

Navigating the Impossibility of Universally Ethical AI - a Practical Guide for Researchers

Anonymous Full Paper
Submission

001 Abstract

002 The area of ethics is full of trade-offs, but this is
003 sometimes ignored in Machine Learning research. I
004 argue that there can not be a single way to create
005 ethical AI, agreed upon by humanity at large, but
006 that this can be handled by considering a broader
007 range of stakeholders and clearly stating assump-
008 tions, values, and trade-offs. I begin this paper with
009 a survey of different current approaches to AI ethics
010 and some of their drawbacks. I go on to highlight
011 the impossibility of a universally ethical AI based
012 on the inherent contradictions in plural human val-
013 ues and directly contradicting definitions, as well
014 as to provide some discussion on the authoritarian
015 connotations of aims to find a singular true morality.
016 Then, some approaches to handling this impossibil-
017 ity through a democratisation of the problem and
018 clear communication are proposed. Lastly, a prac-
019 tical guide to navigating moral issues in machine
020 learning research – based on dialogue, transparency,
021 and conscious trade-offs – is proposed in the form
022 of four basic principles and a checklist.

023 1 Introduction

024 Ethical issues in the development and deployment of
025 AI tools have been increasingly researched in recent
026 years, with specialised journals [1, 2], dedicated con-
027 ferences [3–5], and increasing work on national and
028 international guidelines for ethical (or trustworthy)
029 AI [6–8]. The increased interest is a positive sign
030 that the issues are taken seriously, but most works
031 are either heavily theoretical, suggests a singular
032 solution, or are clearly aimed at the implementation
033 and large scale developments at companies (e.g. the
034 EU guidelines [6]). By this, there tend to be a gap
035 between the philosophy of AI ethics, and the imple-
036 mentation of it [9]. The aim of this paper is not to
037 solve the issue, but to begin bridging the gap by
038 presenting a variety of approaches from philosophy
039 to a more technical audience, together with guidance
040 on how these can be considered in ML research.

041 A complication in applying methods for ethical
042 AI, is that it is impossible to come to a universal
043 agreement on what constitutes ethical AI. The mul-
044 titude of ethical theories and definitions of values,
045 and the complexity of real world issues, makes trade-
046 offs both necessary and non-trivial. This is an issue

present in most ethical considerations, but in an AI
setting the myth of neutral tech can easily mask the
issue, if one is not careful. This impossibility, and
the necessity to acknowledge it, is my first thesis of
the paper.

However, acknowledging the impossibility of uni-
versally ethical AI is not enough in and of itself,
although it is an important piece in accomplishing
responsible AI practices. Even if no single approach
will be perfect by anyone’s standard, there are ways
to make better decisions. Or rather, to let others
make better decisions. I argue that the responsible
researcher need to consider a multitude of different
views on the ethical issues in their research field,
and take responsibility for and clearly communicate
the prioritisations and trade-offs decided upon. This
is my second thesis of the paper.

The paper begins by a survey of a variety of es-
tablished approaches AI ethics. These are presented
together with some examples on how they can be
used in AI development, and some motivations and
common critiques. Thereafter, I argue for the im-
possibility of finding a singular solution: the first
thesis.

Further, I present some democracy based ap-
proaches to AI ethics, and two frameworks for anal-
ysis: ACROCPoLis and the EU trustworthy AI
assessment list. Based on these, I argue that an
important element of AI ethics is to make conscious
trade-offs and communicate these clearly to allow for
wider discussions both within the AI community and
among the larger public affected by AI implemen-
tation. This is my second thesis of the paper. This
is followed by a practical guide for AI researchers,
with some basic principles and a checklist based on
the previous arguments of the paper. I conclude by
a summarising discussion.

084 2 Two main approaches and 085 their critiques

Many works on ethical AI rely on a utilitarian ap-
proach of maximising a utility, or minimising a cost
[10], and including different aspects and stakehold-
ers in this formulation [11, 12]. Another prevalent
approach is to formulate principles for AI develop-
ers to lean against, such as justice, privacy, and
explainability [6, 7]. However, according to works

093 of moral philosophy and philosophy of AI, there are
094 several issues with these approaches, making them
095 unsuitable as universal solutions.

096 2.1 Utilitarianism, unfairness, and 097 immeasurability

098 Utilitarianism is one of the most famous and well-
099 accepted, ethical frameworks in the field of Ethics.
100 It is a version of consequentialism, where the con-
101 sequences are the only thing determining the righ-
102 teousness of an action, with the specific goal of max-
103 imising overall utility [10]. Although utilitarianism
104 is a well-established framework, nothing is prevent-
105 ing the “utility” from being unfairly distributed.
106 A well-known example is the Utilitarian monster [13],
107 a thought experiment where one individual experi-
108 ences significantly more pleasure from resources
109 than anyone else. There is no single person feeling
110 the same amount of pleasure from a piece of food
111 as this utilitarian monster, no person appreciating
112 a cent or a loving word as much. In a utilitarian
113 decision, any ethical action have to be providing
114 for his monster, since this will always maximise the
115 total pleasure of the word.

116 Another critique can be found in the novella *Those*
117 *who walk away from Omelas* by Ursula K. Le Guin
118 [14]. Here we get introduced to the perfect society
119 of Omelas in the middle of a festival. Everyone is
120 happy and content. There are no wars nor soldiers,
121 but plenty of food and love. If you can think of
122 anything to make it better, more pleasurable, it is
123 already there. Not a single person feels pain - or no,
124 a single child feels pain. They are locked in a cellar
125 with just enough food to eat, without any human
126 interaction and living in their own excrement. They
127 are miserable, but have given up screaming for help.
128 And they are the guarantee for the happiness for the
129 rest, if their misery were to be relieved the perfect
130 society would fall. So people accept the necessity
131 and let them be. A utilitarian approach.

132 Another problem with the utilitarian approach
133 lies in deciding whose utility counts. The pleasure
134 of humanity is probably the standard utilitarianism
135 answer, but it is not unproblematic. Varying defini-
136 tions on who counts as part of humanity have been
137 used to discriminate people for a large part of our
138 history, which might make this an unsuitable term.
139 Even without these connotations, there are other
140 complicating aspects.

141 In a machine learning algorithm, it is often easi-
142 est to only include people directly affected by the
143 algorithm; the people whose utility can not be
144 changed by the system do not need to be considered.
145 Nonetheless, there are people affected indirectly
146 that are easily forgotten. For example, ignoring the
147 working conditions of miners and data-labellers and
148 the environmental impact of large-scale computing

makes the development easier. 149

Further, the utility of each actor is not necessarily
150 weighted equally. This can be a deliberate choice,
151 but is often an effect of the accuracy of predicted
152 outcomes being biased due to data access and his-
153 torical discrimination[15]. There is a lot of research
154 currently going into the latter problem [16], but the
155 truth is that there will always be a trade-off between
156 fairness and accuracy on biased training data [15]. It
157 is not obvious how these should be handled properly
158 [15, 17]. 159

Lastly, putting numbers on pleasure is not a sim-
160 ple task. In healthcare, measures such as Quality
161 Adjusted Life Year are used to make decisions about
162 treatment prioritisation but have also received a fair
163 amount of critique, for example, due to treating
164 the life of people with disabilities as less valuable
165 when the years are quality-adjusted [18]. Similarly,
166 it might be hard to put comparable numbers on
167 private data and discovering new knowledge. 168

Further, in many commercial machine-learning
169 settings the utility-measures chosen and trade-offs
170 done are hidden in large code bases and their intri-
171 cacies treated as trade secrets, discouraging public
172 discourse and promoting the idea of neutral tech.
173 When researchers don’t present their decisions on
174 ethical issues during development, similar issues
175 arise. 176

177 2.2 Guidelines, minimum efforts, and 178 loopholes

The approach of guidelines has its own problems.
179 They are often hard, or even impossible to implement
180 in practice, and enable minimal-effort solutions. 181

Due to vague formulations, there is often no spec-
182 ified minimum requirements for the different values.
183 The concept of privacy is a common value, lifted in
184 both the EU and UNESCO guidelines [6, 7], however
185 the level of privacy is not explicitly stated. Maximis-
186 ing privacy would mean to never use any personal
187 data in AI algorithms, but this is not necessarily
188 the goal. There may also be both minimal legal
189 requirements and ethical principles to follow about
190 the same value. For example, the EU have clear
191 laws, the GDPR [19], on how private information
192 can be used, but still include “Privacy and data
193 governance” as a principle under ethical AI in the
194 guidelines for trustworthy AI [6], not only as part
195 of the legality of the AI. How these values are in-
196 terpreted also affect research directions; currently
197 researchers develop ways to share anonymous data
198 between data owners through methods such as fed-
199 erated learning [20], and differential privacy [21],
200 which has a clearer connection to GDPR rules on
201 sharing data with third parties, than to existential
202 questions on how sharing personal data with *any*
203 data owner affects ones sense of integrity (for the
204

205 latter view, see e.g. [22]).

206 Another limitation on the concept of guidelines
207 is found when asking “who writes them, and for
208 whom?”. Both the EU and UNESCO guidelines are
209 written by teams of AI experts, which is highlighted
210 as a strength. In some ways it is, but it is also impor-
211 tant to consider which voices may have been left out
212 when asking highly educated people in a technical
213 field to make the ethical considerations. Secondly,
214 many guidelines seem directed mainly towards com-
215 panies with large scale projects. As research plays
216 an important role in the development of new meth-
217 ods, one might consider the lack of focus on more
218 basic research in the area to be an issue.

219 Finally, the successful usage of guidelines requires
220 the actors to actually aim at being fair and ethical,
221 since otherwise it can lead to fairwashing; the con-
222 cept of making a model seem fair or ethical when it
223 is not. A clear example of this is LaundryML [23],
224 developed to give a fair explanation of any decisions
225 from a black box algorithm, even when the original
226 decision is made on biased data, e.g. salary decisions
227 based on gender being explained with reference to
228 factors such as higher education and marital status.
229 It is probable that other similar loopholes exist.

230 3 Alternative approaches

231 Although the utilitarian and principle-based ap-
232 proaches to AI ethics are common, there are works
233 lifting other ethical considerations in this area.
234 These are not necessarily theories about making
235 the right decisions, but rather frameworks for taking
236 other aspects into account in the ethical analysis.

237 3.1 Power, relations, and care

238 Lacking in utilitarian and guideline based ap-
239 proaches is the study of power and relationships.
240 One could argue that these are not relevant, but
241 just as well that they are. Take the example of
242 the child in Omelas again. A fully utilitarian ap-
243 proach would suggest this is okay since it maximises
244 overall happiness, and that the situation would be
245 the same if the suffering was not done by a pow-
246 erless kid, but by a leader who took the suffering
247 upon themselves to further the prosperity of the peo-
248 ple. However, the power dynamics in these two
249 situations are completely different, and many would
250 argue that suffering pain to help overall happiness
251 is strictly more moral than forcing the same amount
252 of pain upon someone else for the same end. In
253 the AI community, there are several relevant power
254 dynamics between companies, users, data labellers,
255 researchers, mineral miners, and governments, that
256 can affect what we deem ethical. The issues of power
257 dynamics and interpersonal relationships have been

258 thoroughly discussed in various branches of feminist
259 ethics; a rich subfield of ethics philosophy.

260 One theory more comparable to utilitarianism is
261 care ethics, generally seen as a type of feminist frame-
262 work, which says that ethical decisions should be
263 based on care and empathy. There are a few papers
264 in which this is proposed as an alternative frame-
265 work for AI ethics. In [24], care ethics is proposed as
266 a way to handle the moral distance that automated
267 decisions may lead to. When decisions about people
268 are made with little human contact, care has to be
269 explicitly encoded in the moral judgements as they
270 will not naturally arise from the human connections
271 (as they are not present). By incorporating the
272 notion of care and considering the context of each
273 decision, they argue that more stakeholders will be
274 taken serious in the implementation process. In [25]
275 ethics of care is viewed as the opposite to oppres-
276 sive AI development strategies based on capitalism,
277 and power seeking agents. Care is also argued to be
278 important to achieve explainable and transparent
279 AI accessible for more diverse groups [26], and to
280 understand existential aspects of incorporating AI
281 in day to day life [22].

282 3.2 Climate, resources, and nature

283 Environmental ethics is a broad field of research, but
284 the common theme is that nature should be consid-
285 ered in moral judgements. In shallow environmental
286 ethics the value of the environment comes mainly
287 from in which ways it can help human societies to
288 flourish, while the subfield of deep environmental
289 ethics arises from the inherent value of the environ-
290 ment itself and anything living in it [27]. Similar to
291 feminist ethics, environmental ethics is more about
292 what issues to highlight and value, than about ex-
293 actly how to make the final decisions.

294 As large-scale computing requires a lot of power
295 [28] and many natural resources [29], environmental
296 ethics can be directly applied by considering these
297 aspects in more depth when considering the ethical-
298 ity of ML models. There are several works arguing
299 for more significant incorporation of these values,
300 either by calls for environmental ethicists to take a
301 larger interest in AI [30], by highlighting trade-offs
302 between ecological and human centred values [31],
303 and by encouraging clearer measurements of the
304 climate impact of model training [28].

305 3.3 Incorporation of non-WEIRD val- 306 ues

307 The acronym WEIRD stands for Western, Educated,
308 Industrialised, Rich, Democratic, and is used to
309 describe cultures fulfilling these criteria. Current
310 work on aligning AI with human values have received
311 critique for only focusing on WEIRD values, even

312 though these are in a minority (based on population
313 size). For example, it has been shown that ChatGPT
314 have both values and ways of thinking significantly
315 more closely related to countries such as the United
316 States and Sweden, than to less WEIRD countries
317 such as Ethiopia and Libya [32].

318 There have also been works aiming to incorpo-
319 rate non-WEIRD values into the discourse. The
320 Ubuntu philosophy has been suggested as a source
321 of inspiration, with its clear focus on helping the
322 group and valuing its social relationships [33]. The
323 concept of *honour* is common in many different non-
324 WEIRD societies, and has been lifted as a way to
325 view AI ethics in a more global lens. This could be
326 done through considering the preservation of honour,
327 and viewing content moderation more in terms of a
328 guardian protecting the users, rather than a form of
329 censorship, regardless of which values are prioritised
330 [34]. However, grouping such diverse values into a
331 single category of non-WEIRD ethics does have its
332 own problems.

333 4 The inherent contradictions

334 As shown in the previous sections, there are many
335 ways to conceptualise ethical AI, each with its own
336 critiques. The critiques for specific views is not the
337 only complication of a universally ethical AI though.
338 In this section I argue that a universal approach is
339 inherently unethical; my first thesis.

340 4.1 Trade-offs and contradicting val- 341 ues

342 In AI ethics it is common to pick some sort of values
343 to be considered. In the approach of guidelines they
344 are often stated as explicit values, but in a utilitarian
345 approach one can see from Section 2 that decisions
346 have to be made on how things are valued, in terms
347 of how things can be counted as utility. When
348 several things to value are picked however, they may
349 contradict each other. One example is the common
350 values of privacy and explainability, present together
351 in for example the EU guidelines on trustworthy
352 AI [6]. It has been shown that some methods for
353 generating explanations risks leaking information
354 about the training data and thus violating privacy
355 [35]. This may be a more general problem, as a valid
356 explanation require a basis in the training data.

357 4.2 Contradicting definitions

358 Even if all humans could agree to some set of values,
359 implementation is not straight forward, due to in-
360 herent contradictions. As discussed in the previous
361 section, different values can contradict each other,
362 but even a single value have room for contradic-
363 tions, based on the vast number of ways in which

364 they can be defined. A simple case is values that
365 act as constraints, but where the cut-off may be
366 different depending on who is asked. For example,
367 some could say that 100% accuracy on test data is a
368 sign of a too small test data set, or that the model
369 risks being sensitive to slight changes in the data
370 generating process, while someone else might require
371 100% accuracy to deem the model trustworthy.

372 A slightly more complex example is that of the
373 fairness paradox. Three well established definitions
374 in algorithmic fairness are equalized odds (false posi-
375 tive/negative rates do not depend on protected char-
376 acteristics), predictive parity (the predictive power
377 do not depend on protected characteristics), and
378 counterfactual fairness (the model should have no
379 causal effect between protected characteristics and
380 outcome), and it has been proven that these can not
381 occur at the same time [36]. Even definitions that
382 contradict each other can be gradually improved on,
383 but this again requires making trade-offs between
384 them [37].

4.3 The power to decide

385 From Sections 2 and 3, it is clear that human values
386 are not universal, and that there is a multitude of
387 ethical frameworks. This means, inherently, that
388 there is no superior ethical theory. Assuming that
389 there is would ignore millennia of ethical debates
390 and the diversity of human cultures. There have
391 certainly been efforts to define a singular ethics and
392 a universally correct way of living, but historically
393 many such attempts have supported colonialist ef-
394 forts and become an excuse to treat people with
395 other ways of living as morally inferior.

396 This is also a clear problem in current efforts
397 of globally used AI models with “human-aligned”
398 values, which are not representative of all parts of
399 humanity [32]. This is also a clear problem when a
400 select number of experts try to figure out AI ethics
401 among themselves, by developing guidelines or writ-
402 ing papers for exclusive conferences. However many
403 AI experts and ethics professors put their heads
404 together, most people are still not a part of the dis-
405 cussions, and do not have a say in the process. In the
406 AI community, we can support any number of values
407 such as “equity”, “reducing bias”, and “privacy”,
408 but without talking to different stakeholders and
409 considering a pluralistic view on ethics, it will still
410 be an ethics for the few. For the groups represented
411 among the experts, and the people we can imagine.

412 A related question is what happens when a small
413 number of pre trained AI systems is implemented in
414 a large set of tasks, similarly to what we can see with
415 Large Language Models (LLMs) today. Even assum-
416 ing they represent some sort of average or generally
417 agreed upon ethics, small everyday decisions are
418 suddenly taken using a universally adapted ethical
419

420 system, instead of a multitude of more intuitions.

421 **4.4 A future: The universal ethical**
422 **validation layer**

423 Consider a potential future, a few years from now.
424 The field of AI ethics has continued to grow, but the
425 conclusions in the large journals begin to coalesce on
426 a set of clearly defined values. A group of Machine
427 learning researchers see their chance to make an
428 actual improvement in the world, by creating the
429 *Universal Ethics Assessment Layer*.

430 They review the top cited papers on ethical AI
431 and create a layer that can create an ethicality score
432 for any machine learning model output. The score
433 is based on formalised concepts of values such as
434 fairness, honesty, and justice, and the importance
435 of them weighted according to their prevalence in
436 the surveyed papers. The layer is made to be ap-
437 plied to deep learning architectures, taking in the
438 input and the output of the main model, as well
439 as the results of querying the model for different
440 inputs for comparison. Based on this an ethical-
441 ity score is provided, which during training is fed
442 back to the model as a reward until the ethicality
443 score is consistently over a certain threshold. If the
444 threshold is violated during deployment, the model
445 will randomly make small, temporary tweaks to its
446 weights until an output with a valid ethicality score
447 is achieved.

448 The results are presented in a paper in Nature,
449 with validation of the model on a large variety of ar-
450 chitectures and tasks, and showing accomplishment
451 of a higher ethicality score than humans on com-
452 parative tasks, from image generation and chatting,
453 to cancer diagnosis and mortgage approval. The
454 paper is soon one of the most cited in the field of
455 AI ethics, and newspaper headlines read "The align-
456 ment problem is solved". The work receives some
457 critique for setting fixed weights on the different val-
458 ues, but the debate dies quietly with the developers
459 answer that the weights can be easily changed in
460 their open-source code. However, no technical user
461 seems to bother.

462 As a part of efforts for more ethical AI research,
463 many of the large ML venues start requiring usage of
464 the Ethical layer in submitted work. On the website
465 of one of the largest conferences on machine learning,
466 one can read

467 "To make accuracy stay a useful tool in the evalu-
468 ation of machine learning methodology, it is impor-
469 tant that ethicality can not be traded for accuracy.
470 Therefore, each submission aiming to improve accu-
471 racy or developing new architectures must include
472 the ethical layer in their model, with the original
473 weights and an ethicality threshold of 15. Any non-
474 motivated deviations from this will lead to desk
475 rejection."

476 The ethical layer further becomes a simple way
477 for companies to show compliance with legal require-
478 ments on trustworthy and responsible AI, and soon
479 this layer is a basic building block in all machine
480 learning systems. The same ethical basis is used all
481 over, and when a healthcare provider is accused of
482 racism they simply point to the ethicality score of
483 the treatment decision, and the issue is settled.

484 Since models are now verifiably ethical, people
485 start to imitate them in their personal moral deci-
486 sions. As time goes and the universal ethicality is
487 more widely applied, people forget how essentially, a
488 few AI researcher single-handedly defined ethicality
489 all over the world.

490 This may in some senses be an appealing solution,
491 but it also takes away something of the plurality of
492 human views. When the same ethical trade-offs is
493 propagated through all parts of society and AI tools
494 are harder to avoid, minority opinions become less
495 reflected and marginalisation risks increasing.

496 **5 Democratic approaches to**
497 **AI ethics**

498 One way to handle these impossibilities is to take
499 a more political view of the issues, and consider
500 democratic approaches to AI ethics.

501 **5.1 Democratic regulation**

502 One straight forward way to democratise AI is by
503 democratically constructing legal requirements for
504 its development, research, and deployment. Un-
505 common opinions, minority groups, and people not
506 residing in the country may still be missed, but by
507 voting a larger part of the population has a say. One
508 example of legal regulation is the EU AI Act, which
509 regulates how AI systems can be used, together with
510 parts of GDPR when personal data is concerned.
511 Arguments for regulation via legal means have been
512 lifted before, both based on the concept of ethicality
513 as inherently political [38] and the idea that self-
514 regulation from tech companies centralises the power
515 over AI systems even more [39].

516 **5.2 Transparency**

517 A core part of a democratic society is transparency
518 that allows the public to question and analyse po-
519 litical decisions, and the actions of others in power.
520 When AI systems are used in more sensitive situa-
521 tions and given more autonomy, the transparency of
522 them becomes important by the same reason. Note
523 that this transparency is not equivalent to common
524 notions of explainability, which often focuses on ex-
525 plaining specific decisions. Rather, the democratic
526 transparency requires public access to information
527 about what values are encoded in the machine, and

528 how. Visiting the website for *chatGPT*, the informa- 578
529 tion provided about value alignment is significantly 579
530 lacking. This is a part of their explanation:

531 We randomly selected several alternative 581
532 completions, and had AI trainers rank 582
533 them. Using these reward models, we can 583
534 fine-tune the model using Proximal Policy 584
535 Optimization. 585

536 and another states that 586

537 We're using the Moderation API to warn 588
538 or block certain types of unsafe content, ... 589

539 where the goal of the mentioned Moderation API 591
540 is to flag content that is sexual, hateful, violent, or 592
541 promotes self-harm. Nowhere do they give details
542 on who the AI trainers were or what values they
543 prioritised, nor what they define as hateful and why
544 they decided to flag exactly those four as unsafe
545 content.

546 To be able to have thorough discussions on the 594
547 values encoded in models like *chatGPT*, they need to 595
548 be more transparent than this. If the data trainers 596
549 where all male college students in the US, quite many 597
550 human perspectives are missed in the fine tuning 598
551 process. Exactly what data was used to train the 599
552 original model is also not publicly available. 600

5.3 Participatory design 601

554 Another approach to involve more people in AI de- 602
555 velopment is through participatory design. The 603
556 concept is wide and can mean anything from sur- 604
557 veying a few colleagues about an interface design 605
558 choice, to letting a representative (by some defini- 606
559 tion) group fully own the project. A recent survey of 607
560 participatory design in AI development noted that
561 research project often do not enable the extra time
562 needed to fully involve a part of the public, which
563 has led to smaller scale involvements or using prox-
564 ies, such as a stand in person for another group or a
565 ML model trained on different peoples preferences
566 [dem:delgado2023participatory].

6 Frameworks for analysis 608

568 Several different frameworks have been proposed 610
569 to analyse the ethical implications of AI models. 611
570 Here I present two; the assessment list from the 612
571 EU expert panel and ACROCPoLis from a recent 613
572 scientific paper. 614

6.1 Trustworthy AI assessment list 615

574 As a part of the EU guidelines for trustworthy AI, 616
575 an assessment list was released [6]. The list is several 617
576 pages long and consists of concrete questions to ask 618
577 when developing an AI system, such as: “Did you 619

578 assess to what extent the decisions and hence the out- 579
579 come made by the AI system can be understood?”, 580
580 “Did you establish mechanisms to ensure fairness in 581
581 your AI systems? Did you consider other potential 582
582 mechanisms?”, and “Did you assess whether there 583
583 could be persons or groups who might be dispropor- 584
584 tionately affected by negative implications?”. 585

585 The idea is that when developing or implementing 586
586 an AI system the list should be used to verify that all 587
587 aspects of the guidelines have been considered. The 588
588 list is thorough and captures all the different parts 589
589 of the guidelines, but are at points vague, clearly 590
590 focused towards larger corporations, or without spec- 591
591 ified results, e.g. questions of the form “Did you 592
592 assess...?”. 593

6.2 ACROCPoLis 593

594 ACROCPoLis is a recently developed framework 595
595 to assess the ethical implications of an AI system 596
596 [40]. For good and bad, it is significantly broader in 597
597 formulation than the above mentioned assessment 598
598 list. The idea is to identify the categories Actors, 599
599 Context, Resources, Outcome, Criteria, Power, and 600
600 the Links between them. This is in many ways 601
601 based on a feminist ethics with its focus on context 602
602 and relations, but it also aims to facilitate viewing 603
603 issues from different perspectives by singling out 604
604 the differences in context and power between the 605
605 stakeholders (actors). However, it does not provide 606
606 any clear answers about ethicality, but rather a 607
607 language to discuss it.

7 A practical guide 608

609 To conclude, I provide a practical guide to navigat- 609
610 ing the complex questions of AI ethics in everyday 610
610 research activity. A few basic principles are de- 611
611 rived from the arguments this far, which is then 612
612 implemented as a short checklist inspired by the 613
613 ACROCPoLis framework and the EU assessment 614
614 list. Finally the checklist is applied to a hypotheti- 615
615 cal research project. 616

7.1 Basic principles 617

618 Based on the discussion this far, I propose four basic 618
619 principles for responsible AI research. 619

620 **There is no universal morality.** As argued 620
621 in the previous section, a universal ethics for AI is 621
622 inherently impossible. To still perform responsible 622
623 research in the field of machine learning, this needs 623
624 to be both accepted and acknowledged. Assuming 624
625 that there is a singular solution to ethical AI risks 625
626 authoritarian consequences and further marginalisa- 626
627 tion of minority groups. 627

628 **Each stakeholder deserves consideration.** 628
629 This is a major part of the democratic view on ethics, 629

Responsible ML research

Stage 1: Identify potential impact

- Consider the potential applications of your research project, and how likely they are. Can you stand behind them with good conscience?
- Who would deploy or use this type of system, and what are the power dynamics?
- What technologies can this knowledge enable, and who does it affect?
- Where does your data come from, are there data labellers and/or unknowing subjects involved?
- What resources do you use, what is the climate impact, and who does this affect?

Stage 2: Gather different perspectives

- See if you can get a chance to talk to the identified stakeholders through personal contacts, public events, or online communities.
- Otherwise, read works by them or journalistic reports, and imagine yourself in their place.

Stage 3: Make transparent trade-offs

- Make trade-offs based on your ethical values, but remember the position of power this puts you in.
- Clearly present these trade-offs in all communication about your research, whether it is in a scientific paper or public outreach.

Stage 4: Community engagement

- Ask other researchers about their ethical stance and the trade-offs of their research at conferences.
- Discuss ethical issues with colleagues.
- Engage in public dialogue and help unveil hidden values of discussed ML systems.

Figure 1. The checklist

but is also necessary for both determining utility distributions and to follow the feminist approach. One part of this is to work towards identifying all relevant actors; not only those directly affected by an algorithm, but those affected by the data acquisition or resource usage. Animals and the earth itself can also be included here, to use tools from environmental ethics.

There will be trade-offs. Most technology, and most actions, have a wide array of consequences, and often some are wanted and others unwanted. This means that trade-offs are necessary. There is no strictly correct way to choose these, but by communicating the choices they can be questioned, criticised, or reproduced. Ethical choices are hard, but unavoidable. Choosing to step away from a research topic is also a moral judgement.

Silence is a political decision. It is not the job of a researcher to decide what is the right usage of technology, but as the ones at the forefront of the knowledge researchers are needed to explain and point out ethical dilemmas and potential consequences. Especially, decisions on ethical judgements in AI development and research need to be communicated, otherwise it is close to impossible for others to consider or criticise the view. To not communicate these things, or not take part in public discussion, is a political choice, especially considering the importance of knowledge and journalistic investigation in a well functioning democracy.

7.2 The checklist

In Figure 1, the basic principles are combined into a checklist to be used in everyday research activity. The format is inspired by the AI assessment list, but with more focus on actions, fewer points, and a clearer focus towards smaller scale research. The ACROCPoLis framework is incorporated with actions corresponding to identify the contextual and power dynamics information of the problem, with an extra focus on identification of relevant stakeholders as both values like privacy and utilitarian approaches also require knowledge of the actors in the system. The aim of the checklist is not to favour specific ethical values, apart from the basic principles from the previous subsection, but rather to guide the researcher to make their own, informed decisions based on their ethical values, while still emphasising how this puts us in a position of power.

The checklist is divided into four stages, mainly based on when in the research project they would come into play. In the first stage, potential stakeholders and impact on them should be identified. This is not as action focused as the later stages, but is necessary to do beforehand. The questions are mainly aimed at enforcing consideration of different aspects and groups. The second stage is about gath-

686 ering the perspectives of the actors identified in the
687 first stage. This allows for a deeper basis of under-
688 standing to base the later trade-offs on. The third
689 stage is where the researchers own ethical stance
690 comes into play. As has been shown there is no
691 universally correct approach to ethical AI, so here
692 the researcher must use their knowledge from the
693 previous stages together with their own values to
694 decide what trade-offs to make. However, to enable
695 an ongoing debate on the ethical values in machine
696 learning, all of these trade-offs need to be thoroughly
697 documented and communicated. The forth stage is
698 a bit outside of any specific research project, but
699 rather a continuous work to foster a research culture
700 in ML where ethical trade-offs are made open to
701 scrutiny and debate among both researchers and
702 the broader public. In the Appendix, an example is
703 provided where the proposed framework is applied
704 to a hypothetical research project by a fictional re-
705 search group. It is omitted from the main paper
706 for consistency, but may provide a more digestible
707 interpretation of the checklist.

708 8 Discussion and conclusion

709 Ethics is a complex subject without simple answers,
710 in AI just as in any other branch of human ex-
711 periences. In this paper I presented a variety of
712 approaches to ethics and how they can, and have
713 been, applied to the field of AI and machine learning.
714 Contradictions arise both from the plurality of the
715 approaches and disagreements about definitions of
716 fundamental values. What makes the field of ethics
717 for AI somewhat unique is that many ethical deci-
718 sions are implicitly made by researchers developing
719 new algorithms, potentially without realising it. A
720 simple assumption that notions of fairness can be
721 directly incorporated into a utility function is an
722 indirect utilitarian view on the issue, while silence
723 on environmental costs for network training is a de-
724 cision on which actors matter. Being a researcher in
725 this field is also being in a position of power, with
726 agency to shape future research directions as well
727 as the public understanding of the algorithms we
728 implement. I propose that we openly acknowledge
729 the inherent impossibilities of singular technical so-
730 lutions to ethical AI, and instead take responsibility
731 for the trade-offs and values we decide to base our
732 work on. Further, I suggest we give away some of
733 our implicit power by communicating more trans-
734 parently on ethics issues and lifting the voices of
735 groups less often represented in the AI community.

References

- 736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
- [1] *AI and Ethics*. Springer Nature Link. URL: <https://link.springer.com/journal/43681>.
 - [2] *the AI Ethics Journal*. Artificial Intelligence Robotics Ethics Society. URL: <https://www.aiethicsjournal.org/>.
 - [3] *Artificial Intelligence, Ethics, and Society*. Association for the Advancement of Artificial Intelligence. URL: <https://aaai.org/conference/aies/>.
 - [4] *Frontiers of Artificial Intelligence, Ethics, and Multidisciplinary Applications*. Springer. URL: <https://www.faiema.org/>.
 - [5] *Secure and Trustworthy Machine Learning*. IEEE. URL: <https://satml.org/>.
 - [6] E. H.-L. E. G. on AI. *Ethics guidelines for trustworthy AI*. European Union. 2019.
 - [7] *the Ethics of Artificial Intelligence*. the United Nations Educational, Scientific and Cultural Organization. URL: <https://unesdoc.unesco.org/ark:/48223/pf0000381137>.
 - [8] *ARTIFICIAL INTELLIGENCE (AI) ETHICS GUIDE*. United States Agency International Development.
 - [9] C. Geldhauser and H. Diebel-Fischer. “Is diverse and inclusive AI trapped in the gap between reality and algorithmizability?” In: *Northern Lights Deep Learning Conference*. PMLR. 2024, pp. 75–80.
 - [10] J. Driver. “The History of Utilitarianism”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta and U. Nodelman. Winter 2022. Metaphysics Research Lab, Stanford University, 2022.
 - [11] V. X. Chen and J. Hooker. “A guide to formulating equity and fairness in an optimization model”. In: *Preprint* (2021), pp. 162–174.
 - [12] T. Kitagawa and A. Tetenov. “Who Should Be Treated? Empirical Welfare Maximization Methods for Treatment Choice”. In: (2018). DOI: <https://doi.org/10.3982/ECTA13288>.
 - [13] N. Robert. *Anarchy, State, and Utopia*. Basic Books, Harper Colins, 1974. ISBN: 9780465051007.
 - [14] U. K. L. Guin. “The Ones Who Walk Away from Omelas”. In: *New Dimensions 3*. Ed. by R. Silverberg. 1973.
 - [15] N. Mehrabi, F. Morstatter, N. Saxena, K. Lerman, and A. Galstyan. “A survey on bias and fairness in machine learning”. In: *ACM computing surveys (CSUR)* 54.6 (2021), pp. 1–35. DOI: <https://doi.org/10.1145/3457607>.

- [16] D. Pessach and E. Shmueli. “A review on fairness in machine learning”. In: *ACM Computing Surveys (CSUR)* 55.3 (2022), pp. 1–44. DOI: <https://doi.org/10.1145/3494672>. 844
845
- [17] A. F. Cooper, E. Abrams, and N. Na. “Emergent unfairness in algorithmic fairness-accuracy trade-off research”. In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 2021, pp. 46–54. DOI: <https://doi.org/10.1145/3461702.346251>. 846
847
848
849
850
851
- [18] P. Räsänen, E. Roine, H. Sintonen, V. Semberg-Konttinen, O.-P. Ryyänen, and R. Roine. “Use of quality-adjusted life years for the estimation of effectiveness of health care: A systematic literature review”. In: *International Journal of Technology Assessment in Health Care* 22.2 (2006), 235–241. DOI: [10.1017/S0266462306051051](https://doi.org/10.1017/S0266462306051051). 852
853
854
855
856
- [19] EU. *GDPR*. 857
- [20] L. Li, Y. Fan, M. Tse, and K.-Y. Lin. “A review of applications in federated learning”. In: *Computers & Industrial Engineering* 149 (2020), p. 106854. DOI: <https://doi.org/10.1016/j.cie.2020.106854>. 858
859
860
861
862
863
864
865
- [21] C. Dwork. “Differential privacy: A survey of results”. In: *International conference on theory and applications of models of computation*. Springer, 2008, pp. 1–19. 866
867
- [22] A. Lagerkvist, M. Tudor, J. Smolicki, C. M. Ess, J. Eriksson Lundström, and M. Rogg. “Body stakes: an existential ethics of care in living with biometrics and AI”. In: *AI & SOCIETY* 39.1 (2024), pp. 169–181. DOI: <https://doi.org/10.1007/s00146-022-01550-8>. 868
869
870
871
872
- [23] U. Aivodji, H. Arai, O. Fortineau, S. Gambs, S. Hara, and A. Tapp. “Fairwashing: the risk of rationalization”. In: *Proceedings of the 36th International Conference on Machine Learning*. Ed. by K. Chaudhuri and R. Salakhutdinov. Vol. 97. Proceedings of Machine Learning Research. PMLR, 2019, pp. 161–170. URL: <https://proceedings.mlr.press/v97/aivodji19a.html>. 873
874
875
876
877
- [24] C. Villegas-Galaviz and K. Martin. “Moral distance, AI, and the ethics of care”. In: *AI & society* 39.4 (2024), pp. 1695–1706. DOI: <https://doi.org/10.1007/s00146-023-01642-z>. 878
879
880
- [25] P. Ricaurte. “AI for/by the majority world: From technologies of dispossession to technologies of radical care”. In: *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*. AIES ’23. Montréal, QC, Canada: Association for Computing Machinery, 2023, 3–4. ISBN: 9798400702310. DOI: [10.1145/3600211.3607544](https://doi.org/10.1145/3600211.3607544). URL: <https://doi.org/10.1145/3600211.3607544>. 881
882
883
884
885
886
887
888
889
890
- [26] O. Tracey and R. Irish. “Explainability for All: Care Ethics for Implementing Artificial Intelligence”. In: *2023 IEEE International Symposium on Technology and Society (ISTAS)*. 2023, pp. 1–8. DOI: [10.1109/ISTAS57930.2023.10306199](https://doi.org/10.1109/ISTAS57930.2023.10306199). 891
892
893
894
895
896
- [27] A. Brennan and N. Y. S. Lo. “Environmental Ethics”. In: *The Stanford Encyclopedia of Philosophy*. Ed. by E. N. Zalta and U. Nodelman. Summer 2024. Metaphysics Research Lab, Stanford University, 2024. 897
- [28] S. A. Budenny, V. D. Lazarev, N. N. Zakharenko, A. N. Korovin, O. Plosskaya, D. V. Dimitrov, V. Akhripkin, I. Pavlov, I. V. Osledets, I. S. Barsola, et al. “Eco2ai: carbon emissions tracking of machine learning models as the first step towards sustainable ai”. In: *Doklady Mathematics*. Vol. 106. Suppl 1. Springer, 2022, S118–S128. DOI: <https://doi.org/10.1134/S1064562422060230>. 898
899
900
901
902
903
904
905
906
- [29] V. Joler and K. Crawford. *Anatomy of an AI system*. 2018. URL: <https://anatomyof.ai/>. 907
- [30] S. D. Baum and A. Owe. “Artificial intelligence needs environmental ethics”. In: *Ethics, Policy & Environment* 26.1 (2023), pp. 139–143. DOI: <https://doi.org/10.1080/21550085.2022.2076538>. 908
909
910
911
912
- [31] C. Moyano-Fernández and J. Rueda. “AI, sustainability, and environmental ethics”. In: *Ethics of artificial intelligence*. Springer, 2024, pp. 219–236. DOI: https://doi.org/10.1007/978-3-031-48135-2_11. 913
914
915
916
917
- [32] M. Atari, M. J. Xue, P. S. Park, D. Blasi, and J. Henrich. “Which humans?” In: *PsyArXiv* (2023). 918
919
920
- [33] V. Dignum. “Relational artificial intelligence”. In: *arXiv preprint arXiv:2202.07446* (2022). DOI: [10.48550/arXiv.2202.07446](https://doi.org/10.48550/arXiv.2202.07446). 921
922
923
- [34] S. T.-I. Wu, D. Demetriou, and R. A. Husain. “Honor ethics: The challenge of globalizing value alignment in AI”. In: *Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency*. 2023, pp. 593–602. DOI: <https://doi.org/10.1145/3593013.359402>. 924
925
926
927
928
929
930
931
932
933
934
935
936
- [35] R. Shokri, M. Strobel, and Y. Zick. “On the privacy risks of model explanations”. In: *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*. 2021, pp. 231–241. DOI: <https://doi.org/10.1145/3461702.346253>. 937
938
939
940
941
942
943
944
945
946
947
948
949
950

897 [36] F. Beigang. “Yet Another Impossibility The- 951
 898 orem in Algorithmic Fairness”. In: *Minds* 952
 899 *and Machines* 33.4 (2023), pp. 715–735. DOI: 953
 900 [https://doi.org/10.1007/s11023-023-](https://doi.org/10.1007/s11023-023-09645-x) 954
 901 [09645-x](https://doi.org/10.1007/s11023-023-09645-x). 955

902 [37] C. Hertweck and T. Ráz. “Gradual (in) com- 956
 903 patibility of fairness criteria”. In: *Proceedings* 957
 904 *of the AAAI Conference on Artificial Intelli-* 958
 905 *gence*. Vol. 36. 11. 2022, pp. 11926–11934. DOI: 959
 906 [https://doi.org/10.1609/aaai.v36i11.](https://doi.org/10.1609/aaai.v36i11.21450) 960
 907 [21450](https://doi.org/10.1609/aaai.v36i11.21450). 961

908 [38] G. van Maanen. “AI ethics, ethics washing, 962
 909 and the need to politicize data ethics”. In: 963
 910 *Digital Society* 1.2 (2022), p. 9. DOI: <https://doi.org/10.1007/s44206-022-00013-3>. 964

912 [39] R. Ochigame. “The invention of ‘ethical AI’: 965
 913 How big tech manipulates academia to avoid 966
 914 regulation”. In: *Economies of virtue* 49 (2019). 967

915 [40] A. Aler Tubella, D. Coelho Mollo, A. Dahlgren 968
 916 Lindström, H. Devinney, V. Dignum, P. Er- 969
 917 icson, A. Jonsson, T. Kampik, T. Lenaerts, 970
 918 J. A. Mendez, et al. “ACROCPoLis: A de- 971
 919 scriptive framework for making sense of fair- 972
 920 ness”. In: *Proceedings of the 2023 ACM Con-* 973
 921 *ference on Fairness, Accountability, and Trans-* 974
 922 *parency*. New York, NY, USA: Association for 975
 923 Computing Machinery, 2023, pp. 1014–1025. 976
 924 ISBN: 9798400701924. DOI: [10.1145/3593013.](https://doi.org/10.1145/3593013.3594059) 977
 925 [3594059](https://doi.org/10.1145/3593013.3594059). URL: [https://doi.org/10.1145/](https://doi.org/10.1145/3593013.3594059) 978
 926 [3593013.3594059](https://doi.org/10.1145/3593013.3594059). 979

927 **A Application of the checklist**
 928 **to a hypothetical research**
 929 **project**

930 **Disclaimer:** *This is a hypothetical research group*
 931 *and research project, meant to illustrate how the*
 932 *checklist can be used. The concerns arising is not*
 933 *based on any actual persons opinion, and the actions*
 934 *taken are not necessarily the best choices by any*
 935 *definition.*

936 The research group X has an idea for a new
 937 method to detect anomalies in sparse sequential
 938 data. They hope this can be used to improve auto-
 939 mated health care suggestions for patients that take
 940 regular tests during pregnancy, cancer treatment,
 941 or type 2 diabetes controls. Before implementing
 942 anything, they go to stage 1 of the checklist.

943 **Stage 1:** When considering potential applica-
 944 tions, they realise that these can be very broad;
 945 from detection of credit card fraud to extraction
 946 of personal information. The first would be a rela-
 947 tively easy change of application, while the latter
 948 would require significant additional efforts. With
 949 the healthcare application they initially had in mind,
 950 they think hospitals would be the main users, or

more probable companies selling the tools to the hos-
 pitals. Here there are power dynamics both between
 the companies and hospital, the hospital steering
 committee and the doctors, and between doctors and
 patients, with the patients in many senses furthest
 down the hierarchy. The data would come from a
 nearby hospital that the group is collaborating with,
 where data comes from patients and is labelled by
 doctors and nurses. The plan is to use anonymised
 data so that they do not need to ask for patients per-
 mission and go through ethics reviews. The model
 training do not require excessive compute power.

Stage 2: The group schedules a meeting with
 some doctors and nurses at the hospital, all related
 to type 2 diabetes care. Some of the doctors are
 very positive to decision tools to help them make
 correct decisions under time pressure, while some of
 the nurses are afraid the patients will feel less com-
 comfortable when their data is fed to a machine, even
 if it stays locally. The doctors that have had col-
 laborations with the research group before is mostly
 excited by the potential of revealing new patterns
 in the data they already use for decisions. Patients
 are invited to attend a workshop on the usage of the
 proposed methods, and the views are mixed. Many
 do not like the idea of their data being used, even af-
 ter understanding that it is anonymised. Some hope
 the method can allow them to take tests at their
 local health centre instead of going to the hospital,
 if decisions are taken by a machine anyway, while a
 few others says they would stop going to controls if
 AI is in any way involved.

Stage 3: It is now time to make trade-offs, and
 the research group has trouble reaching consensus.
 Some think the potential knowledge gain outweighs
 the patient concerns about using anonymous data
 and propose they should use the data and train
 the model, but collaborate with some medical re-
 searchers to see if the model can identify previously
 unknown patterns. This would avoid using the
 model to make actual decisions, and thus not de-
 crease the willingness of some patients to attend
 their controls. However, to meet the patients with
 the privacy concerns half way, they prepare an in-
 formation folder on how the data anonymisation
 works hoping this will decrease their fears. Another
 group is sceptical whether any new medical knowl-
 edge would be gained and instead focuses on how the
 proposed model scheme could be used to improve
 accessibility of healthcare for those living far from
 the hospital. This second group additionally do not
 want to use the patient data, even anonymised, due
 to the patient concerns. Instead they want to work
 with the doctors at the hospital on how to encode
 the patterns already used into a simple and explain-
 able model that could be used by type 2 diabetes
 patients to determine how often they need to attend
 controls at the hospital, based on self reported blood

1009 sugar measurements and medical history. In the end,
1010 the two subgroups decide to pursue their respective
1011 research direction, and both present the trade-offs
1012 they have done to the patients and doctors they
1013 talked with before, as well as in their papers.