, J.-C. Laprie, B. Randell, and C. Landwehr. 2004. Basic concepts and taxonomy of dependable and secure computing. *IEEE Transactions on Dependable and Secure Computing* 1, 1 (2004), 11–33.

- 417 [4] D. Barone and A. Damiani. 2016. Esperienza pratica nella applicazione delle analisi SIL (IEC 61508/61511) relative ai sistemi di sicurezza ad alta
418 affidabilità, per uno stabilimento a rischio di incidente rilevante. (2016). <http://conference.ing.unipi.it/vgr2016/images/papers/133.pdf> Valutazione e
419 Gestione del Rischio negli Insediamenti Civili ed Industriali.
- 420 [5] Roger Boza. 2019. *Subtle Process Anomalies Detection using Machine Learning Methods*. Technical Report. U.S. Department of Energy, Office of
421 Nuclear Energy.
- 422 [6] F. Calivá, F. S. De Ribeiro, A. Mylonakis, C. Demazi'ere, P. Vinai, G. Leontidis, and S. Kollias. 2018. A Deep Learning Approach to Anomaly Detection
423 in Nuclear Reactors. In *International Joint Conference on Neural Networks*. 1–8.
- 424 [7] V. Chandola, A. Banerjee, and V. Kumar. 2009. Anomaly Detection: A Survey. *ACM Comput. Surv.* 41, 3, Article 15 (2009).
- 425 [8] S. Chiaradonna, G. Masetti, F. Di Giandomenico, F. Righetti, and C. Vallati. 2021. Enhancing sustainability of the railway infrastructure: Trading
426 energy saving and unavailability through efficient switch heating policies. *Sustainable Computing: Informatics and Systems* 30 (2021), 100519.
- 427 [9] European Commission. 2021. Proposal for a Regulation of the European Parliament and of the Council laying down harmonised rules on Artificial
428 Intelligence and amending certain union legislative acts. <https://eur-lex.europa.eu/legal-content/EN/TXT/?uri=CELEX%3A52021PC0206>.
- 429 [10] D. Daniels and N. Tudor. 2022. Software Reliability and the Misuse of Statistics. *Safety-Critical Systems eJournal* 1, 1 (2022).
- 430 [11] European Committee for Electrotechnical Standardization 2020. *Railway applications - Communication, signalling and processing systems - Software
431 for railway control and protection systems*. European Committee for Electrotechnical Standardization.
- 432 [12] World Economic Forum. 2017. Digital transformation initiative chemistry and advanced materials industry. [http://reports.weforum.org/digital-
433 transformation/wp-content/blogs.dir/94/mp/files/pages/files/white-paper-dti-2017-chemistry.pdf](http://reports.weforum.org/digital-transformation/wp-content/blogs.dir/94/mp/files/pages/files/white-paper-dti-2017-chemistry.pdf).
- 434 [13] M. Gharib and A. Bondavalli. 2019. On the Evaluation Measures for Machine Learning Algorithms for Safety-Critical Systems. In *15th European
435 Dependable Computing Conference*. 141–144.
- 436 [14] M. Gharib, P. Lollini, M. Botta, E. Amparore, S. Donatelli, and A. Bondavalli. 2018. On the Safety of Automotive Systems Incorporating Machine
437 Learning Based Components: A Position Paper. In *48th Annual IEEE/IFIP International Conference on Dependable Systems and Networks Workshops*.
438 271–274.
- 439 [15] M. Gharib, T. Zoppi, and A. Bondavalli. 2021. Understanding the Properness of Incorporating Machine Learning Algorithms in Safety-Critical
440 Systems. In *Proceedings of the 36th Annual ACM Symposium on Applied Computing*. 232–234.
- 441 [16] M. Ghassemi, L. Oakden-Rayner, and A. L. Beam. 2021. The false hope of current approaches to explainable artificial intelligence in health care. *The
442 Lancet, digital health* 3 (2021), E745–E750. Issue 11.
- 443 [17] Elizabeth Gibney et al. 2016. Google AI algorithm masters ancient game of Go. *Nature* 529, 7587 (2016), 445–446.
- 444 [18] M. Gomez-Fernandez, K. Higley, A. Tokuhito, K. Welter, W.-K. Wong, and H. Yang. 2020. Status of research and development of learning-based
445 approaches in nuclear science and engineering: A review. *Nuclear Engineering and Design* 359 (2020), 110479.
- 446 [19] Kaiping He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2015. Delving deep into rectifiers: Surpassing human-level performance on imagenet
447 classification. In *Proceedings of the IEEE international conference on computer vision*. 1026–1034.
- 448 [20] J. Henriksson, M. Borg, and C. Englund. 2018. Automotive safety and machine learning: initial results from a study on how to adapt the ISO 26262
449 safety standard. 47–49.
- 450 [21] Andreas Holzinger. 2018. From machine learning to explainable AI. In *2018 world symposium on digital intelligence for systems and machines (DISA)*.
451 IEEE, 55–66.
- 452 [22] International Electrotechnical Commission 2006. *Nuclear power plants – Instrumentation and control systems important to safety – Software aspects
453 for computer-based systems performing category A functions*. International Electrotechnical Commission.
- 454 [23] International Electrotechnical Commission 2010. *Functional safety of electrical/electronic/programmable electronic safety-related systems*. International
455 Electrotechnical Commission.
- 456 [24] International Electrotechnical Commission 2016. *Functional safety - Safety instrumented systems for the process industry sector*. International
457 Electrotechnical Commission.
- 458 [25] International Electrotechnical Commission 2021. *Safety of machinery - Functional safety of safety-related control systems*. International Electrotechnical
459 Commission.
- 460 [26] IUC16 International Union of Railways. [n. d.]. Technologies and Potential Developments for Energy Efficiency and CO2 Reduction in Rail
461 Systems. [https://uic.org/IMG/pdf/_27_technologies_and_potential_developments_for_energy_efficiency_and_co2_reductions_in_rail_systems.
462 _uic_in_colaboration.pdf](https://uic.org/IMG/pdf/_27_technologies_and_potential_developments_for_energy_efficiency_and_co2_reductions_in_rail_systems_uic_in_colaboration.pdf). Online; accessed 15 January 2019.
- 463 [27] B. Johnson. 2022. Metacognition for artificial intelligence system safety – An approach to safe and desired behavior. *Safety Science* 151 (2022),
464 105743.
- 465 [28] Z. Kurd, T. Kelly, and J. Austin. 2006. Developing Artificial Neural Networks for Safety Critical Systems. *Neural Comput. Appl.* 16, 1 (2006), 11–19.
- 466 [29] J. Lahtinen, M. Johansson, J. Ranta, H. Harju, and R. Nevalainen. 2010. Comparison between IEC 60880 and IEC 61508 for Certification Purposes in
467 the Nuclear Domain. In *Computer Safety, Reliability, and Security*. 55–67.
- 468 [30] M. Liao, K. Lan, and Y. Yao. 2022. Sustainability implications of artificial intelligence in the chemical industry: A conceptual framework. *Journal of
469 industrial ecology* 26, 1 (2022), 164–182.
- 470 [31] R. Marcinkevičs and J. E. Vogt. 2020. Interpretability and Explainability: A Machine Learning Zoo Mini-tour. *ArXiv* 2012.01805 (2020).
- 471 [32] S. Mohseni, H. Wang, Z. Yu, C. Xiao, Z. Wang, and J. Yadawa. 2021. Taxonomy of Machine Learning Safety: A Survey and Primer. [https://
472 arxiv.org/abs/2106.04823](https://arxiv.org/abs/2106.04823)

- 469 [33] Mark Munro, Jacob Whiton, and Robert Maxim. 2019. What jobs are affected by AI? (2019).
- 470 [34] United Nations. 2016. *The Sustainable Development Goals 2016*. Technical Report. eSocialSciences.
- 471 [35] E. Nissan. 2019. An Overview of AI Methods for in-Core Fuel Management: Tools for the Automatic Design of Nuclear Reactor Core Configurations
- 472 for Fuel Reload, (Re)arranging New and Partly Spent Fuel. *Designs* 3, 3 (2019).
- 473 [36] A. Parida P. Norbbin, J. Lin. 2016. Energy efficiency optimization for railway switches & crossings: a case study in Sweden. In *WCRR 2016, 11th*
- 474 *World Congress on Railway Research*. SPARK knowledge sharing portal.
- 475 [37] Radio Technical Commission for Aeronautics 2012. *Software Considerations in Airborne Systems and Equipment Certification*. Radio Technical
- 476 Commission for Aeronautics.
- 477 [38] Pranav Rajpurkar, Jeremy Irvin, Kaylie Zhu, Brandon Yang, Hershel Mehta, Tony Duan, Daisy Ding, Aarti Bagul, Curtis Langlotz, Katie Shpanskaya,
- 478 et al. 2017. Chexnet: Radiologist-level pneumonia detection on chest x-rays with deep learning. *arXiv preprint arXiv:1711.05225* (2017).
- 479 [39] CEN-CENELEC Focus Group Report. 2020. Road Map on Artificial Intelligence (AI).
- 480 [40] C. Rudin. 2019. Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead. *ArXiv*
- 481 1811.10154 (2019).
- 482 [41] E. L. Russell, L. H. Chiang, and R. D. Braatz. 2000. *Data-driven Methods for Fault Detection and Diagnosis in Chemical Processes*.
- 483 [42] R. Salay, R. Queiroz, and K. Czarnecki. 2017. An Analysis of ISO 26262: Using Machine Learning Safely in Automotive Software. <https://arxiv.org/abs/1709.02435>
- 484 [43] Julian Schrittwieser, Ioannis Antonoglou, Thomas Hubert, Karen Simonyan, Laurent Sifre, Simon Schmitt, Arthur Guez, Edward Lockhart, Demis
- 485 Hassabis, Thore Graepel, et al. 2020. Mastering atari, go, chess and shogi by planning with a learned model. *Nature* 588, 7839 (2020), 604–609.
- 486 [44] T Sivageerthi, Bathrinath Sankaranarayanan, Syed Mithun Ali, and Koppihraj Karuppiah. 2022. Modelling the Relationships among the Key
- 487 Factors Affecting the Performance of Coal-Fired Thermal Power Plants: Implications for Achieving Clean Energy. *Sustainability* 14, 6 (2022), 3588.
- 488 [45] David J Smith and Kenneth GL Simpson. 2020. *The Safety Critical Systems Handbook: A Straightforward Guide to Functional Safety: IEC 61508 (2010*
- 489 *Edition), IEC 61511 (2015 Edition) and Related Guidance*. Butterworth-Heinemann.
- 490 [46] David J. Smith and Kenneth G. L. Simpson (Eds.). 2020. *The Safety Critical Systems Handbook* (fifth edition ed.).
- 491 [47] DKE standards. 2020. German standardization roadmap on artificial intelligence. <https://www.din.de/resource/blob/772610/e96c34dd6b12900ea75b460538805349/normungsroadmap-en-data.pdf>
- 492 [48] N. Tamascelli, N. Paltrinieri, and V. Cozzani. 2020. Predicting chattering alarms: A machine Learning approach. *Computers & Chemical Engineering*
- 493 143 (2020), 107122.
- 494 [49] R. Vinuesa, H. Azizpour, I. Leite, M. Balaam, V. Dignum, S. Domisch, A. Felländer, S. D. Langhans, M. Tegmark, and F. Fuso Nerini. 2020. The role of
- 495 artificial intelligence in achieving the Sustainable Development Goals. *Nature communications* 11, 1 (2020), 1–10.
- 496 [50] Hao W. and Jinsong Z. 2020. Fault detection and diagnosis based on transfer learning for multimode chemical processes. *Computers & Chemical*
- 497 *Engineering* 135 (2020), 106731.
- 498 [51] F. Wastin, K. Simola, Tanarro C. J., A. Liessens, Z. Simic, and O. Eulaerts. 2019. *Horizon Scanning for Nuclear Safety, Security and Safeguards Yearly*
- 499 *Report*. Technical Report. European Union. EUR 30153 EN.
- 500 [52] T. Wu, Y. Dong, Z. Dong, A. Singa, X. Chen, and Y. Zhang. 2020. Testing Artificial Intelligence System Towards Safety and Robustness: State of the
- 501 Art. *International Journal of Computer Science* 47, 3 (2020).

502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520