

Please cite the Published Version

Crockett, Keeley Alexandra, Gerber, Luciano, Latham, Annabel and Colyer, Edwin (2021) Building Trustworthy AI Solutions: A Case for Practical Solutions for Small Businesses. IEEE Transactions on Artificial Intelligence. ISSN 2691-4581

DOI: https://doi.org/10.1109/tai.2021.3137091

Publisher: Institute of Electrical and Electronics Engineers (IEEE)

Version: Published Version

Downloaded from: https://e-space.mmu.ac.uk/629069/

Usage rights: (cc) BY Creati

Creative Commons: Attribution 4.0

Additional Information: This is an Open Access article published in IEEE Transactions on Artificial Intelligence.

Enquiries:

If you have questions about this document, contact openresearch@mmu.ac.uk. Please include the URL of the record in e-space. If you believe that your, or a third party's rights have been compromised through this document please see our Take Down policy (available from https://www.mmu.ac.uk/library/using-the-library/policies-and-guidelines) 2

3

4

Building Trustworthy AI Solutions: A Case for Practical Solutions for Small Businesses

Keeley Crockett[®], *Senior Member, IEEE*, Edwin Colyer[®], *Senior Member, IEEE*, Luciano Gerber[®], and Annabel Latham[®], *Senior Member, IEEE*

5 Abstract—Building trustworthy artificial intelligence (AI) solutions, whether in academia or industry, must take into considera-6 tion a number of dimensions including legal, social, ethical, public 7 8 opinion, and environmental aspects. A plethora of guidelines, principles, and toolkits have been published globally, but have seen 9 limited grassroots implementation, especially among small- and 10 11 medium-sized enterprises (SMEs), mainly due to the lack of knowledge, skills, and resources. In this article, we report on qualitative 12 SME consultations over two events to establish their understanding 13 14 of both data and AI ethical principles and to identify the key barriers SMEs face in their adoption of ethical AI approaches. 15 We then use independent experts to review and code 77 published 16 toolkits designed to build and support ethical and responsible AI 17 practices, based on 33 evaluation criteria. The toolkits were evalu-18 ated considering their scope to address the identified SME barriers 19 20 to adoption, human-centric AI principles, AI life cycle stages, and key themes around responsible AI and practical usability. Toolkits 21 were ranked on the basis of criteria coverage and expert intercoder 22 23 agreement. Results show that there is not a one-size-fits-all toolkit 24 that addresses all criteria suitable for SMEs. Our findings show few exemplars of practical application, little guidance on how to 25 use/apply the toolkits, and very low uptake by SMEs. Our analysis 26 27 provides a mechanism for SMEs to select their own toolkits based 28 on their current capacity, resources, and ethical awareness levels -29 focusing initially at the conceptualization stage of the AI life cycle and then extending throughout. 30

31 Impact Statement-In parallel to the recent acceleration in development and adoption of artificial intelligence, there has been 32 intense and worldwide discourse around the ethics of such sys-33 34 tems. This debate has highlighted that without good governance, transparency and monitoring, indiscriminate use of AI could lead 35 36 to significant harms, discrimination, and injustice. Consensus has 37 settled on a broad set of overarching principles for ethical AI; now 38 myriad resources and toolkits exist to assist with embedding ethical practices along the research-development-deployment value chain. 39 40 Our evaluation of 77 toolkits reveals the breadth and depth of the themes they cover and barriers to their use, including a lack of 41 adoption case studies. We provide organizations, especially SMEs, 42 43 with an easy-to-use lookup table (Table V) to help them select a set of

Manuscript received July 30, 2021; revised October 6, 2021; accepted December 4, 2021. This article was recommended for publication by Associate Editor F. Chowdhury upon evaluation of the reviewers' comments. (*Corresponding author: Keeley Crockett.*)

Keeley Crockett, Luciano Gerber, and Annabel Latham are with the Department of Computing and Mathematics, Manchester Metropolitan University, M1 5GD Manchester, U.K. (e-mail: k.crockett@mmu.ac.uk; L.Gerber@mmu.ac.uk; A.Latham@mmu.ac.uk).

Edwin Colyer is with Research and Knowledge Exchange, Manchester Metropolitan University, M1 5GD Manchester, U.K. (e-mail: E.Colyer@mmu.ac.uk).

This article has supplementary downloadable material available at https://doi.org/10.1109/TAI.2021.3137091, provided by the authors.

Digital Object Identifier 10.1109/TAI.2021.3137091

toolkits to ensure that as well as addressing all key ethical themes,44they can also match their resources, skills and priority areas for45implementing ethical best practice.46

Index Terms—Artificial intelligence (AI), business, ethics, responsible, toolkits, trustworthy.

I. INTRODUCTION

HE ethical, social, and legal landscape of artificial intel-50 ligence (AI) driven systems is rapidly changing. Since 51 the General Data Protection Regulation 2018 [1], stakeholders 52 developing AI systems have faced numerous challenges in the 53 interpretation and implementation of Article 22, specifically 54 concerning an individual's rights in the context of automated 55 decision-making, the ability to explain AI decisions, explanation 56 of the logic involved, and the development of models using only 57 "correct" data. This has caused major challenges because of the 58 lack of legal guidance, case law, and ethical principles about the 59 use of AI in different contexts. For small- and medium-sized 60 enterprises (SMEs), these challenges are even greater due to a 61 lack of specific skills, budget, and human resource. The interna-62 tional policy and impact landscape of AI is still fragmented in 63 approaches to regulation, frameworks, guidelines, and standards 64 (i.e., P7000), with numerous ethical principles being circulated 65 which all convey broadly similar messages [2]-[15]. 66

These "guidelines" often focus on the AI technology or 67 service rather than organizational processes and human behav-68 iors, providing little to no mechanisms for accountability and 69 compliance (audit), and ignore the benefits of coproduction 70 and public scrutiny [16]. From an SME perspective, practical 71 implementation is difficult if not impossible. There has been 72 significant "bad press" around poor design, poor rationale, and 73 unethical applications of AI, which has fueled public mistrust. 74 Pownall [17] provides an excellent, regularly updated repository 75 of news stories that challenge whether the use of AI is ethical, for 76 example, the use of face tracking tablets which profile customers 77 and deliver relevant advertisements in UBER. As the public 78 gains knowledge and understanding of issues around the use 79 and application of AI (including bias, fairness, accountability, 80 responsibility, etc.) coupled with an increased awareness of data 81 privacy, both public services and the private sector will have to 82 become more accountable if they win public trust and secure the 83 vital public "license to operate." Reputational damage as a result 84 of insufficient or ineffective data and AI governance can cause 85 significant harm to a business, with greater impact on SMEs 86 [17]. There is still a significant gap between top-down theory 87

This work is licensed under a Creative Commons Attribution 4.0 License. For more information, see https://creativecommons.org/licenses/by/4.0/

47

48

and practical adoption of robust ethical practices across the entire
AI value chain [15], [18], [19], but our research suggests that
this is more prevalent in SMEs.

91 In this article, we adopt the European Commission's definition of an SME which is an enterprise with fewer than 250 92 employees, a turnover below €50 million or a balanced sheet 93 total below €43 million [20]. A small business has fewer than 94 50 employees and a micro business fewer than ten employees 95 [20]. Global business will have different definitions on the size 96 97 of SMEs, for example, in USA, an SME may have up to 500 employees dependent on the sector [21]. The World Bank states 98 that globally, SMEs represent 90% of businesses and account 99 for over 50% of employment, and in emerging markets, seven 100 out of ten jobs are created by SMEs [22]. In many countries, 101 SMEs are able to access competitive public funding to support 102 103 growth acceleration and drive innovation in the AI space, but to date there has been little to no focus on responsible innovation. 104 These programs have generally ignored the need for strong AI 105 106 and data governance, and not provided training and upskilling in the domains. Fortunately, over the past few years numerous 107 108 organizations and academics have published "ethical toolkits" to help organizations adopt and embed processes and practices 109 that mitigate risks and "do AI ethically." These toolkits help 110 organizations ensure their innovative systems adhere to the key 111 112 pillars of "ethical tech" around beneficence, nonmaleficence, autonomy, justice, and explicability [19]. 113

The overall aim of this article is to evaluate the thematic 114 and AI life cycle coverage of these toolkits. We also assess the 115 usability of the toolkits from an SME perspective and identify 116 117 which toolkits are least onerous to adopt and address the barriers to adoption highlighted by SMEs. By categorizing the toolkits 118 against ethical AI themes and adoption/usability, we provide 119 120 organizations of all sizes, but especially SMEs, with an easy way to identify the most suitable tools, methods, and processes 121 to implement. Our study is divided into two parts. First, we 122 conducted qualitative SME consultations over two events to es-123 tablish their understanding of both data and AI ethical principles 124 and to identify the key barriers SMEs face in their adoption. As 125 the collaboration between business and universities is a highly 126 important mechanism for R&D activities and for stimulating 127 innovation, it is important that academics make the good ethical 128 research practices from within their institutions integral to con-129 tract research and knowledge exchange activities. Second, we 130 conducted a review of available toolkits (published in academic, 131 organizational, government, and gray literature) that support 132 ethical and responsible AI practices. We evaluated these toolkits 133 using criteria partly informed by our SME consultations across 134 four aspects of ethical AI: 1) human-centric ethical principles; 135 2) applicability across the AI development life cycle; 3) barriers 136 to adoption; and 4) key ethics themes covered. 137

In this article, we define a toolkit as a document or resource including guidelines (provided the described methods, techniques,
or instructions for implementation), checklists, methodologies,
activities, processes, frameworks, workflows, or approaches
where the content focus is on responsible or ethical data (data
ethics) or AI (ethical/responsible/trustworthy/trusted AI). We
expand the definition of toolkit defined by Morley *et al.* [23]

which focuses only on technical toolkits designed for data 145 scientists and developers up to 2018. 146

This research aims to address the following research questions. 147

- 1) What are the barriers to ethical AI adoption by SMEs? 149
- 2) What is the current state of the market in practical toolkits
 for embedding AI ethical frameworks and governance into
 an SME culture?
 151

153

158

159

182

The main contributions of this article are as follows.

- An analysis of the viewpoints of SMEs on ethical data and AI practices established through two engagement events which are useful to those organizations which are developing toolkits.
- 2) Identification of barriers to adoption of ethical principles, practices, and toolkits for SMEs.
- 3) A review and evaluation of recent toolkits against four groups of criteria (common ethical principles, stages of the AI product life cycle, responsible AI aspects and practical application aspects) designed to facilitate practical application of data and AI ethical practices.
- 4) An easy-to-use lookup table of ranked toolkits based on expert intercoder agreements of criteria coverage suitable for SMEs to use.
 165
- 5) Recommendations to the research community on the role
 of data and AI ethics in business knowledge exchange.
 169

The rest of this article is organized as follows. Section II 170 presents a summary of the core risk factors associated with 171 AI and an overview of the latest legal frameworks and current 172 ethical guidelines and principles. In Section III, we present our 173 two-part methodology; first, describing two SME events leading 174 to the identification of barriers to adoption of ethical toolkits and 175 second, our method for conducting a review and coding of the 176 state-of-the-art toolkits against a range of criteria. We perform an 177 analysis of these toolkits and SME events in Section IV, which 178 leads to a series of recommendations, conclusions, and the wider 179 implications of findings in Section V. 180

II. BACKGROUND 181

A. Risk Factors in AI

When conceptualizing, creating, and implementing an AI 183 system, it is important to consider the risk factors associated with 184 the data used, the model(s) built, and the life span of the model 185 [18], [19]. Furthermore, the societal outcomes and impacts 186 (negative or positive; helpful or harmful) arising during the life 187 span of application should also be considered. From a business 188 perspective, there is a clear relationship between perceived risk 189 in an AI system in a given context and how much trust users have 190 in the decisions it makes [24], [25]. The majority of risk factors 191 are well documented. Bias is one of the most complex factors 192 as consideration must be given to bias that is embedded into 193 organizational or industrial cultures, personal, unconscious, and 194 human bias and data representation bias [26], [27]. For example, 195 data that have been labeled by humans for training a model may 196 be subjective, even among experts. Different models may need to 197 be developed for different genders, cultures, etc., as it is rarely 198 possible to generalize models to an entire human population 199

based on limited training data. Fairness is about treating people 200 equally through developing models that encapsulate moral stan-201 dards in the decision-making process. Explainability is required, 202 203 so all stakeholders, including people impacted by the decisions of automated systems, can understand how a decision is made 204 and the user knows why a system has made a decision [28], 205 [29]. Societal impacts (potential benefits and harms) must be 206 considered by a business, not only just to mitigate reputational 207 damage in case of legal complaints but also to meet or exceed 208 209 minimum standards of business ethics. Businesses must question where responsibility (tasks and obligations) lies within their AI 210 governance framework and define accountability (oversight and 211 liability) to roles across the design/development/ deployment 212 life cycle. With AI legislation changes on the horizon, deep 213 thinking and consensus surrounding these risk factors is required 214 by both academics and industry regardless of size to assess 215 the risk of an AI solution to both individuals and society. The 216 problem is now bridging the gap between principles and practice, 217 so there is some assurance that AI systems comply with the 218 agreed principles. 219

220 B. Principles and Guidelines

Over the past five years, governments, corporations, and inter-221 222 national bodies have produced a significant amount of guidance on the ethical dimensions of AI and data driven technologies. 223 To understand how crowded this space is and the difficultly 224 of choice for SMEs with regard to which guidelines to follow, 225 this section provides a brief overview. In 2019, Jobin et al. [4] 226 conducted a survey of global ethical guidelines comprised of 84 227 documents and analyzed their thematic coverage over 11 ethical 228 principles identified by keywords. This work provides a good 229 understanding of the coverage of ethical AI principles and guide-230 lines between 2011 and April 2019. However, the landscape 231 is very dynamic. In 2019, the Beijing Academy of Artificial 232 Intelligence published the Beijing AI Principles advocating eth-233 ical AI [5], OECD proposed five value-based principles for the 234 responsible stewardship of trustworthy AI [7], and the European 235 Commission issued ethical guidelines for Trustworthy AI [2]. In 236 2020, the U.S. Office of Management and Budget issued Guid-237 ance for Regulation of Artificial Intelligence Applications [11]. 238 In June 2021, The General Conference of the United Nations 239 Educational, Scientific and Cultural Organization (UNESCO) 240 presented the Draft Text of the Recommendation on the Ethics 241 of Artificial Intelligence, which focuses on a human-centered 242 approach to AI, recommending that "AI must be for the greater 243 interest of the people, not the other way around" [8]. The 244 U.K. government provided an updated summary of data and AI 245 ethical principles developed by both the public sector and the 246 government in 2020 [9], which included a joint publication on 247 AI procurement guidelines developed with the World Economic 248 Forum [30], and specific guidelines and a checklist for using AI 249 in health care [31]. In 2021, the U.K. AI Council published an AI 250 road map [32], further "guidance" on procurement [33] and its 251 national data strategy [34]. A brief analysis of the commonality 252 of ethical principles can be found as shown by Crockett [35], 253 254 from which a subset of our toolkit evaluation criteria is derived.

C. Legal Frameworks

Legal frameworks in the space of AI and data driven technolo-256 gies are relatively new and rapidly emerging. The GDPR 2018 257 [1] first introduced Article 22, a series of safeguards and infor-258 mation obligations in relation to automated decision-making. 259 These included empowering the data subject as stated in Recital 260 71 "not to be subject to a decision based solely on automated 261 processing, including profiling, which produces legal effects 262 concerning him or her or similarly significantly affects him or 263 *her*" [1], the right to ask for human intervention, explanation 264 of how the automated decision was made "the logic involved." 265 Recital 71 states that the data controller should use appropriate 266 mathematical and statistical procedures for profiling and that 267 data should be accurate in order to minimize the risk of errors [1]. 268 In 2018, the EU also published its AI strategy which promoted a 269 human-centric approach, which focused on respecting European 270 values and human rights. Recently, the EU has published the 271 proposed Regulatory Framework on AI [36], which contains 272 a framework to assess the risk of any AI product, service, or 273 system. Four risk levels are defined as follows. 274

- Unacceptable risk: AI systems considered a clear threat to the safety, livelihoods, and rights of people will be banned. 276
- High risk: AI systems identified as high risk (including law enforcement, credit scoring, and border control management) are subject to a deep risk assessment, mitigation strategy, high quality datasets, traceability, documentation, clear explainability protocols to the user, and a high level of robustness, security, and accuracy.
 277
 278
 279
 280
 281
 281
 282
 282
 283
 284
 284
 285
 286
 286
 287
 287
 288
 289
 281
 281
 281
 282
 282
 283
 284
 284
 285
 286
 286
 287
 287
 288
 288
 289
 281
 281
 282
 282
 284
 284
 285
 286
 286
 287
 287
 288
 288
 289
 281
 281
 281
 282
 282
 284
 284
 285
 286
 286
 287
 288
 288
 288
 288
 288
 288
 288
 289
 289
 281
 281
 282
 284
 284
 285
 286
 286
 286
 287
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
 288
- 3) *Limited risk*: This includes chatbots where human-machine transparency is a requirement.
- *Minimal risk*: This includes applications such as AIenabled video games or spam filters [36].
 286

An excellent primer on the principles and priorities required 287 for a legal framework can be found in [37], produced by the 288 Council of Europe's Ad Hoc Committee on Artificial Intelligence. Leslie *et al.* [37] also provide suggestions on options for a 290 legal framework and a mapping between substantive human and 291 legal rights and key obligations of AI developers when building 292 AI systems and services. 293

III. METHODOLOGY

294

283

284

This article comprises a two-part methodology. The first part 295 is an analysis of a series of practical SME engagement events. 296 These events took place between July 2020 and June 2021 and 297 were designed to capture the "SME voice" on their understand-298 ing of ethical AI, its practical implementation, awareness of eth-299 ical toolkits, and the barriers to adopting good ethical practices. 300 The aim of the analysis was to establish which themes associated 301 with ethical AI that SMEs are most aware of, and the perceived 302 barriers to ethical AI adoption. Part two is a review of a range 303 of practical toolkits designed to support the implementing into 304 practice of ethical AI principles. These toolkits were evaluated 305 and coded against the common themes and barriers from the 306 SME events and against a range of criteria relating to coverage 307 of the AI life cycle, and general ethics themes. 308

333

309 A. Part 1: SME Engagement and Consultation Study

This section outlines the methodologies for two distinct SME engagement events which explored the need for and barriers to ethical AI.

313 1) Event 1: Our Place Our Data: To understand the landscape for local businesses and local authorities in ethical AI 314 understanding and practice, a qualitative research study took 315 place in June and July 2020, comprising two roundtables and 316 follow-up interviews. The study was initiated by Manchester 317 Metropolitan University (MMU), designed in collaboration with 318 an independent think tank and with the support of the U.K.'s 319 All-Party Parliamentary Group on Data Analytics (APPGDA). 320 During the roundtables, participants were provided with an 321 overview of a proposed model for place-based support for ethical 322 AI to build a local ecosystem in which ethical and responsible AI 323 development could be nurtured and thrive. The theme for the first 324 roundtable (n = 20) was "Data and Public: Creating a data-driven 325 future for Greater Manchester" and sought to capture responses 326 to a series of key questions, which included the following. 327

- 1) How can the public be better engaged with policies aroundethical data use?
- What are the current challenges and shortcomings associated with ethical guidelines and principles for the use of data by public and private-sector bodies?
 - 3) What does an effective local data ecosystem looks like?

The second roundtable was at U.K. national level, featuring 334 not only local SMEs and Policy Makers but also Members of 335 Parliament and the House of Lords, and key national stakehold-336 337 ers such as the Centre for Data Ethics and Innovation (CDEI), Visa, British Standards Institute, and the Greater Manchester 338 Combined Authority (GMCA). The second roundtable (n =339 18) focused on how parliament and government could work to 340 develop local data strategies as part of a wider effort to make the 341 342 U.K. a world leader in ethical, data-driven technologies. It also 343 analyzed current links between central government, regulators, local and combined authorities, and industry, and considered 344 how those links could be developed over the coming years. The 345 discussion focused on how to develop place-based approaches 346 347 to data ethics; the role for regulators and government bodies; the feasibility of an "Ethical AI kitemark," which organizations 348 should lead on ethical AI policies at the national and regional 349 level; and what challenges exist with regard to bringing these 350 bodies together. 351

352 Following the roundtables (between August 2020 and March 2021), a series of supplementary follow-up interviews were 353 conducted by Policy Connect with selected participants to ex-354 plore some of the emergent themes in greater depth. Summary 355 reports from both roundtable events and the interviews were 356 produced by Policy Connect and cross-checked by this study's 357 358 authors (Crockett and Colyer) for accuracy, identified emergent themes, and indicators of agreement, disagreement, and consen-359 sus among participants. 360

2) Event 2. Greater Manchester AI Foundry: The Greater Manchester AI Foundry [41], with £3 million ERDF funding, is a three-year research and innovation project which commenced in July 2020. The aim of the Foundry is to increase SME performance by placing AI research and innovation at the center of business growth through practical knowledge transfer from AI 366 academic research into industry. SMEs go through two phases: 367 1) Phase 1 is a series of workshops on AI development from 368 a business perspective and 2) Phase 2 is a technical assist to 369 develop a prototype AI solution. The objective is that research 370 acts as a technology accelerator for new products and services 371 based on AI. Given the importance of the development of ethical 372 technology, a pilot workshop was given in early 2021 to the 373 first cohort of SME participants (n = 20) to enable SMEs to 374 gain an understanding of ethical, social, and legal perspectives 375 of AI and data privacy, and also to facilitate practical ethics 376 into the technical assists. The workshop was not intended to 377 provide any legal advice, rather it was designed to showcase 378 best practice in ethics and regulatory compliance. The first 379 workshop was positively received and a full workshop was 380 developed and embedded with a second cohort in June 2021. 381 In the full workshop, SMEs were actively encouraged to look at 382 the impact and assess the risks of their AI product or service in 383 light of the newly proposed EU regulation [36]. The workshops 384 introduced a variety of ethical toolkits and activities with SMEs 385 including datasheets for datasets [42], consequence scanning 386 [43], conducting a data privacy impact assessment [44], and 387 examining the risk to stakeholders of an AI recruitment tool 388 using padlet [45]. Feedback on adoption of potential tools and 389 barriers to use was obtained through Q and A and discussion 390 during and after the workshop. Workshop members were also 391 asked to complete a longitudinal ethical AI practice survey [46]. 392 Feedback was anonymized and collated and thematic coding 393 was undertaken to identify ethical concerns and barriers. 394

B. Part 2: Review of Practical "Ethical" Toolkits

Our review of toolkits covers academic, organizational, gov-396 ernment, and gray literature sources. The search strategy em-397 ployed the following primary keywords: (toolkit, resource, 398 guidelines, guidance, checklist, methodology, method, activity, 399 process, framework, workflow, approach); (ethical, responsible, 400 trustworthy, trusted, data, data ethics, tech ethics); and [artificial 401 intelligence (AI), machine learning (ML)]. Our toolkit dataset 402 was created by using the primary keywords to perform searches 403 on Google Scholar and Scopus and gray online literature on 404 Google from 2017 to July 5th, 2021. Our toolkit dataset was 405 also cross-checked with work published by Morley et al. [23] 406 and Moltzau [38], who produced a full typology of identified 407 methods and tools (up to mid-July 2019) which were limited 408 to helping developers, engineers, and designers of ML apply 409 ethics within their roles. In comparison, our review takes on a 410 more holistic view in analyzing toolkits that are also used to 411 initiate engagement with wider public stakeholders to explain 412 decisions and build trust. Inclusion criteria were documents 413 (checklists, guidelines, activities) including those published by 414 public and private sectors, governments, and international bodies 415 and the toolkit language was English. Exclusion criteria were 416 legal frameworks, opinion articles and speeches. Once a list of 417 toolkits that met the inclusion criteria was obtained (referred to 418 as the EAI toolkit dataset), each toolkit was evaluated and coded 419 independently by expert researchers in the field of AI and ethics 420

TABLE I GROUP B: COMMON ETHICAL PRINCIPLES

Criterion	Ethical Principal
No	
B_{I}	AI should not be used to harm or kill any human and respect
	human rights
B_2	AI must always be fair, unbiased and transparent in the
	decision-making process
B_3	AI systems and solutions should always operate within the
	law and have human accountability
B_4	Data Governance and Data Privacy should be incorporated
	into the AI life cycle
B_5	Humans should always know when they have interactions
	with an AI system
B_6	AI systems should be inclusive to all human-centered AI
	design
B_7	Appropriate levels of explainability should always be
	provided on AI decision making
B_{δ}	Humans must always be in the loop when an AI is making a
	decision that affect other humans
B_{9}	Humans responsible for designing, developing and operating
	AI systems should be competent in the skills and knowledge
	required
B_{10}	AI systems should be sustainable and work to benefit humans,
	the society and the environment

 B_{11} AI systems should be inclusive to all

 TABLE II

 GROUP C: STAGES OF THE AI PRODUCT LIFE CYCLE

Criterion	Criterion Name/Description
No	
C_I	Conceptualization : includes imagineering, defining aims, objectives, desiderata, cost/benefit of new AI products and services and conducting a risk assessment
C_2	Data Preparation and Exploration : e.g., collection, curation, feature engineering, cleaning, feature selection, and sampling
C_3	Model Building and Evaluation
C_4	Deployment and Monitoring

based on four groups of criteria, shown in Tables I–IV. For each
toolkit, its source (academic, organizational, business, and gray)
was recorded, along with publication year, whether it was open
source, and the country of origin.

425 Criteria in Group *E* were determined on the basis of the find426 ings of the two SME engagement events reported in Section IV
427 – analysis of SME engagement events.

A modified nominal group approach to coding was adopted 428 [39], [40]. The first round of coding involved three experts in 429 the fields of AI, ethics, and business engagement, independently 430 evaluating two-thirds of the EAI toolkit dataset with each toolkit 431 being evaluated by two experts initially. A structured spreadsheet 432 containing links to the toolkits and the 33 criteria for coding 433 was given to each expert to evaluate and code independently. 434 Each criterion was coded according to a three-point Likert scale 435 with values in (01, 2) indicating, respectively, weak, moderate, 436 and strong levels of support by a toolkit for a given criterion. 437 For example, if a toolkit strongly addressed B_{10} – AI systems 438 should be sustainable and work to benefit humans, the society, 439 and the environment - then it was scored as 2; if it moderately or 440 partially addressed that criterion, it was scored 1; and if support 441

TABLE III GROUP D: RESPONSIBLE AI THEMES

Criteria	Criterion Name/Description
No	
D_1	Robustness
D_2	Fairness (includes bias)
D_3	Transparency
D_4	Accountability
D_5	Explainability
D_6	Privacy
D_7	Safety
D_8	Impact (both positive and negative, on society)
D_{9}	Inclusivity of the toolkit (in general): incorporation of needs
	from stakeholders with different roles (e.g., managerial, data
	protection officer), motivations, technical expertise (e.g.,
	machine learning engineers, senior management), and cognitive
	equity (for example, that it was inclusive to people with varying
	levels of educational attainment)
D_{10}	Inclusivity w.r.t to General Public: as D_9 but, more
	specifically, the extent to which the conception of the toolkits
	included and offered consultation with the general public

TABLE IV GROUP E: PRACTICAL APPLICATION ASPECTS

Criteria	Criterion Name/Description
No	
E_I	Exemplars: case studies, examples of what-good-looks-like,
	among others.
E_2	Quick Read e.g. too-long-didn't-read; short, accessible,
	practical, quick-start type of guidance for application of the
	principles.
E_3	Stakeholders Inclusivity: does the toolkit address different
	types of stakeholders such as technical, managerial, and end
	user (e.g., customer)?
E_4	Feasibility of applying the toolkit with respect to a typical
	SME skillset.
E_5	Feasibility of applying the toolkit with respect to resources
	such as workload, personnel, budget at SMEs.
E_6	Recommendations of AI Techniques: e.g., does the toolkit
	make concrete recommendations for data management and
	machine learning methods?
E_7	Recommendations on Personnel Training
E_{s}	Evidence of adoption of the toolkit by an SME

for the criterion was largely or completely absent, then it was scored as 0.

The first round of independent coding revealed a 72% agree-444 ment across 33 criteria; 18% of criteria indicated that there was 445 a disagreement with one expert coding 0 and another scoring 446 1 or 2; in 10% of cases, both experts agreed that the toolkit 447 contained at least some evidence of the criteria, but the experts 448 disagreed on how much (scoring 1 or 2). When adopting a 449 percentage agreement approach [39] there is no agreed threshold 450 for consensus, and it is up to the researchers to judge what 451 represents acceptable agreement for a particular study. A second 452 round of independent expert coding was then instigated for 453 all toolkits where there was significant disagreement for any 454 criteria, defined as when one expert scored 0 and the other 455 expert either 1 or 2; these toolkits were fully coded by a third 456 expert in an attempt to establish majority agreement. The level 457 of agreement between the three experts was then recorded in a 458

442

structured spreadsheet for 77 toolkits. There was a good majority 459 agreement between the two experts for 89% of the 33 criteria 460 scored across the 77 toolkits. Experts were unable to reach a 461 462 majority agreement on all criteria across all toolkits in only 1% of cases. The most common disagreement between the coders was 463 on the interpretation of B_{10} – AI systems should be sustainable 464 and work to benefit humans, the society, and the environment (6 465 out of 77 toolkits) and on the toolkit coverage of C_4 – deployment 466 and monitoring (6 out of 77 toolkits). 467

468

IV. ANALYSIS AND DISCUSSION

A. Analysis of SME Engagement Events 469

Event 1: For event 1, analysis of the first roundtable revealed 470 that ethical and legal issues surrounding "data" and not "AI" 471 needed to be resolved first before the wider ethical aspects 472 of AI could be addressed. This was true for both public and 473 private sector organizations. The key themes emerging from the 474 roundtables were as follows: 475

- 1) ethical guidelines and principles should be simple and 476 477 flexible and should be much more than a checklist;
- 2) practical guidance on how to apply data and ethical AI 478 principles should be usable; 479
- 3) mechanisms were needed to support practical guidance 480 481 (training, resource support) in partnership with local authorities: 482
- 4) data-driven technology strategies should be developed in 483 partnership with all stakeholders; 484
- 5) SMEs should have access to "resource knowledge shar-485 ing" to make effective and ethical use of AI and ML. 486

487 The main output of the Event 1 study was a report Our Place, Our Data: Involving Local People in Data and AI-Based Recov-488 489 ery [47], which made five recommendations to the U.K. government, including that local authorities should work in partnership 490 with businesses (including SMEs) and academic institutions to 491 develop data-driven technology strategies to develop innovative 492 AI services and products which have citizen engagement at the 493 heart of the creation process. 494

Event 2: The analysis of Event 2 was based on Q and A 495 during the two cohort sessions and follow-ups in 1:1 virtual 496 meetings. SMEs referred to the following Information Commis-497 sioner's Office (ICO) guidance: What are the accountability and 498 governance implications of AI? [48], guidance on AI and data 499 protection [44], data protection impact assessments [44], what 500 do we need to do to ensure lawfulness, fairness, and transparency 501 in AI systems? [45], and how do we ensure individual rights in 502 our AI systems? [49]. They noted these documents as long and 503 complicated, and provided no practical advice or methods on 504 how to apply them. The key message was that toolkits/guidance 505 needed to be simpler. One SME data scientist stated that they 506 "did not know some of this existed" emphasizing the general lack 507 508 of awareness. SMEs thought that training or free consultancy was required to help them understand and apply legal guidance 509 in relation to AI and data. Three SMEs also thought that in 510 general, ICO guidance was "subject to interpretation." Positive 511 feedback was received about the use of consequence scanning 512 513 [43] as a useful way to think about harms and risks of a product at conceptualization, but in general SMEs said whether they would 514 be used in practice was based on whether they had available 515 resource. They had no strong opinion about the benefits of 516 involving the public, for example, as a stakeholder in an activity 517 such as consequence scanning. Despite growing consensus on 518 the benefits of public involvement to build trust in AI tech [50], 519 [51], SMEs indicated that they were not sure how to involve 520 the public and that the real benefits of consulting with the 521 public was not clear. Two SMEs suggested that successful case 522 studies would benefit them. The SMEs thought that the toolkits 523 presented were useful, but they needed time to learn how to use 524 them – not only just one-off training but also how to practically 525 apply them in their own business. 526

Summary: From these two events, the barriers to SMEs adopting toolkits were identified as follows.

527

528

530

543

- 1) Availability of resources to SMEs (people and time), cur-529 rent skills, and training requirements.
- 2) Skepticism about the benefits of public stakeholder in-531 volvement in the design of new products and services. 532
- 3) Lack of understanding around governance of responsi-533 bility and accountability regarding AI development and 534 implementation outcomes. 535
- 4) The lack of audit and compliance and legal frameworks. 536
- 5) Need for practical training and upskilling regarding ethics, 537 data and legal frameworks, and managing liabilities. 538
- 6) Challenges associated with communication with users -539 different language for different stakeholders. 540
- 7) Serious implications for a business in terms of liability. 541 What are the consequences of noncompliance? 542

B. Toolkit Analysis

Following the methodology described in Section III, a total of 544 77 toolkits were identified which met the inclusion criteria. 30 of 545 these toolkits were from 2021, while the earliest was from 2017. 546 A total of 51% of toolkits were from the US, 23% were from the 547 U.K. and there was representation from South America, China, 548 Denmark, Saudi Arabia, Germany, and Ireland, in addition to 549 three toolkits which were classed as global. The process for 550 analyzing toolkits can be defined as follows. 551

- 1) All toolkits were scored using the groups of criteria B552 to E (see Tables I to IV) according to a three-point Likert 553 scale with values in (0, 1, 2) indicating, respectively, weak, 554 moderate, and strong level of support by a toolkit for 555 a given criterion. As explained in Section III, these are 556 the combined scores from the interannotator coding and 557 agreement process. 558
- 2) For the analysis of the criteria, we derived an *n* by *m* matrix 559 *R* (see supplementary material), where *n* is the number of 560 toolkits (n = 77) and m is the number of criteria considered 561 (m = 33).562
- 3) Each cell in R contains one of (0, 1, 2, D), with D standing 563 for a disagreement among coders. 564
- 4) From R, we derive a mean score for a toolkit (i.e., a 565 row) or a criterion (i.e., a column) by taking the mean 566 of its empirical probability distribution (epdf) (excluding 567 disagreements). More specifically, let X be either a row or 568

a column in *M*, which is assumed to be a discrete random variable. Then, $epdf(X) = (p_0, p_1, p_2)$, where *pi* is the probability of the score *i* in (0, 1, 2).

Table V located in the appendix, displays the statistical summary of scores across the 77 toolkits, ranked on the basis of their coverage of criteria groups C, D, and E, where p_0 , p_1 , and p_2 are the values of the epdf, shown as percentages, of the Likert scores on the criteria, and *m* is the number of criteria assessed. Group B is not included in Table V as it considerably overlaps with responsible AI themes in Group D. We opted for the latter, given that it provides a more fine-grained analysis of tool coverage. For example, B_2 – AI must always be fair, unbiased, and transparent in the decision-making process – is covered by D_2 – fairness (including bias) and D_3 (transparency).

The top-ranking toolkit was Microsoft's Responsible Inno-vation: A Best Practices Toolkit [111]. While this toolkit was targeted at developers, it had a strong focus on identifying potential negative consequences of technology on humans. The toolkit features three elements. The first, judgment call - a game and team-based activity that explores all of Microsoft's AI principles [128] through scenario imagining where the aim is for participants to write product reviews for different stake-holders accessing the impact and harms. Harms modeling – a framework for product teams based on the four pillars of respon-sible innovation ("injuries, denial of consequential services, infringement on human rights, and erosion of democratic and societal structures"[111]) – is designed for teams to look at real world impacts of technology. Finally, community jury, defined as an adaptation of the citizen jury [111] brings together the product team and user stakeholders to discuss various product artifacts, deliberate and cocreate new technologies over a 2-3-h session. This toolkit had moderate to strong coverage across all criteria B, C, and D. However, it did not contain any exemplars E_1 , and had no training guides E_7 , which is a key requirement for SMEs. That said, its uniqueness is its ability to engage the public, seek consensus, and opinion, and it is forward-thinking in terms of providing practical guidance that is applicable to a wide range of businesses/organizations. Ranked second was the U.K. government's Data Ethics Framework Guidance, published in 2020, which focuses on responsible and ethical use of data in the public sector [114]. While the emphasis is on the public sector, the guidance is targeted at all stakeholders who use or interact with data, including policy makers and data scientists. Similar to [111], the emphasis is on defining and understanding the public benefit of any "data project" including human rights, understanding potential consequences, compliance with law and diversity in the development team. The toolkit provides a set of questions which are scored on a Likert scale based on clarity and understanding with respect to a specific project. The framework also covers algorithms and outputs in relation to AI and is applicable to all stages of the AI life cycle. This toolkit also did not provide any examples of practical application E_1 and is less inclusive in its approach by not involving wider publics as stakeholders E_{3} . The toolkit did not offer any specific training E_8 .

Table V also highlights the lowest ranking toolkits [70], [97], and [125], none of which provided strong evidence of coverage across any of the criteria. For example, Covington is a global



Fig. 1. Boxplot showing mean score distributions of independent expert ranked criteria over Likert scale [0, ..., 2].

law firm, based in USA. Its toolkit [125] claims to provide practical guidance for "the evolving regulatory landscape" with an emphasis on USA, U.K., and EU. The guidance is in the form of overviews, summaries of news articles, and a white paper with links to recent AI legislation articles and to the ICO/Alan Turing Explaining AI Decisions' toolkit [83]. On the basis of our findings across the two SME engagement events, SMEs requested more training in order to understand the implications of legal frameworks and this toolkit would be difficult for them to practically apply as it is more a means of monitoring evolving regulation and legislation.

Fig. 1 shows the distribution of mean scores by groups of criteria. For example, one can see that criteria E (the practical application aspects for SMEs) has the lowest median and overall coverage by the toolkits (each, represented as a data point). Each plot represents one toolkit. This confirms the largely consensus view arising from our two events that in spite of the existence of toolkits to support responsible and ethical AI, most still lack adequate instructions and training to facilitate adoption. Many require significant time and specialist skills for implementation due to their length

Analysis has shown that no single toolkit covers all criteria, as indicated in Table V ($p_0 > 0$ in all columns). Consequently, each set of criteria will now be analyzed independently to assess criterion coverage and highlight those toolkits with the highest ranked coverage. This will help SMEs to select toolkits that best align with their business culture and values, and the stage they are at in developing their own ethical policies and procedures.

1) Common Ethical Principles (Group B): Fig. 2 shows the toolkit coverage of the ethical principles $B_1, ..., B_{11}$. Clearly, $B_2 - AI$ must always be fair, unbiased, and transparent in the decision-making process receives the highest coverage across all toolkits. This is closely followed $B_3 - AI$ systems should always $B_2 - AI$ systems should always



Fig. 2. Ranking of Group *B* criteria on mean score.



Fig. 3. Ranking of Group C criteria on mean score.

operate within the law and have human accountability and B_4 660 - data governance and data privacy should be incorporated 661 into the AI life cycle. These findings align with predominant 662 global ethical principles [4]. Of least coverage was B_5 , humans 663 664 should always know when they have interactions with an AI system, which is only highlighted by toolkits [74], [116], [118], 665 [120], and [126] and B_8 – a human should always be in the 666 loop for automated decision-making, covered by [101], [112], 667 668 and [126]. Toolkit [126] (ranked 33 overall) stands out in this group. Titled "Application Guide for the Ethical Assessment 669 of AI for Actors within the Entrepreneurial Ecosystem," the 670 toolkit is an open source guide published by the Inter-America 671 Development Bank in May 2021. Its interdisciplinary approach 672 to ethical self-assessment covers all stages on the AI life cycle, 673 governance, and security with a focus on human involvement in 674 AI systems. The guide has a three-stage assessment to determine 675 the level of human involvement based on the impact that the 676 system has on a human's life. The toolkit helps organizations 677 define associated key performance indicators, risk mitigation, 678 and even develop emergency responses following analysis of all 679 conceivable scenarios. 680

681 2) Stages of AI Product Life Cycle (Group C): Fig. 3 shows 682 the toolkit coverage for the four stages of the AI life cycle: 1) 683 conceptualization C_1 ; 2) data preparation C_2 ; 3) exploration, 684 model building, and evaluation C_3 ; and 4) deployment and 685 monitoring C_4 . Analysis showed that toolkits were less likely



Fig. 4. Ranking of Group D criteria on mean score.

to cover the audit and compliance stage of the life cycle, com-686 pared to the other stages, presumably because few regulatory 687 frameworks or standards are yet approved. For example, to 688 date, out of the IEEE P7000 standards in development, only the 689 IEEE 7010-2020 – IEEE Recommended Practice for Assessing 690 the Impact of Autonomous and Intelligent Systems on Human 691 Well-Being [14] is available on subscription only. Only toolkits 692 [55], [56], [65], [70], [83], [85], [95], [101], [104], and [107] 693 covered the whole life cycle, but to varying degrees. Toolkits 694 [56] and [107] ranked, respectively, third and fifth overall against 695 all criteria (see Table II). Agile ethics for AI (HAI) [56] is a 696 Trello board which contains a series of boards covering scope, 697 data audit, training, analysis, feedback, calibrate (optimal AI for 698 increased uptake), augmentation (e.g., upskilling and training), 699 and "people and the environment" which addresses accountabil-700 ity in AI deployment. Each board contains a series of "TO DOs" 701 with specific resources, all available as open source. The World 702 Economic Forum's AI Procurement in a Box: Workbook [107] is 703 a lengthy tool kit (54 pages) that features a series of questions and 704 risk matrices and mapping tools covering the full AI life cycle. 705 It is intended for businesses seeking to procure AI solutions. It 706 also features a user manual with a strong emphasis on how to 707 define the public benefit of AI while assessing risks in the early 708 stages of conceptualization. The toolkit provides guidance on 709 how to address both the technical and ethical limitations of data, 710 clearly addressing the impact of bias. 711

3) Responsible AI Themes (Group D): Fig. 4. shows the 712 toolkit coverage for the responsible AI themes: Robustness D_1 , 713 fairness D_2 , transparency D_3 , accountability D_4 , explainability 714 D_5 , privacy D_6 , safety D_7 , impact D_8 , inclusivity of the toolkit 715 (in general) D_9 , and inclusivity w.r.t. general public inclusion as 716 a stakeholder D_{10} . Examination of Group D criteria allows for 717 more fine-grained analysis than within the more general ethical 718 principles (see Fig. 2) and we expected to see the similarity with 719 ethical principle B_2 and fairness D_2 with regard to coverage. 720 Ninety-five percent of all toolkits moderately or strongly ad-721 dressed the issue of fairness, with 88% also addressing the 722 impact of AI technology on society D_8 . Accountability D_4 , both 723 in terms of the processes of developing responsible technology 724



Fig. 5. Ranking of Group *E* criteria on mean score.

and the decision outcome, quality of the data and the model pro-725 duced, also had moderate to strong coverage in 89% of toolkits. 726 More than half (53%) of the toolkits failed to include the public 727 728 voice, in any codesign or coproduction process to seek their opin-729 ions (D_{10}) and only 62% of toolkits were moderately inclusive to the requirements and needs of a wide range of stakeholders (i.e., 730 data scientists, software developers, managers, CEOs) (D_9) . The 731 Action-Oriented AI Policy Toolkit for Technology Audits by 732 733 Community Advocates and Activists [122], Agile Ethics for AI (HAI) [56], the JUST AI reflection prototype [82], Microsoft's 734 - Responsible Innovation: A Best Practices Toolkit [111], U.K. 735 governments, Data Ethics Framework Guidance [114], and the 736 Royal Society – Democratizing decisions about technology 737 toolkit [120] were the only toolkits to have strong coverage of 738 public inclusivity embedded within the toolkit objectives. As 739 reported in Ouchchy et al. [129], public opinion is critical in the 740 acceptance and adoption of new technology. Other work [130] 741 has recommended that businesses including ethical value state-742 ments on trusted webpages; the inclusion of both ethicists and 743 the public in new technology discussions could avert negative 744 media responses and reputational damage to businesses. The im-745 portance of the role of the public stakeholder is also highlighted 746 in policy road maps [32] and proposed regulation [36]. 747

4) Practical Application Aspects (Group E): Fig. 5 displays 748 the ranked criteria in relation to different aspects regarding the 749 practical application of the toolkits. Only 27% of the toolkits 750 were coded as being equivalent to "quick start" guidance E_2 . 751 Sixty-nine percent of toolkits and their associated websites 752 provided no exemplars or case studies of how to practically 753 apply the toolkit; only 6% provided at least one example of 754 adoption E_1 . Coverage of stakeholders' inclusivity E_4 within 755 the toolkit was scored as weak (27%), moderate (56%), and 756 strong (17%). Analysis showed that toolkits were designed with 757 specific audiences in mind, for example, the technical commu-758 nity (data scientists, programmers, and data analysts) where the 759 focus was on criteria such as bias and fairness in both data 760 quality and model generation. There were few toolkits that had 761 end users and public inclusivity in mind, suggesting that the 762 trajectory of practical application of toolkits is behind emerging 763 legislation and wider discourse around building trust through 764

public involvement [120]. The feasibility of practical application 765 of toolkits w.r.t. to SME resources (workload, personnel, and 766 budgets) E_5 was ranked similar to E_4 . This indicated that SMEs 767 would have to make a moderate to high investment to apply 768 toolkits and embed ethical values and processes into business 769 operations. Eighty-three percent of toolkits provided no training 770 opportunities such as step-by-step instructions, user guides or 771 checklist on how to practically use the toolkit. A strong emphasis 772 on training E_7 could only be found in IEEE Ethical Aligned 773 Design [65] and The Royal Society – Democratizing decisions 774 about technology toolkit [120]. The following toolkits covered 775 some aspects of training: [56], [60], [70], [88], [99], [102], [104], 776 [107], [108], [114], and [120]. An observation was that toolkits 777 that were focused on the conceptualization stage of the AI life 778 cycle and/or had more stakeholder inclusivity included some 779 form of training. 780

Finally, evidence of adoption of a specific toolkit by SMEs' 781 E_8 was barely evident to nonexistent in 90% of toolkits. This 782 suggests that either toolkits have not been designed with SMEs in 783 mind, the barriers to practical application are too high, or toolkits 784 are simply not being evaluated and publicized through practical 785 use cases. Digital Catapult's Machine Intelligence for Business 786 [88] (ranked 24th in Table II) has published a short case study on 787 Loomi - an AI assistant which builds trust through ethical trans-788 parent design [129]. Loomi, also the name of the SME featured 789 in the case study, utilized Digital Catapult's ethics framework 790 to reposition "the product using ethics as a key differentiator." 791 IDEO's toolkit (ranked 16th in Table II) highlights the benefits 792 of human-centered design using its Design Kit [64] in a series 793 of humanitarian case studies. 794

Across the criteria in this category E_1, \dots, E_8 , DotEveryone's 795 Consequence Scanning toolkit [43], ranked 21st (Table II), 796 exhibited moderate to strong coverage of all criteria. This 797 open-source toolkit, developed in U.K., allows businesses and 798 organizations (regardless of size) to examine, debate, risk assess, 799 and mitigate the potential consequences of their product/service 800 on society, communities, and the environment. A manual is 801 provided (27 pages), with minimal resources required. The tool 802 is employed at the conceptualization stage, with all stakeholders 803 taking part, although public stakeholders are not specifically 804 mentioned (D_{10}) . A strong facilitator is needed which may be a 805 barrier for SMEs, but a session can last as little as 90 min. The 806 tool has been reportedly adopted by SalesforceUX [130] as a 807 way to bring design risks out into the open. 808

C. Discussion

This article has evaluated and analyzed 77 toolkits that cover 810 different aspects of the ML/AL life cycle and common ethical 811 principles, responsible AI themes, such as bias and fairness, 812 and degrees of practical application. Consequently, every or-813 ganization should be able to find one or more toolkits that fit 814 with their working practices, culture, and to complement their 815 organizational values. Although Table II ranked Microsoft's 816 Responsible innovation: A Best Practices Toolkit [111] as the 817 number one toolkit with regard to our criteria (C, D, and E), it still 818

has limitations in its practical application by SMEs. Therefore, 819 this research concludes that there is not a toolkit currently in 820 existence that overcomes all the barriers and fully meets all the 821 822 needs of SMEs identified in the analysis of the two SME engagement events. SMEs struggle with long, wordy, and technical 823 documents. They require case studies, clear compelling stories 824 of benefits, and step-by-step instruction manuals on how to use 825 and embed toolkits into operations (and how much time/cash it 826 will cost). 827

828 There was a good distribution across the toolkits of all the ethical principles (criteria B). Greatest coverage (mean of 1.64) 829 was the Data Ethics Impact Assessment (ranked 17th in Table II) 830 [91] which comprised a 16-page questionnaire designed for 831 organizations to integrate the assessment of data ethics and the 832 impacts of their AI on humans and society within their develop-833 834 ment and operational processes. The 56 questions cover aspects of transparency, equality, data governance, sustainability, ac-835 countability, and human-centered design and centered, drawing 836 837 on DataEthics.eu's principles of data ethics. In contrast, Nesta's Civic Al Toolkit [121], which focused on using AI and data to ad-838 839 dress climate crisis and the Online Ethics Canvas [127], had little to no coverage. Results concluded that few toolkits addressed all 840 11 principles, and none were considered to fully address all 11 841 by any expert coder. Therefore, organizations will probably need 842 843 to use more than one toolkit to get comprehensive coverage.

Detailed analysis in Section IV revealed that toolkits [55], 844 [56], [65], [70], [83], [85], [95], [101], [104], and [107] covered 845 the whole AI life cycle, but to varying degrees. Experts agreed 846 that 24% of toolkits did not cover audit and compliance and 847 848 this may be due to the current lack of AI legislation, regulation, 849 and ethics standards. However, the proposed EU Regulation on AI [132] is likely to have a significant impact on future toolkit 850 851 development, as it is being described by the Global Centre for Data Innovation as the "most restrictive regulation of AI" in the 852 world. The expert coders agreed that 80% of toolkits analyzed in 853 this study placed emphasis on getting things right the first time, 854 i.e., at the point of AI product or service conceptualization, and 855 can be seen as proactive in determining the consequences and 856 harms a potential product could have on humans and society. 857

Analysis across the responsible AI themes (criteria D) in-858 dicates that the vast majority of toolkits covered aspects of 859 fairness and the impact of AI. While these are core values 860 in developing ethical and responsible AI, SMEs do need to 861 ensure that they address all themes across the AI life cycle 862 through culture change, rather than becoming fixated on bias 863 and fairness to the detriment of other themes. It is unsurprising 864 that so few toolkits strongly emphasize the importance of citizen 865 representation in their toolkit application. Only 8% of all toolk-866 its strongly advocated the participation of citizens, with 53% 867 868 relying only on internal stakeholders to take part. An absence of public involvement, especially in the new AI product/service 869 870 conceptualization phase, leads to flaws in design thinking due to a lack of diversity and inclusivity, which leads to narrower 871 perspectives. Consequently, a great business idea, with no public 872 license to operate, can ultimately lead to reputational damage 873 and loss of revenue. For example, Deloitte reported that a lack 874 875 of inclusivity in the conceptualization stage of a smart city design 876 resulted in a negative impact as people in wheelchairs were

unable to access eye-level retina scanners that require the person 877 to be standing [133]. Section IV highlighted only six toolkits 878 featuring citizen inclusivity. SMEs urgently need to find ways to 879 engage and involve more diverse teams including people outside 880 of their organizations, such as the general public. Our SME 881 engagement events found that this activity is typically beyond 882 their resources and skillset; they also raised concerns about 883 intellectual property rights and trade secrets being disclosed. 884 Put simply, SMEs need support and advice on how to engage 885 effectively. The Community Jury proposed within Microsoft's 886 Responsible innovation: A best practices toolkit [111] is a good 887 example of citizen engagement in the AI life cycle. The caveat 888 is that it was designed by and for a large corporate and not 889 an SME. Setting up such a jury may be daunting and resource 890 intensive for an SME; we propose setting up city or regional 891 juries, focused on ethical AI tech, as part of collective approach, 892 where SMEs could present novel ideas and seek public opinion 893 on design solutions. Ultimately, SMEs should seek to cocreate 894 and codesign with citizens to build trust and obtain the public 895 license to operate, but this is a significant step change to current 896 operations. 897

Our analysis also highlighted the lack of exemplars or case 898 studies by those organizations who have developed the toolkits. 899 There was little evidence of adoption and virtually none involv-900 ing SMEs. This is not to say they haven't been involved, but 901 stories, outcomes, analyses, benefits, and outcomes are not in 902 the public domain. This is a key knowledge gap that should 903 be addressed to close the gap between ethical principles and 904 practice. Toolkit developers could produce publicly accessible 905 case studies to thoroughly document the journey and the impacts 906 of adopting ethical practices. This is crucial to lower resistance, 907 leverage investment, and gain the trust and attention of SMEs 908 to invest their limited resources in upskilling and training their 909 employees on AI ethics. 910

Guidance on how to train people to use the toolkits is another 911 significant challenge. Our analysis indicated that 83% of toolkits 912 did not provide any training material on how to practically 913 implement the tool within the organization. While the overall 914 majority of toolkits are open source and in the public domain, 915 some organizations did offer consultation opportunities for a fee 916 [113], [124], [125]. However, this is not enough, particularly for 917 SMEs, if they do not come with comprehensive training and 918 support materials. 919

It is important to note that many of these toolkits have been 920 designed for specific and narrow purposes, with no intention to 921 support all possible dimensions of ethical AI, not least because 922 many were produced while ethical frameworks were still under 923 development. For example, IBM's 360 Fairness tool [78] was 924 conceived to focus on evaluating bias and the fairness of algo-925 rithms, with no explicit regard for any assessment of eventual 926 outcomes from decisions supported by said algorithms. At the 927 other end of the spectrum, AINow's Algorithmic Impact Assess-928 ment toolkit [53] is "designed to support affected communities 929 and stakeholders as they seek to assess the claims made about 930 these systems, and to determine where – or if – their use is 931 acceptable." It is therefore good to bear in mind that SMEs 932 may need to deploy two or more toolkits to fully capture all 933 dimensions of ethical operations. 934

V. CONCLUSION

This research aimed to address two research questions as follows:

935

- 938 1) first, to understand the AI ethics landscape from the SME
 939 perspective (and uncover any existing barriers to adop 940 tion);
- 941 2) second, to evaluate and identify existing toolkits that are942 suitable for practical application by SMEs.

Two SME engagement events were conducted that identified 943 a number of common barriers to ethical AI adoption by SMEs 944 on the themes of: 1) resources (people and time); 2) practical 945 business-focused training and upskilling on ethical and respon-946 sible AI; 3) data and AI governance infrastructures; 4) citizen 947 engagement; 5) applicability of legal frameworks (data and AI) 948 and how to apply them; and 6) audit, compliance, and liability. 949 Next, a comprehensive review provided a picture of the current 950 state of the market in availability of toolkits for embedding AI 951 ethical frameworks and governance into an SME culture. Our 952 key findings are summarized as recommendations to both the 953 SME and academic communities. 954

There is no one-size-fits-all toolkit that provides guidance sufficient to cover all ethical principles and themes around responsible and ethical AI. Toolkits vary in their feasibility to implement. It is recommended that SMEs select toolkits based on their current capacity, resources, and ethical awareness levels – focusing initially at the conceptualization stage of the AI life cycle and then extending throughout.

Academics engaged in knowledge transfer projects with busi-962 963 nesses should also share good ethical practices, policies, procedures and approval templates from their universities. While 964 established processes governing research ethics are different, 965 for example, in terms of the data processed and controlled, 966 and differences in legal basis according to GDPR, they can 967 968 help inform the private sector and provide cross pollination of 969 good ethical practices. In this article, ethical AI toolkits have been analyzed from an SME perceptive; however, evaluation of 970 criteria B, C, and D provides a useful reference to the academic 971 community, who may wish to embed the use of toolkits into their 972 ethics approvals and evaluations of research projects. Finally, 973 this analysis contributes a useful teaching resource for courses 974 that include AI ethics and/or data and AI governance, to enable 975 future data scientists and analysts to operationalize practical data 976 and AI ethics within their future employment settings. 977

978 Our next step is to produce an easy online tool to help SMEs select the best toolkits to implement/inform practice based on 979 coverage, ease of implementation, and stage in their ethical AI 980 evolution as a company. Our proposed online selection tool 981 will be a curated database that will allow SMEs to provide 982 their own rating across different categories following a similar 983 984 methodology to ours in this article. They will also be able to propose and categorize new toolkits to add to the database as 985 and when they become available, given the high level of activity 986 in this domain. The tool will be cocreated with SMEs and citizen 987 stakeholders and be flexible to incorporate legislation changes 988 and provide a go-to resource kit. 989

APPENDIX

TABLE V TOOLKIT COVERAGE OF CRITERIA C, D, and E

п	0	Ref	m	\mathbf{p}_0	\mathbf{p}_1	\mathbf{p}_2	mean	rank
73	2 [111]	22	13	5	82	1.69	1
7:	5 Í	114	22	13	14	73	1.6	2
6		561	22	0	27	64	1.55	2
7	, , , ,	1161	22	2	27	60	1.55	3
	0 I 0 I	1071	22	14	18	68	1.54	4
6	8 [10/]	22	5	36	59	1.54	4
24	4	[70]	22	4	41	55	1.51	6
7.	3 [112]	22	13	32	55	1.42	7
70] 0	109]	22	18	23	59	1.41	8
1	8	[65]	22	18	23	59	1.41	8
1	3	601	22	5	50	45	1.4	10
5	8	001	22	4	50	41	1.7	10
	0 4 [112	22	4	55	41	1.37	11
14	4 I	113]	22	14	36	50	1.36	12
0	I I	102]	22	18	32	50	1.32	13
6	4 [104]	22	14	41	45	1.31	14
6] 0	101]	22	14	41	45	1.31	14
11	7	[64]	22	9	55	36	1.27	16
4	9	91]	22	19	36	45	1.26	17
4	1	841	22	10	26	45	1.26	17
	, , , ,	1191	20	1.7	30	40	1.20	10
		110	20	15	45	40	1.25	19
41	U	83]	22	31	14	55	1.24	20
9)	[43]	22	18	41	41	1.23	21
7	'	[57]	22	27	23	50	1.23	21
8	2 [120]	19	37	5	58	1.21	23
4	5	88]	21	9	62	29	1.2	24
3	8	811	22	23	36	41	1.18	25
25	8	741	22	10	45	26	1.17	26
2			22	19	45	30	1.17	20
3		[70]	22	19	45	36	1.17	26
8	1 [119]	22	19	45	36	1.17	26
53	5	[96]	22	19	45	36	1.17	26
2	3	[69]	22	19	45	36	1.17	26
30	6	[79]	21	28	29	43	1.15	31
2	6	721	22	27	32	41	1 14	32
8	-, 8, 1	126	22	22	22	45	1.17	22
4	7 1	1201	22	32	25	45	1.15	22
-	<u> </u>	201	22	23	41	30	1.13	34
8		58]	22	32	27	41	1.09	35
6	9 [108]	22	27	41	32	1.05	36
14	4	[61]	22	32	32	36	1.04	37
2	7	[73]	22	23	50	27	1.04	37
1	5	[62]	22	22	55	23	1.01	39
3	3	771	22	22	55	23	1.01	30
4	3	1861	22	10	55	10	1.01	41
		[00]	22	18	04	18	1	41
-	,	[55]	21	43	14	43	1	41
5.	2	[94]	22	32	36	32	1	41
2	2	[68]	22	18	64	18	1	41
5	7	[98]	22	32	36	32	1	41
5	0	[92]	22	28	45	27	0.99	46
2	5	[71]	22	28	45	27	0.99	46
5	3	[95]	22	36	32	32	0.96	48
8	5 1	1231	20	25	55	20	0.90	40
6	7 1	1021	20	20	35	20	0.95	49
-		1051	22	32	41	27	0.95	49
1	9 I	11/]	20	20	65	15	0.95	51
6	5 [105]	22	50	9	41	0.91	52
6	6 [106]	22	41	27	32	0.91	52
3	4	[78]	21	38	33	29	0.91	54
8	6 [124]	22	37	36	27	0.9	55
2	2	[52]	22	36	41	23	0.87	56
4	6	[89]	21	33	48	19	0.86	57
8	4 I	1221	22	41	26	22	0.00	59
1	9	[66]	22	-11	24	23	0.02	20
1	0	[00]	22	41	30	25	0.82	38
-	1.	1101	22	41	36	23	0.82	58
7	I I	110]	22	37	45	18	0.81	61
7	7 [115]	22	41	45	14	0.73	62
5	1	[93]	22	46	36	18	0.72	63
4	4	[87]	22	45	41	14	0.69	64
1	6	[63]	22	50	32	18	0.68	65
3	7	801	21	33	67	0	0.67	66
1	2	[42]	22	55	27	19	0.63	67
	-	[95]	22	55	2/	18	0.03	0/
4.	4 3 ·	102]	22	68	14	18	0.5	68
6	4	103]	22	68	14	18	0.5	68
5	9 [100]	22	68	18	14	0.46	70
3	5	[53]	22	59	36	5	0.46	71
2	9	[75]	22	64	27	9	0.45	72
3	9	[82]	22	64	27	9	0.45	72
8	3 I	121]	22	64	27	9	0.45	72
7	6 Î	1271	22	73	27	0	0.27	75
5	6	[97]	22	01	~ / 0	0	0.00	76
	7 1	1251	22	71 05	7	0	0.09	/0
	, I	[دم.	22	73	3	0	0.05	11

ACKNOWLEDGMENT

The authors would like to thank Policy Connect and the 992 APPGDA for their work in the inquiry that led to the Our Place 993 Our Data Report [47] and the open source communities that we 994

990

used to conduct the data processing, analysis, and visualization 995 such as Seaborn [133], Matplotlib [134], Pandas [135], and 996 Jupyter Lab. 997

REFERENCES

- 999 [1] European Commission, "General data protection regulation," Recital, vol. 71, pp. 119-114, 2018. [Online]. Available: https://gdpr-info.eu/ 1000
- [2] European Commission, "Ethics guidelines for trustworthy AI," 2019. 1001 1002 [Online]. Available: https://digital-strategy.ec.europa.eu/en/library/ 1003 ethics-guidelines-trustworthy-ai
- [3] United Nations, "A framework for ethical AI at the United Nations," 1005 2021. [Online]. Available: https://unite.un.org/sites/unite.un.org/files/ unite_paper_-_ethical_ai_at_the_un.pdf
- 1007 [4] A. Jobin, M. Ienca, and E. Vayena, "The global landscape of AI ethics guidelines," Nat. Mach. Intell., vol. 1, pp. 389–399, 2019. BAAI, "Beijing AI principles," 2019. [Online]. Available: https://www. 1008
- 1009 [5] 1010 baai.ac.cn/news/beijing-ai-principles-en.html
 - [6] R. Vought, "Regulation of artificial intelligence applications," 2020. [Online]. Available: https://www.whitehouse.gov/wp-content/uploads/ 2020/11/M-21-06.pdf
 - [7] OECD, "Principles on AI," 2019. [Online]. Available: https://www.oecd. org/going-digital/ai/principles/
- [8] UNESCO, "Draft text of the recommendation on the ethics of artificial 1016 1017 intelligence," UNESCO Digit. Library, 2021. [Online]. Available: https: //unesdoc.unesco.org/ark:/48223/pf0000377897 1018
- U.K.-Gov, "Data ethics and AI guidance landscape," 2020. [Online]. 1019 [9] 1020 Available: https://www.gov.uk/guidance/data-ethics-and-ai-guidance-1021 landscape
- 1022 [10] P. Cihon, M. J. Kleinaltenkamp, J. Schuett, and S. D. Baum, "AI CERTIFICATION: Advancing ethical practice by reducing infor-1023 1024 mation asymmetries," IEEE Trans. Technol. Soc., to be published, 1025 doi: 10.1109/TTS.2021.3077595.
 - [11] U. S. Government, "Guidance for regulation of artificial intelligence applications," 2020. [Online]. Available: https://www.whitehouse.gov/wpcontent/uploads/2020/01/Draft-OMB-Memo-on-Regulation-of-AI-1-7-19.pdf
- 1030 [12] Australia Government, "Australia's artificial intelligence ethics frame-1031 work," 2019. [Online]. Available: https://www.industry.gov.au/data-and-1032 publications/australias-artificial-intelligence-ethics-framework
- 1033 [13] IEEE, "The ethics certification program for autonomous and intelligent 1034 systems (ECPAIS)," 2020. [Online]. Available: https://standards.ieee. 1035 org/industry-connections/ecpais.html
- [14] IEEE, "Ethics in action in autonomous and intelligent systems," 1036 1037 IEEE P7000 Standards Projects, 2021. [Online]. Available: https:// 1038 ethicsinaction.ieee.org/p7000/
- D. Schiff, J. Borenstein, J. Biddle, and K. Laas, "AI ethics in the 1039 [15] 1040 public, private, and NGO sectors: A review of a global document col-1041 lection," IEEE Trans. Technol. Soc., vol. 2, no. 1, pp. 31-42, Mar. 2021, 1042 doi: 10.1109/TTS.2021.305212.
- 1043 [16] A. Kumar, B. Finley, B. T, S. Tarkoma, and P. Hui, "Sketching an AI 1044 marketplace: Tech, economic, and regulatory aspects," IEEE Access, vol. 9, pp. 13761-13774, 2021, doi: 10.1109/ACCESS.2021.3050929. 1045
- [17] C. A. Pownall, "Understanding the reputational risks of AI," 1046 1047 IAAIC Repository, 2021. [Online]. Available: https://docs.google.com/ 1048 spreadsheets/d/1Bn55B4xz21
- [18] European Commission, "SME definition," [Online]. Available: https:// 1049 ec.europa.eu/growth/smes/sme-definition_en 1050
- 1051 [19] North American Industry Classification System, "SME definition," US 1052 Census Bur, 2020. [Online]. Available: https://www.census.gov/eos/ 1053 www/naics/development_partners/devpartners.html
- 1054 [20] The World Bank, "Small and medium enterprises (SMEs) finance," 2021. 1055 [Online]. Available: https://www.worldbank.org/en/topic/smefinance
- 1056 [21] D. Leslie, "Understanding artificial intelligence ethics and safety," The Alan Turing Institute, 2019. [Online]. Available: 1057 1058 https://www.turing.ac.uk/sites/default/files/2019-06/understanding_ 1059 artificial_intelligence_ethics_and_safety.pdf
- [22] L. Floridi et al., "AI4People-An ethical framework for a good AI 1060 1061 society: Opportunities, risks, principles, and recommendations," Minds 1062 Mach., vol. 28, pp. 689-707, 2018.
- 1063 [23] J. Morley, L. Floridi, L. Kinsey, and L. A. Elhalal, "From what to how: 1064 An initial review of publicly available AI ethics tools, methods and research to translate principles into practices," Sci. Eng. Ethics, vol. 26, 1065 pp. 2141-2168, 2020. 1066
- [24] M. Astobiza, M. Toboso, M. Aparicio, and D. López, "AI ethics for 1067 1068 sustainable development goals," IEEE Technol. Soc. Mag., vol. 40, no. 2, pp. 66-71, Jun. 2021. 1069

- [25] T. Araujo, N. Helberger, S. Kruikemeier, and C. H. De Vreese, "AI 1070 we trust? Perceptions about automated decision-making by artificial 1071 intelligence," AI Soc., vol. 35, no. 3, pp. 611-623, 2020. 1072
- [26] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan, 1073 "A survey on bias and fairness in machine learning," ACM Comput. 1074 Surv., vol. 54, no. 6, 2021, [Online]. Available: https://doi.org/10.1145/ 1075 3457607 1076
- [27] Centre for Data Ethics and Innovation, Review into bias in algorithm 1077 decision-making, 2020. [Online]. Available: https://assets.publishing. 1078 service.gov.uk/government/uploads/system/uploads/attachment_data/ 1079 file/957259/Review_into_bias_in_algorithmic_decision-making.pdf 1080
- [28] A. Bibal, M. Lognoul, A. De Streel, and A. B. Frénay, "Legal require-1081 ments on explainability in machine learning," Artif. Intell. Law, vol. 29, 1082 no. 2, pp. 149–169, 2021. 1083
- [29] N. Burkart and M. F. Huber, "A survey on the explainability of supervised 1084 machine learning," J. Artif. Intell. Res., vol. 70, pp. 245-317, 2021. 1085
- Office for AI and World Economic Forum, "AI procurement guide-[30] 1086 lines," 2020. [Online]. Available: https://www.gov.uk/government/ 1087 publications/guidelines-for-ai-procurement 1088
- [31] U.K.-Gov Department of Health and Social Care, "A guide to good prac-1089 tice for digital and data-driven health technologies - updated 2021," 2021. 1090 [Online]. Available: https://www.gov.uk/government/publications/code-1091 of-conduct-for-data-driven-health-and-care-technology/initial-code-1092 of-conduct-for-data-driven-health-and-care-technology 1093
- [32] U.K. Government, "AI roadmap," 2021. [Online]. Available: https:// 1094 www.gov.uk/government/publications/ai-roadmap 1095
- [33] U.K. Government, "Public sector guidance: Ethics, transparency and 1096 accountability framework for automated decision-making," 2021. 1097 Available: https://www.gov.uk/government/publications/ [Online]. 1098 ethics-transparency-and-accountability-framework-for-automated-1099 decision-making/ethics-transparency-and-accountability-framework-1100 for-automated-decision-making 1101
- [34] U.K. Government, "Government response to the consultation on 1102 the National Data Strategy," 2021. [Online]. Available: https: 1103 //www.gov.uk/government/consultations/uk-national-data-strategy-1104 nds-consultation/outcome/government-response-to-the-consultation-1105 on-the-national-data-strategy 1106
- [35] K. Crockett, Adaptive Psychological Profiling from Non-Verbal Be-1107 haviour - "Why are Ethics Just not Enough to Build Trust?", 2021, 1108 Women in Computational Intelligence, A. Smith, Eds. New York, NY, 1109 USA: Springer, 2021. 1110
- [36] European Union, "Regulation of the European parliament and of the 1111 council, laying down harmonised rules on artificial intelligence (Artifi-1112 cial intelligence act) and amending certain union legislative acts, 2021. 1113 [Online]. Available: https://eur-lex.europa.eu/legal-content/EN/TXT/ 1114 ?uri=CELEX:52021PC0206 1115
- [37] D. Leslie, C. Burr, M. Aitken, J. Cowls, M. Katell, and M. Briggs, "AI, 1116 human rights, democracy and the rule of law: A primer prepared for the 1117 Council of Europe," The Alan Turing Institute, 2021. [Online]. Avail-1118 able: https://www.turing.ac.uk/research/publications/ai-human-rights-1119 democracy-and-rule-law-primer-prepared-council-europe 1120
- [38] Al. Moltzau, "A typology of AI ethics tools, methods and research," 1121 2019. [Online]. Available: https://towardsdatascience.com/a-typology-1122 of-ai-ethics-tools-methods-and-research-cacdea134503 1123
- [39] J. Saldaña, The Coding Manual For Qualitative Researchers, Thousand 1124 Oaks, CA, USA: Sage, 2021. 1125
- [40] K. M. MacQueen, E. McLellan-Lemal, K. Bartholow, and B. Milstein, 1126 "Teambased codebook development: Structure, process, and agreement," 1127 in Handbook For Team-Based Oualitative Research, G. Guest and K. M. 1128 MacQueen, Eds. Lanham, MD, USA: AltaMira Press, pp. 119-135, 2008. 1129
- [41] GM AI Foundry, 2021. [Online]. Available: https://gmaifoundry.ac.uk/ 1130 about/ 1131
- [42] T. Gebru et al., "Datasheets for datasets," 2018, arXiv:1803.09010. 1132 [Online]. Available: https://arxiv.org/abs/1803.09010 1133
- [43] DotEveryone, "Consequence scanning," 2018. [Online]. Available: https:// 1134 //doteveryone.org.uk/project/consequence-scanning/ 1136

1135

1137

1138

1139

- [44] ICO, "Data protection impact assessments," 2020. [Online]. Available: https://ico.org.uk/for-organisations/guide-to-data-protection/guideto-the-general-data-protection-regulation-gdpr/accountability-andgovernance/data-protection-impact-assessments/
- [45] ICO, "What do we need to do to ensure lawfulness, fairness, 1140 and transparency in AI systems?," 2020. [Online]. Available: 1141 https://ico.org.uk/for-organisations/guide-to-data-protection/key-1142 data-protection-themes/guidance-on-ai-and-data-protection/what-do-1143 we-need-to-do-to-ensure-lawfulness-fairness-and-transparency-in-ai-1144 systems/ 1145

998

1004

1006

1011 1012

1013

1014

1015

1026

1027

1028

- [46] MMU, "Ethical AI practice survey," 2021. [Online]. Available: https: 1146 1147 //mmu.onlinesurveys.ac.uk/ai-ethics-survey
- 1148 [47] Policy Connect, "Our place our data: Involving local people in data and AI based recovery," 2021. [Online]. Available: 1149 1150 https://www.policyconnect.org.uk/research/our-place-our-datainvolving-local-people-data-and-ai-based-recovery 1151
- ICO, "What are the accountability and governance implications of AI?," 1152 [48] 1153 2020. [Online]. Available: https://ico.org.uk/for-organisations/guide-to-1154 data-protection/key-data-protection-themes/guidance-on-ai-and-data-1155 protection/what-are-the-accountability-and-governance-implications-1156 of-ai/
- [49] ICO, "How do we ensure individual rights in our AI systems?," 1157 1158 2020. [Online]. Available: https://ico.org.uk/for-organisations/guide-1159 to-data-protection/key-data-protection-themes/guidance-on-aiand-data-protection/how-do-we-ensure-individual-rights-in-our-ai-1160 1161 systems/
- [50] N. Aoki, "The importance of the assurance that 'humans are still in the 1162 1163 decision loop' for public trust in artificial intelligence: Evidence from 1164 an online experiment," in Computers in Human Behavior. New York, 1165 NY, USA: Elsevier, 2021, Art. no. 106572. [Online]. Available: https: //doi.org/10.1016/j.chb.2020.106572 1166
- [51] A. Jacovi, A. Marasović, T. Miller, and Y. Goldberg, "Formalizing trust 1167 1168 in artificial intelligence: Prerequisites, causes and goals of human trust in AI," in Proc. 2021 ACM Conf. Fairness, Accountability, Transparency, 1169 1170 2021, pp. 624-635.
- [52] a3i, "The trust in AI framework," 2018. [Online]. Available: http://a3i. 1171 1172 ai/trust-in-ai
- 1173 [53] AI Now Institute, "Algorithmic accountability policy toolkit," 2018. [Online]. Available: https://ainowinstitute.org/aap-toolkit.pdf 1174
- 1175 [54] The Institute for Ethical AI & ML, "AI-RFX procurement framework," 1176 2018. [Online]. Available: https://github.com/EthicalML/XAI

1177

1178

1179

1203

1204

1207

1208

- [55] T. Arnold and M. Scheutz, "The 'big red button' is too late: An alternative model for the ethical evaluation of AI systems," Ethics Inf. Technol., vol. 2, no. 1, pp. 59-69, 2018, doi: 10.1007/s10676-018-9447-7.
- 1180 [56] HAI, "Agile ethics for AI trello board," 2021. [Online]. Available: https: 1181 //trello.com/b/SarLFYOd/agile-ethics-for-ai-hai
- 1182 [57] FAT/ML, "Principles for accountable algorithms and a social impact 1183 statement for algorithms," 2020. [Online]. Available: https://www.fatml. 1184 org/resources/principles-for-accountable-algorithms
- N. Diakopoulos, D. Trielli, and G. Lee, "Algorithm tips," 2018. [Online]. 1185 [58] Available: http://algorithmtips.org/ 1186
- Z. Epstein et al., "TuringBox: An experimental platform for the evaluation [59] 1187 1188 of AI systems," in Proc. 27th Int. Joint Conf. Artif. Intell. Demos, 2018, pp. 5826-5828. 1189
- 1190 [60] J. Glenn, "Futures wheel," 2018. [Online]. Available: http://ethicskit.org/ 1191 futures-wheel.html
- P. Hall and N. Gill, "An introduction to machine learning interpretability," 1192 [61] 1193 O'Reilly. 2019. [Online]. Available: https://www.h2o.ai/wpcontent/uploads/2019/08/An-Introduction-to-Machine-Learning-1194 1195 Interpretability-Second-Edition.pdf
- 1196 [62] Ethics Kit, "Ethics toolkit," 2021. [Online]. Available: http://ethicskit. 1197 org/tools.html
- 1198 The Data Nutrition Project. [Online]. Available: https://datanutrition.org/ [63]
- IDEO.ORG, "Design kit," [Online]. Available: https://www.designkit. [64] 1199 org/case-studies 1200
- 1201 [65] IEEE, "Ethically aligned design," 2018. [Online]. Available: https:// 1202 ethicsinaction.ieee.org/
 - OPAL, "Open algorithms," 2021. [Online]. Available: https://www. [66] opalproject.org/about-opal
- 1205 Moral Machines, 2018. [Online]. Available: https://www.moralmachine. [67] 1206 net/
 - [68] M. Mitchell et al., "Model cards for model reporting," in Proc. Conf. Fairness, Accountability, Transparency, 2019, pp. 220-229.
- 1209 [69] The Federation, "New economic impact model," 2019. [Online]. Avail-1210 able: http://ethicskit.org/downloads/economy-impact-model.pdf
- 1211 [70] ODI, "Data ethics canvas," 2021. [Online]. Available: https://theodi.org/ 1212 article/the-data-ethics-canvas-2021/
- C. Oxborough, E. Cameron, A. Rao, A. Birchall, A. Townsend, and C. 1213 [71] 1214 Westermann, "Explainable AI: Driving business value through greater 1215 understanding," Price Waterhouse Cooper, 2018. [Online]. Available: 1216 https://www.pwc.co.uk/audit-assurance/assets/explainable-ai.pdf
- D. Peters and R. A. Calvo, "Beyond principles: A process for responsible 1217 [72] 2019. [Online]. Available: https://medium.com/ethics-of-1218 tech." 1219 digital-experience/beyond-principles-a-process-for-responsible-tech-1220 aefc921f7317

- [73] D. Peters, R. A. Calvo, and R. M. Ryan, "Designing for motivation, 1221 engagement and wellbeing in digital experience," Front. Psychol., vol. 9, 1222 2018, Art. no. 797. 1223
- [74] D. Reisman, J. Schultz, K. Crawford, and M. Whittaker, "Algorithmic 1224 impact assessments: A practical framework for public agency account-1225 ability," AINow, 2018. [Online]. Available: https://ainowinstitute.org/ 1226 aiareport2018.pdf 1227
- [75] Responsible AI Licenses, 2021. [Online]. Available: https://www. 1228 1229 licenses.ai/
- [76] Royal Society and British Academy, "Data management and use: Governance in the 21st century," 2018. [Online]. Avail-1230 1231 able: https://royalsociety.org/-/media/policy/projects/data-governance/ 1232 data-management-governance.pdf
- B. C. Stahl and D. Wright, "Ethics and privacy in AI and big data: [77] 1234 Implementing responsible research and innovation.," IEEE Secur. Priv., 1235 vol. 16, no. 3, pp. 26-33, May/Jun. 2018. [Online]. Available: https: 1236 //doi.org/10.1109/MSP.2018.2701164 1237
- [78] IBM, "IBM 360 fairness," 2019. [Online]. Available: https://www.ibm. com/blogs/research/2019/08/ai-explainability-360/
- [79] EthicalML, "XAI library," 2018. [Online]. Available: https://github.com/ EthicalML/xai
- [80] W. Zhao, "Improving social responsibility of artificial intelligence by 1242 using ISO 26000," in Proc. Published Under License by IOP Publishing 1243 Ltd IOP Conf. Ser.: Mater. Sci. Eng., Vol. 428, 3rd Int. Conf. Automat., 1244 Control Robot. Eng., 2018, Art. no. 012049. [Online]. Available: https: 1245 //iopscience.iop.org/article/10.1088/1757-899X/428/1/012049 1246
- [81] M. Zook et al., "Ten simple rules for responsible big data research," 2017. 1247 [Online]. Available: https://collaborate.princeton.edu/en/publications/ 1248 ten-simple-rules-for-responsible-big-data-research1249 1250
- [82] Ada Lovelace Institute, "JUST AI reflection prototype," 2021. [Online]. Available: https://www.adalovelaceinstitute.org/project/justai-reflection-prototype/
- [83] Information Commissioners Office and Alan Turing Institute, "Guidance on explaining AI decisions," 2020. [Online]. Available: https://ico.org.uk/for-organisations/guide-to-data-protection/keydata-protection-themes/explaining-decisions-made-with-ai/
- [84] Unbiased, "AI4DM toolkits," 2020. [Online]. Available: http://proboscis. org.uk/6473/special-offer-on-unbias-ai4dm-toolkits/
- [85] OECD, "Public consultation on the OECD framework for classifying AI systems," 2021. [Online]. Available: https://oecd.ai/classification
- [86] E. Blasch, J. Sung, and T. Nguyen, "Multisource AI scorecard table for system evaluation," 2021, arXiv:2102.03985. [Online]. Available: https: //arxiv.org/abs/2102.03985
- [87] Microsoft, "Allovus design, fairlearn: A toolkit for assessing and improving fairness in AI," 2020. [Online]. Available: https://www.microsoft.com/en-us/research/uploads/prod/2020/05/ Fairlearn_WhitePaper-2020-09-22.pdf
- [88] Digital Catapult, "Machines for machine intelligence," 2021. [Online]. https://www.digicatapult.org.uk/for-startups/acceleration-Available: programmes/machine-intelligence-garage
- [89] ACLU Washington, "Algorithmic equity toolkit," 2019. [Online]. Available: https://www.aclu-wa.org/AEKit
- [90] AI Global, "Responsible AI design assistant beta," 2021. [Online]. Available: https://oproma.github.io/rai-trustindex/
- [91] Data ethics, "Data ethics impact assessment," 2021. [Online]. 1275 Available: https://dataethics.eu/wp-content/uploads/dataethics-impact-1276 assessment-2021.pdf
- [92] Deon, "An ethics checklist for data scientists," 2021. [Online]. Available: https://deon.drivendata.org/
- [93] M. Arnold et al., "FactSheets: Increasing trust in AI services through supplier's declarations of conformity," IBM J. Res. Develop., vol. 63, no. 4/5, pp. 6:1-6:13, Jul.-Sep. 2019, doi: 10.1147/JRD.2019.2942288.
- [94] Google, "Playing with AI fairness tool," 2021. [Online]. Available: https: //pair-code.github.io/what-if-tool/ai-fairness.html
- [95] Google, "Explainable AI beta tools and frameworks," 2021. [Online]. Available: https://cloud.google.com/explainable-ai/
- VDE Bertelsmann Stiftung, "From principles to practice an in-[96] 1287 terdisciplinary framework to operationalise AI ethics," 2020. [On-1288 line]. Available: https://www.bertelsmann-stiftung.de/fileadmin/files/ 1289 BSt/Publikationen/GrauePublikationen/WKIO_2020_final.pdf
- [97] B. Goefring, F. Rossi, and D. Zaharchuk, "Advancing AI ethics beyond 1291 compliance from principles to practice," IBM, 2020. [Online]. Available: 1292 https://www.ibm.com/downloads/cas/J2LAYLOZ 1294
- [98] ICO, "Guidance on AI and data protection," 2021. [Online]. Available: https://ico.org.uk/for-organisations/guide-to-data-protection/keydata-protection-themes/guidance-on-ai-and-data-protection/

1233

1238

1239

1240

1241

1251

1252

1253

1254 1255

1256

1257

1258

1259

1260

1261

1262

1263

1264

1265

1266

1267

1268

1269

1270

1271

1272

1273

1274

1277

1278

1279

1280

1281

1282

1283

1284

1285

1286

1290

1293

1295

- 1297 [99] EthicalOS, "Ethical OS starter checklist," 2021. [Online]. Available: 1298 https://ethicalos.org/
- 1299 [100] I. D. Raji et al., "Closing the AI accountability gap: Defining an end-to-1300 end framework for internal algorithmic auditing," in Proc. Conf. Fairness, 1301 Accountability, Transparency, 2020, pp. 33-44.
- [101] The Institute for Ethical AI and Machine Learning, "The AI-RFX pro-1302 curement framework," 2021. [Online]. Available: https://ethical.institute/ 1303 1304 rfx.html
- 1305 [102] Design Ethically, "A library of resources and toolkits to help you integrate 1306 ethical design into your practice," 2021. [Online]. Available: https://www. 1307 designethically.com/toolkit
- 1308 [103] M. Madaio, L. Stark, J. Vaughan, and H. Wallach, "Co-designing check-1309 lists to understand organizational challenges and opportunities around 1310 fairness in AI," in Proc. CHI Conf. Humans Factors Comput. Syst., 2020, 1311 pp. 1–14.
- [104] Price Water House Cooper, "Responsible AI diagnostic tool," 2021. [On-1312 1313 line]. Available: https://www.pwc.com/gx/en/issues/data-and-analytics/ 1314 artificial-intelligence/what-is-responsible-ai.html
- [105] Smart Dubai AI Systems, "Ethics self-assessment tool," 2021. [Online]. 1315 1316 Available: https://www.smartdubai.ae/self-assessment
- [106] Aequitas, "Bias and fairness audit toolkit," 2021. [Online]. Available: 1317 1318 https://github.com/dssg/aequitas
- 1319 [107] World Economic Forum, "AI procurement in a box: Workbook," 2020. [Online]. Available: http://www3.weforum.org/docs/WEF_AI_ 1320 1321 Procurement_in_a_Box_Workbook_2020.pdf
- [108] World Economic Forum, "Empowering AI leadership, an oversight 1322 toolkit for boards of directors," 2020. [Online]. Available: https://spark. 1323 1324 adobe.com/page/RsXNkZANwMLEf/
- [109] 510, "F.A.C.T score for responsible AI," 2020. [Online]. Available: https: 1325 1326 //www.510.global/f-a-c-t-score-for-responsible-ai/
- 1327 [110] Microsoft, "InterpretML a toolkit for understanding machine learning models," 2020. [Online]. Available: https://www.microsoft.com/en-us/ 1328 1329 research/uploads/prod/2020/05/InterpretML-Whitepaper.pdf
- [111] Microsoft, "Responsible innovation: A best practices toolkit," 2020. [On-1330 1331 line]. Available: https://docs.microsoft.com/en-us/azure/architecture/ 1332 guide/responsible-innovation/
- [112] S. Vallor, "An ethical toolkit for engineering/design practice, 1333 1334 2018. [Online]. Available: https://www.scu.edu/ethics-in-technologypractice/ethical-toolkit/ 1335
- 1336 [113] AI Ethics Lab, "Ethics training," 2019. [Online]. Available: https:// 1337 aiethicslab.com/training/
- U.K. Government, "Data ethics framework," 2018. [Online]. Available: 1338 [114] 1339 https://www.gov.uk/government/publications/data-ethics-framework
- [115] Linklaters, "A toolkit for artificial intelligence (AI) projects," 2021. 1340 1341 [Online]. Available: https://www.linklaters.com/en/insights/thought-1342 leadership/artificial-intelligence-toolkit/ethical-safe-legal---a-toolkit-1343 for-artificial-intelligence-projects
- [116] Rolls Royce, "The Aletheia framework," 2021. [Online]. Available: 1344 https://www.rolls-royce.com/sustainability/ethics-and-compliance/the-1345 1346 aletheia-framework.aspx
- 1347 [117] NetHope, "Artificial intelligence (AI) ethics for nonprofits toolkit," 2020. [Online]. Available: https://solutionscenter.nethope.org/artifici-1348 1349 intelligence-ethics-for-nonprofits-toolkit
- [118] Open Roboethics Institute, "AI ethics assessment toolkit," 2019. [On-1350 1351 line]. Available: https://openroboethics.org/ai-toolkit/
- 1352 [119] D. Anderson, J. Bongaguro, M. McKinney, A. Nicklin, and J. Wise-1353 man, "Ethics & algorithms toolkit: A risk management framework for governments (and other people too!)," 2018. [Online]. Available: https:// 1354 1355 //ethicstoolkit.ai/
- 1356 [120] RSA, "Democratizing decisions about technology: A toolkit," 2019. 1357 [Online]. Available: https://www.thersa.org/globalassets/reports/2019/ democratising-decisions-tech-report.pdf 1358
- 1359 [121] Nesta, "Civic AI toolkit," 2021. [Online]. Available: https://www.nesta. 1360 org.uk/toolkit/civicai/
- [122] P. M. Krafft et al., "An action-oriented AI policy toolkit for 1361 1362 technology audits by community advocates and activists," in Proc. ACM Conf. Fairness, Accountability, Transparency, 2021, 1363 pp. 772-781. 1364
- 1365 [123] Omidyar Network, "Ethical explorer," 2020. [Online]. Available: https: 1366 //ethicalexplorer.org/additional-resources-for-ethical-explorers/
- [124] IDEO, "AI & ethics: Collaborative activities for designers," 2019. 1367 1368 [Online]. Available: https://www.ideo.com/post/ai-ethics-collaborative-1369 activities-for-designers
- Covington, "Artificial intelligence toolkit," 2021. [Online]. Available: 1370 [125] 1371 https://www.cov.com/en/practices-and-industries/industries/artificial-1372 intelligence/toolkit

- [126] Inter-American Development Bank, "Ethical assessment of AI for 1373 actors within the entrepreneurial ecosystem," 2021. [Online]. Available: 1374 https://publications.iadb.org/publications/english/document/Ethical-1375 Assessment-of-AI-for-Actors-within-the-Entrepreneurial-Ecosystem-1376 Application-Guide.pdf 1377
- [127] EthicsCanvas.org, "Online ethics canvas," 2021. [Online]. Available: 1378 https://www.ethicscanvas.org/index.html
- [128] Microsoft, "Responsible AI principles," 2021. [Online]. Available: https://www.microsoft.com/en-us/ai/responsible-ai?activetab=pivot1% 3aprimaryr6
- [129] L. Ouchchy, A. Coin, and V. Dubljević, "AI in the headlines: The portrayal of the ethical issues of artificial intelligence in the media," AI Soc., vol. 35, pp. 927-936, 2020.
- [130] SalesforceUX, "How to run a consequence scanning workshop," 2020. [Online]. Available: https://medium.com/salesforce-ux/how-to-run-aconsequence-scanning-workshop-4b14792ea987
- [131] Digital Catapult, "Loomi: The artificial intelligence assistant," 2021. [On-1389 line]. Available: https://www.digicatapult.org.uk/for-startups/success-1390 stories/loomi 1391
- [132] Center for Data Innovation, "How much will the artificial intelligence act 1392 cost Europe?," 2021. [Online]. Available: https://www2.datainnovation. 1393 org/2021-aia-costs.pdf 1394
- [133] Y. Murphy, S. Garg, B. Sniderman, and T. Buckley, "Ethical tech-1395 nology use in the fourth industrial revolution," 2019. [Online]. 1396 Available: https://www2.deloitte.com/us/en/insights/focus/industry-4-1397 0/ethical-technology-use-fourth-industrial-revolution.html 1398
- [134] M. L. Waskom, "Seaborn: Statistical data visualization," J. Open Source 1399 Softw., vol. 6, no. 60, 2021, Art. no. 3021. 1400
- [135] J. D. Hunter, "Matplotlib: A 2D graphics environment," Comput. Sci. 1401 Eng., vol. 9, no. 03, pp. 90-95, 2007. 1402 1403
- [136] W. McKinney, "Data structures for statistical computing in python," in Proc. 9th Python Sci. Conf., 2010, pp. 51-56.



Keeley Crockett (Senior Member, IEEE) has over 1405 20 years' experience in research and development 1406 1407 of computational intelligence algorithms and applications, including adaptive psychological profiling, 1408 fuzzy systems, dialogue systems, and educational 1409 tutoring systems. She is currently a Professor of com-1410 putational intelligence with Manchester Metropoli-1411 tan University, Manchester, U.K., and Coacademic 1412 lead with the ERDF-funded Greater Manchester AI 1413 Foundry. 1414

Prof. Crockett is the current Chair of the IEEE 1415 Taskforce on Ethical and Social Implications of Computational Intelligence and 1416 IEEE Women in Computational Intelligence and a STEM Ambassador. 1417





Edwin Colver (Senior Member, IEEE) is an Im-



science with Manchester Metropolitan University, 1429 Manchester, U.K. He is currently working on de-1430 veloping AI solutions for SMEs within the ERDF-1431 funded Greater Manchester AI Foundry. His research 1432 interests include fundamental and domain-agnostic 1433 applied data science and machine learning. 1434 1435

Annabel Latham (Senior Member, IEEE) is a Se-1436 nior Lecturer with the Department of Computing 1437 and Mathematics, Manchester Metropolitan Univer-1438 sity, Manchester, U.K. Her research interests include 1439 conversational agents, intelligent tutoring systems, 1440 affective computing, and ethics of AI in education, 1441 user profiling, and computational intelligence. 1442

Dr. Annabel is the Chair of IEEE U.K. and Ire-1443 land Women in Engineering, Chair of the IEEE CIS 1444 1445 Education Repository subcommittee, and a STEM 1446 Ambassador. 1447

1379 1380 1381

1382 1383

1384 1385

1386 1387

1388

1404

1419