

---

# Grounding and Validation of Algorithmic Recourse in Real-World Contexts: A Systematized Literature Review

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1       The aim of algorithmic recourse (AR) is generally understood to be the provision  
2       of “actionable” recommendations to individuals affected by algorithmic decision-  
3       making systems, in an attempt to offer the capacity for taking actions that may  
4       lead to more desirable outcomes in the future. Over the past few years, AR  
5       literature has largely focused on theoretical frameworks to generate “actionable”  
6       counterfactual explanations that further satisfy various desiderata, such as diversity  
7       or robustness. We believe that algorithmic recourse, by its nature, should be seen  
8       as a practical problem: real-world socio-technical decision-making systems are  
9       complex dynamic entities involving various actors (end users, domain experts,  
10       civil servants, system owners, etc.) engaged in social and technical processes.  
11       Thus, research needs to account for the specificities of systems where it would  
12       be applied. To evaluate how authors envision AR “in the wild”, we carry out a  
13       systematized review of 127 publications pertaining to the problem and identify the  
14       real-world considerations that motivate them. Among others, we look at the ways  
15       to make recourse (individually) actionable, the involved stakeholders, the perceived  
16       challenges, and the availability of practitioner-friendly open-source codebases.  
17       We find that there is a strong disconnect between the existing research and the  
18       practical requirements for AR. Most importantly, the grounding and validation of  
19       algorithmic recourse in real-world contexts remain underexplored. As an attempt  
20       to bridge this gap, we provide other authors with five recommendations to make  
21       future solutions easier to adapt to their potential real-world applications.

## 22   1 Introduction

23   Algorithmic decision-making (ADM) tools are frequently seen as a way to improve decision processes  
24   in a variety of high-stakes domains such as public administration [47, 146] or healthcare [45, 87].  
25   Deep learning models have attracted much attention due to their perceived high performance, but  
26   the predictions of such models cannot be interpreted by humans, hence end users – both individuals  
27   subjected to algorithmic decisions and decision-makers operating on them – are placed in a position  
28   where they are unable to understand the grounds of a prediction, act on it, or trust it [159].

29   To help address this problem, a variety of explanation methods has been proposed. Of particular  
30   interest for this paper are counterfactual explanations (CEs) that attempt to explain the predictions for  
31   individual instances of data, taking the form of conditional statements such as “*if the value of feature*  
32   *x was a instead of b, the model would have predicted class y instead of z*”. They are perceived to be  
33   an attractive approach to explanation that does not require “opening the black box” [151] and have  
34   been argued to align with the ways that humans naturally reason about events [84].

35 CEs are also seen as the go-to method for algorithmic recourse (AR), or the generation of actionable  
36 recommendations that provide people with the knowledge needed to achieve more desirable predic-  
37 tions in ADM systems. Recourse is distinct from the “explanation” or “justification” of algorithmic  
38 decisions, and more closely related to the notion of contestability of Artificial Intelligence [7] in that  
39 it aims not only to improve the *trust* in the algorithm, but also embrace human *agency* [142].

40 Algorithmic recourse is an inherently practical problem in that it resembles a bureaucratic complaint  
41 process: an individual unhappy with some decision engages with a representative of the issuing  
42 organization, in an attempt to overturn it. Yet, we observe that much of the existing work is highly  
43 theoretical, with little consideration of whether it could be applied in organizational settings [see  
44 also 18]. Deploying AR in realistic systems without analyzing its mechanics in a broader context  
45 and without knowing what types of dynamics are expected to arise is bound to lead to unanticipated  
46 outcomes. Many of them will be undesirable and even potentially unsafe, and impossible to validate  
47 with respect to a set of requirements because the requirements for AR are *necessarily* socio-technical.

48 **Societal and institutional components of algorithmic recourse are the focal point of our work,**  
49 as we look beyond the typical technical considerations to assess the practical aspects of the problem.

50 To that end, we contribute a *systematized review* of 127 publications that address the goals of  
51 algorithmic recourse and we evaluate to what extent they incorporate such practical considerations.  
52 We characterize our approach as *systematized* because we follow a fully systematic approach to the  
53 collection of publications, but their selection is not necessarily exhaustive [46] as many impactful  
54 ideas in computer science are published only in the form of pre-prints. Based on our analysis, we also  
55 provide other authors with five recommendations on how to improve the practicality of AR research.

56 The rest of the manuscript is structured as follows. In Section 2 we elaborate on the background of  
57 our work. Then, in Section 3 we describe our approach to this review. Next, Section 4 introduces  
58 our findings. Section 5 provides a discussion of our results, introduces our recommendations, and  
59 addresses the limitations of the current work. Finally, Section 6 forms the conclusion to this paper.

## 60 2 Background

### 61 2.1 On algorithmic recourse

62 Algorithmic – or actionable, individual – recourse was introduced in [138] as “*the ability of a person*  
63 *to change the decision of the model through actionable input variables*”, building on the earlier  
64 work of [151] who argued that CEs are a psychologically-grounded way to (1) help decision-subjects  
65 understand an algorithmic decision, (2) provide them with information needed to contest it, and (3)  
66 inform about actions that could be taken to overturn it. For instance, consider a person who has  
67 unsuccessfully applied for a loan; they may then receive AR such as “*if you requested \$5000 less,*  
68 *you would qualify for this loan*”. The key consideration for AR is “actionability”, which entails that  
69 the recipient of the recommendation should be capable of implementing it. If they had been informed  
70 “*if you were 10 years younger, you would qualify for the loan*”, they would have still received a  
71 valid CE, *but not* recourse. More recently [69] has recast the problem as reasoning about minimal  
72 interventions on the structural causal model. This formulation (at least theoretically) addresses an  
73 important shortcoming of “correlational” recourse. Without accounting for the downstream causal  
74 effects of actions, an individual may exert more effort than necessary and still fail to achieve the  
75 target outcome. Indeed, counterfactuals are an inherently causal concept [103].

76 We note that problems similar to AR have been studied under a variety of different names: *actionable*  
77 *knowledge discovery* [e.g., 2], *action rules mining* [e.g., 110], *inverse classification* [e.g., 5], *why*  
78 *not questions* [e.g., 58], or *actionable feature tweaking* [134]. These alternative formulations have  
79 generally focused on “business” knowledge, rather than individual recommendations, but ultimately  
80 the goal of all these approaches is to extract information from a (black-box) model that allows the  
81 user – an individual or a decision-maker – to act. We highlight them to emphasize that AR does  
82 not have to be achieved through the means of CEs. Rather CEs should be seen as *one of the means*  
83 to achieve AR, particularly promising in that they do not require expert-level understanding of the  
84 model to be useful. Nonetheless, we decide to distinguish between the literature on AR (commonly  
85 equated with actionable CEs), and these alternative formulations in our work.

86 Existing research has generally considered AR in simplistic settings that are far removed from  
87 real-world socio-technical decision-making systems, where it would be implemented as a process.

88 For example, such systems are dynamic [113, 137], must support the implementation of AR at scale  
89 [9, 94], and involve various stakeholders beyond the end users [17, 151]. Moreover, if the intended  
90 goal of AR is to help individuals subjected to algorithmic decisions in an effective manner, research  
91 must entail a rich understanding of “actionability” to account for the differences between them [142].

## 92 2.2 On the position of our review

93 Several groups of authors have previously surveyed the landscape of counterfactual explanations in  
94 general, and algorithmic recourse specifically. Perhaps the most relevant to our work is [71], which  
95 discusses five deficits of research on CEs, with a special focus on the (lack of) psychological grounding.  
96 Another pertinent publication is [70], which attempts to unify the definitions and formulations of  
97 AR in existing literature, but the work primarily focuses on technical aspects. Next, [143] develops  
98 a rubric to compare counterfactual explainers (equated with AR) and identifies 21 research challenges.  
99 While these also remain mostly technical, several of them are relevant to our work, for instance, CEs  
100 “as an interactive service to the applicants” or reinforcing “the ties between machine learning and  
101 regulatory communities”. More recently, [48] reviewed and benchmarked a number of CE generators,  
102 but AR is only a secondary consideration in the work. We also highlight [130], which is the only  
103 systematic review of counterfactual and contrastive approaches to date. The authors understand CEs  
104 as a way to justify model predictions (i.e., they are different from AR). We agree with this distinction  
105 in that CEs can be useful for reasons other than recourse, such as model debugging [e.g., 1, 122].  
106 Finally, although not reviews, [13] and [142] are particularly relevant to our work, offering critical  
107 perspectives on AR and addressing multiple shortcomings of recourse literature.

## 108 3 Methods

109 In this section, we briefly discuss our approach to the literature review following the SALSA – Search,  
110 Appraisal, Synthesis, Analysis – framework introduced in [46]. We also provide a more detailed  
111 description to allow for the reproduction of our process in the supplementary materials. Figure 1  
112 presents our process in the form of a PRISMA flow diagram [97].

### 113 3.1 Search

114 We make use of three search engines to collect the initial set of studies: ACM Digital Library, IEEE  
115 Xplore, and SCOPUS. Given the previously mentioned blurry distinction between AR and CEs,  
116 we consider the papers discussing either problem. In a small scoping review, we identify several  
117 keywords common to publications on recourse, as well as several equivalent terms to build the query.  
118 We search in titles, abstracts, and keywords, arriving at 3092 records after de-duplication. To facilitate  
119 the screening process, we employ the open-source ASReview tool, which makes use of an active  
120 learning approach to re-order the set of publications, such that the most relevant ones are always  
121 “at the top of the stack” [139]. The researchers behind the tool suggest employing a stopping rule  
122 measured in the number of consecutive irrelevant records, which we set to 30, or 1% of the entire  
123 dataset. We accept all papers that focus on algorithmic recourse and counterfactual explanations,  
124 completing the screening after evaluating 1040 abstracts, leading to 499 relevant records.

125 We observe that some important publications may be missing from our results. For instance, [151]  
126 was published in a legal journal that is not indexed by computer science search engines. Thus, we  
127 decide to augment the set of records by applying snowballing, which has been shown as a good  
128 alternative to databases in systematic reviews in software engineering [162]. We collect the references  
129 for the top 50 (10%) “most impactful” publications, measured by the number of citations. While this  
130 introduces several pre-prints into our result set [52, 61, 91, 113, 143, 150], we decide not to exclude  
131 them. Our review remains primarily concerned with peer-reviewed work. After adding the snowballed  
132 references to our dataset, we are left with 2018 records for the second screening with ASReview.  
133 This time, we look for publications that specifically refer to the problem of AR, “actionable” CEs, or  
134 modifying outcomes of automated decision-making systems. We employ a stricter stopping rule to  
135 minimize the risk of false negatives, completing the screening after 60 consecutive irrelevant records  
136 with 203 records considered for full-text appraisal. To allow for complete reproducibility of the  
137 search process, we provide an extended discussion (including queries) in the technical Appendix A.

138 **3.2 Appraisal**

139 We were able to retrieve all of the remaining 203 documents. For each document, we require that the  
 140 authors explicitly cite recourse as the center of interest, or look at (1) explanations (2) provided for  
 141 individual instances (3) with the goal of acting upon them (4) in an attempt to modify the predictions  
 142 (5) of a classification model. We exclude 51 publications as they are not on topic, primarily because  
 143 they focus on CEs for the sake of explanation. Four works in this category look at (what they  
 144 call) recourse but extend the problem to settings beyond the scope of this review: recommender  
 145 systems [31, 43, 145], text classification [37], and anomaly detection [27]. Further 15 publications  
 146 are duplicates, typically pre-prints of other documents that were included in the review. Next, 8  
 147 documents were published before [151] that sparked the research on AR, and thus we exclude them as  
 148 well. These look at the alternative formulations discussed earlier in Section 2.1. Finally, 2 documents  
 149 are not publications: one is an abstract of a talk, and the other is a student poster. For each document,  
 150 we answer a number of questions relating to the practical considerations introduced by the authors.

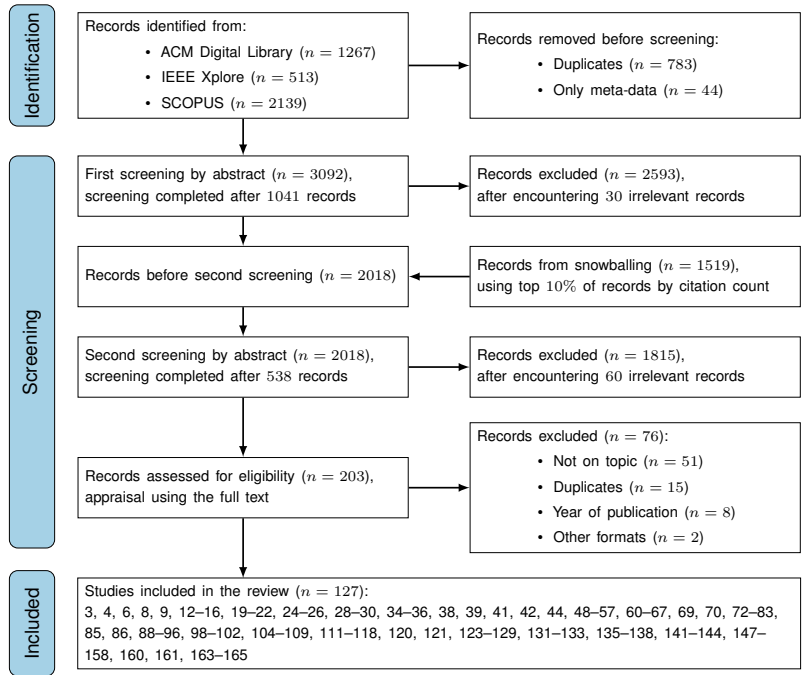


Figure 1: Identification of studies via databases and snowballing

151 **3.3 Synthesis**

152 To compile the results we carry out a standard thematic content analysis following the approach  
 153 presented in [40]. First, we explore the data extracted from the set of publications relevant to each  
 154 question to find the commonalities, which serves as the grounds for creating the initial set of codes.  
 155 We evaluate the documents against these codes and keep track of any other considerations. If such  
 156 considerations appear in multiple documents, we create new codes for them. Afterward, we re-  
 157 evaluate all documents against the new code. As the coding exercise is carried out by one author, they  
 158 do a third pass over all documents to double-check for potential errors. Finally, where relevant, we  
 159 cluster the codes into larger themes. In this analysis we only look at the explicit statements provided  
 160 by the authors, we do not attempt to infer their understanding of the problem. Thus, the numbers  
 161 provided in Section 4 should be understood as describing how algorithmic recourse is *discussed* in  
 162 the literature. For brevity, we focus our discussion on the main themes, but we still highlight specific  
 163 publications if we observe that the authors introduce novel, highly relevant considerations that do not  
 164 fit into other themes. Finally, even though we also evaluated the technical aspects of the proposed  
 165 solutions – requirements for methods and datasets used in evaluations – they are not covered in this  
 166 review. Instead, we point the interested readers to [48, 70, 143].

## 167 4 Results

168 The following nine sections introduce the results of the thematic analysis. For each question, we  
169 explain why it is relevant to the analysis and examine the main themes. We also highlight highly  
170 important but underexplored themes. We start with the general points such as contributions and  
171 definitions in Sections 4.1 to 4.3. Then, in Sections 4.4 to 4.7 we investigate the societal components  
172 of AR research. Finally, in Sections 4.8 and 4.9 we look at the aspects relevant to practitioners.

### 173 4.1 What types of contributions do the authors choose to make to the AR research?

174 We start by looking at the main goals of the collected publications to validate our assumption that  
175 AR literature is primarily concerned with technical solutions. We annotate each entry with at most  
176 two codes based on the form of contributions. By far the largest group is *propose methods*, which  
177 applies to 88 (69.3%) out of the 127 publications. These are primarily generators for individual CEs,  
178 but we also find 18 (14.2%) documents that propose other methods. Next, 20 (15.7%) publications  
179 *develop theoretical frameworks*, for instance by grounding AR in user studies or providing critical  
180 perspectives on the problem. Further, 15 (11.8%) focus on *empirical or theoretical analyses* of the  
181 properties of AR and another 15 publications *apply* it in a variety of domains. We did not identify  
182 any applications evaluated with humans in the loop. Then, 5 (3.9%) publications *benchmark* existing  
183 methods, while 3 (2.4%) *review* them. We make our annotations available in technical Appendix B.

### 184 4.2 What are the criteria covered in the authors' definitions of AR?

185 We also evaluate what is understood as the problem to be addressed by AR mechanisms. In particular,  
186 what are the criteria to satisfy authors' definitions of recourse. A similar question was posed by [70]  
187 who combined six definitions into "*recourse can be achieved by an affected individual if they can*  
188 *understand and accordingly act to alleviate an unfavorable situation, thus exercising temporally-*  
189 *extended agency*", but this approach was far from systematic. Instead, we are interested in the  
190 underlying concepts. 74 (58.3%) publications explicitly define AR, 16 (12.6%) mention it but do not  
191 include a definition, while 37 (29.1%) do not mention AR, even though they align with its (overall)  
192 goals. The most common theme is *overturning undesirable decisions*, present in 47 definitions (63.5%  
193 of all definitions), but specifically *overturning algorithmic decisions* is mentioned only 43 (58.1%  
194 times). It is generally understood that AR is *provided to affected individuals* (44, or 59.5%) but 4 (5.4%  
195 definitions *consider stakeholders* more broadly. *Actionability* as a requirement for recourse is noted  
196 in only 39 (52.7%) definitions. Then, 20 (27.0%) publications specifically mention counterfactual  
197 explanations as means to AR, while 26 (35.1%) include various other technical considerations in the  
198 definitions, such as "changes to actionable input variables" or "desired classes".

199 We also point to several themes that are, interestingly, underrepresented. Only 18 (24.3%) documents  
200 mention *explanation, justification, or understanding of a decision* as the pre-requisite for AR. Next,  
201 10 (13.5%) highlight *future-orientation or other temporal aspects* of the provided recommendations.  
202 Although "*consequential settings*", typically bank lending, are given as examples in nine (12.2%  
203 definitions, they are never explicitly mentioned as the scenarios where recourse ought to be provided,  
204 which may be akin to the "enjoyment of recourse" as defined by [142] where people are aware that  
205 there exists a way to reverse undesirable decisions.<sup>1</sup> 8 publications (10.8%) promote *AR as an ability*.  
206 Finally, only 2 (2.7%) publications require that recourse accounts for the *preferences* of its recipients.

### 207 4.3 What are the criteria covered in the authors' definitions of actionability?

208 As we observe, "actionability" is a concept that underpins AR but we discover that, in general, its  
209 understanding is limited. 91 (71.6%) publications attempt to define what it means (for a CE) to be  
210 actionable. Most commonly, in 48 (52.7%) out of 91 definitions, it is understood as *acting only on*  
211 *directly-mutable features*, 6 (6.6%) distinguish that *features may be indirectly-mutable* but still not  
212 actionable, while 22 (24.2%) also highlight that *feature values may need to be constrained*. Next, 19  
213 (20.9%) definitions rely on a tautology that actionability means *people can take actions*, 11 (12.1%  
214 emphasize that these *actions must be successful or lead to change*, and 3 (3.3%) further require  
215 that they are *aligned with people's real-world objectives*. Only 14 (15.4%) definitions put users

<sup>1</sup>Financial domain dominates the evaluations as well, with 90 of 116 evaluations on non-synthetic data making use of at least one finance-related dataset, most commonly German Credit Data [59] with 51 uses.

216 at the center stage, indicating that actionability *depends on the user or their preferences*, while 2  
217 (2.2%) highlight the *importance of the context* [144, 156], for instance, that the ability to act on a  
218 recommendation may change over time. Importantly, ethical considerations are never mentioned as  
219 the pre-requisite for actionability, but we find some broader discussions about this [e.g., 142].

#### 220 4.4 What is the role of end users? What other stakeholders are envisioned in the AR process?

221 Given that AR is to be implemented in socio-technical systems that include a variety of actors, we  
222 are interested in the types of stakeholders acknowledged in the literature. A total of 105 publications  
223 provide explicit consideration of this type. In general, end users subject to algorithmic decisions  
224 are envisioned to be the recipients of AR, but this is not always the case: it may also be provided to  
225 experts [e.g., 21, 22, 76] or organizations [e.g., 65, 72, 147], which highlights that in some cases AR  
226 may be carried out on behalf of the affected individuals. In any case, 47 (44.8%) publications in the  
227 subset agree that end users should inform actionability, but it is rarely clear *how* these preferences  
228 should be specified. User-friendly (interactive) interfaces are a consideration in only 14 (13.3%)  
229 documents. A total of 29 (27.6%) publications envision domain experts as someone who inform  
230 the recourse process. They are either expected to inform actionability in the AR system or provide  
231 other forms of knowledge, typically in the form of a causal structure. Besides the experts, authors  
232 of 35 (33.3%) papers have discussed a variety of stakeholders. Most commonly system owners  
233 [e.g., 20, 34, 38, 89], but also auditors [e.g., 138, 158], data scientists [e.g., 28, 82], developers [e.g.,  
234 22, 131], practitioners [e.g., 100, 156], regulators [e.g., 28, 120], or even potential attackers [102].

#### 235 4.5 What types of real-world considerations motivate existing research?

236 With the multitude of challenges that stand ahead of real-world AR, we are interested in the considera-  
237 tions that motivate existing work. The main theme we find is *ensuring proper individual actionability*,  
238 which is addressed in 46 (37.4%) of 123 publications relevant to this question. This is typically  
239 achieved with the encoding of user preferences as constraints, but other means include providing  
240 diverse CEs. In fact, *tackling specific desiderata for AR* (beyond actionability) is the second largest  
241 area of research with 28 (22.8%) publications. Various *other technical challenges* are considered  
242 in 24 (19.5%) documents, for example, integrating background knowledge [e.g., 16, 62, 64, 98], or  
243 incorporating feature importance [e.g., 4, 6, 96, 116]. We also find 19 (15.4%) publications that  
244 discuss the problem of *communicating recourse to the end users*. 16 (13.0%) focus on the *dynamics*  
245 *of real-world systems*, typically addressing the robustness of AR [e.g., 75, 91, 93, 137], while 14  
246 (11.4%) look at recourse in *multi-agent systems*. This also relates to *performance considerations*  
247 emphasized in 15 (12.2%) of documents. *Causality* drives research in 14 (11.4%) cases. We also  
248 find several themes that are under-emphasized: only 9 (7.3%) publications are directly *motivated by*  
249 *research in psychology*, while *ethics of AR* are emphasized in only 7 (5.7%) documents.

#### 250 4.6 What types of real-world considerations are seen as challenges for future work?

251 While the previous section looked at the considerations that drive existing research, in this section we  
252 distill the recommendations for *future* research going beyond the improvement of own work, which  
253 are provided in 74 documents. *Causality* is highlighted as a challenge in 22 (29.7%) of them, while  
254 *other technical considerations* are given in 20 (27.0%) cases. These range from robustness [e.g.,  
255 51, 117, 137], support for categorical features [e.g., 36, 157], or distinguishing between valid CEs and  
256 adversarial examples [101]. Next, 19 (25.6%) documents highlight the importance of *ensuring proper*  
257 *individual actionability*, which also relates to *communicating recourse to the end users* (9, or 12.2%)  
258 and *supporting realistic cost functions* (8, or 10.8%). *Ethics of AR* are highlighted in 11 (14.9%)  
259 publications, for example, that AR research may detract from other obligations of model owners  
260 [77, 133]. The same number of publications emphasize the need to (1) *ground research in user studies*,  
261 and (2) accommodate for the *dynamics of real-world systems*. *Privacy or security* is highlighted in 10  
262 (13.5%) documents, while the *abuse of recourse*, such as strategic behaviors, surfaces in 7 (9.4%)  
263 papers. Other challenges include improving *performance* (8, or 10.8%), considering *multi-agent*  
264 *systems* (4, or 5.4%), and developing *legal frameworks* (4, or 5.4%) for recourse. We also highlight  
265 several challenges particularly relevant to our work: (the usefulness of) recourse is perceived as  
266 difficult to evaluate in practice [41, 60, 115], it must account for individual, contextual, societal, and  
267 even cultural factors [123], which further means that engagement with recourse mechanisms and the  
268 likelihood of its implementation are context-dependent [e.g., 6, 42, 128].

#### 269 4.7 What types of (emergent) group-level dynamics are addressed in the existing research?

270 Real-world systems entail the implementation of recourse by more than one agent, which may  
271 introduce group-level dynamics. Nonetheless, out of 119 documents relevant to this question, 93  
272 (78.2%) seem to understand recourse as a purely individual phenomenon. Among the remaining  
273 26 documents we find considerations for several different group-level effects. Various perspectives  
274 on the problem of fair AR, covering both individual and group formulations are addressed by  
275 [12, 36, 52, 120, 121, 131, 149, 154]. Next, [9] shows that the implementation of AR on a large scale  
276 may lead to domain and model shifts, which introduce unexpected costs for the stakeholders.<sup>2</sup> In [42]  
277 the authors focus on another negative consequence of AR at scale, showing that it may reinforce  
278 social segregation. The impact of the “right to be forgotten”, where data deletion requests trigger  
279 model retraining that may invalidate existing recourses is addressed in [75]. Then, [94] develop a  
280 game-theoretic framework for AR in multi-agent settings, attempting to optimize for “social welfare”  
281 rather than the profits of individual agents. We find two further similar perspectives on recourse:  
282 [38] proposes auditing and subsidies to minimize the risks of strategic behaviors in a multi-agent  
283 setting, while [136] attempts to incentivize actual improvement for a population of agents. Finally,  
284 [65] provides a framework that generates transparent and consistent recourses for a sub-population.  
285 We also note two other lines of research that account for the remaining documents with group-level  
286 considerations. First, in a causal setting [e.g., 68, 73] subpopulations are necessary to estimate  
287 the interventional effects on individuals. Second, several works highlight the importance of global  
288 insights into the data [22, 41, 44, 78, 108, 112, 152], such as recourse summaries [78, 112].

#### 289 4.8 What are the approaches to the realistic evaluation of proposed methods?

290 We now explore the different forms of “real-world” evaluations, going beyond quantitative experi-  
291 ments, which are present in 51 publications. Most commonly, in 28 (54.9%) of those, the authors  
292 make use of *case studies* presenting the methods in an end-to-end manner. Among those, the appli-  
293 cation of recourse in the `Hired.com` marketplace goes furthest in simulating real-world conditions  
294 for AR [89], but the recommendations are still not evaluated with humans in the loop. Further, 9  
295 (17.6%) documents include other forms of *short walk-through examples*. We also identify 14 (27.5%)  
296 papers that evaluate the methods with *user experiments*, 10 of which involve non-expert users and  
297 4 involve expert users. While we do not observe any interviews with non-expert users, we find 1  
298 (2.0%) publication where *experts are interviewed* [22]. *Other involvement of non-experts* applies to  
299 [116], where they inform the development of methods. *Other involvement of experts* is featured in  
300 two documents where they evaluated the outputs of methods [25, 132]. Altogether, end users were  
301 involved in 17 publications, which is only 13.3% of all publications covered in our study, even more  
302 striking than the 21% of CE methods evaluated with user studies as reported in [71].

#### 303 4.9 What are the open source and documentation practices in AR research?

304 Finally, we note that the lack of availability of well-documented open-source code may be an important  
305 obstacle to the application of AR in real-world systems. For all 116 publications that involve some  
306 form of computational experiments, we verify whether the source code is publicly available. If the  
307 authors do not explicitly link to their code in the paper, we attempt to find it independently. Ultimately,  
308 we collect open-source implementations for 64 (55.2%) publications. Then, for each of them, we  
309 evaluate the quality of documentation. The *instructions on the general usage* (such as installation and  
310 workflow) are provided with 27 (41.5%) repositories, while *instructions on the reproduction of results*  
311 in 23 (35.4%). In 19 (29.2%) cases we find *walk-through tutorials*, typically in the form of Jupyter  
312 Notebooks, although we note that they differ in quality. For instance, 5 repositories include code-only  
313 notebooks with no further textual explanation that could guide the practitioner. Implementations  
314 for 4 papers include more “professionalized” *documentation* [9, 86, 100, 156]. The latter sets a  
315 golden standard as it further includes a tutorial video and a live demo. We do not find *any* additional  
316 materials for practitioners for 13 (20.0%) of the available implementations.

---

<sup>2</sup>Such “endogenous dynamics” were postulated earlier in the first version of [113] dated December 22<sup>nd</sup> 2020, but this discussion has been completely removed from the subsequent versions of the pre-print.

## 317 5 Discussion

318 Regardless of whether AR can be normatively expected or not [77], many systems can genuinely  
319 benefit from recourse mechanisms, especially when the interests of the system owner and the end users  
320 are aligned [72], such as in the healthcare system to improve the well-being of patients [76, 96, 155],  
321 or on the online platforms that attempt to improve the experience of their users [89, 134]. Nonetheless,  
322 the values and norms underlying recourse – trust, agency, fairness, safety, and so on – are emergent  
323 properties of systems where recourse mechanisms would be introduced. Such norms can only be  
324 understood and evaluated when accounting for the technical, social, and institutional components of  
325 the system [32], but the latter two remain largely unexplored in the recourse literature.

326 Recourse is not inherently safe or unsafe, *but* its (incorrect) implementation may lead to the emer-  
327 gence of unsafe dynamics, such as the unexpected costs to stakeholders as discussed by [9] or the  
328 reinforcement of social segregation addressed in [42]. While it may be too challenging to provide  
329 accurate system-level evaluations at this stage of research, authors can still expand the boundaries  
330 of their analyses to account for global effects or look at the position of recourse mechanisms in the  
331 broader context of a complete socio-technical AI system [33]. As AR is a “reality-centric AI” problem  
332 [140] by its nature, working towards its integration into existing systems will require a design-oriented  
333 approach, potentially with *specific* systems in mind. The “Abstraction Traps” discussed by [119] in  
334 the context of research on fair machine learning apply here: that technical solutions designed for one  
335 social context cannot be directly repurposed for another application, that values to which they are  
336 expected to adhere to cannot be captured with mathematical formulas, that their insertion into an  
337 existing process will impact its behavior, or that the best solutions may not necessarily be technical.

338 It is perhaps most telling that only 12% of surveyed publications attempt to apply recourse in realistic  
339 settings. We will discuss two of these settings to highlight the stark differences in system properties.  
340 Most of the applications included in our review focus on the provision of actionable individual  
341 recommendations to students [3, 4, 24, 109, 126, 135, 160]. In this relatively low-stakes domain  
342 almost any recourse will be actionable in that following a personalized set of learning activities  
343 does not require any resources other than time. Even then, the system involves multiple actors  
344 – students, teachers, parents – whose interactions will impact the process, for example, because  
345 students may fail to benefit from certain learning activities without additional support. Conversely,  
346 we find several publications where authors attempt to provide recourse in the high-stakes medical  
347 domain [76, 96, 155]. Here, recommendations must be tailored to the preferences, resources, or  
348 lifestyles of patients in order to have a chance of being actionable. Moreover, certain aspects of their  
349 implementation fully rely on other actors, such as a clinician prescribing the medications. Finally, it  
350 may happen that recourse does not exist at all when the outcomes of a patient cannot be improved.

### 351 5.1 Recommendations for future research

352 We distill our findings into five key recommendations. First, in Sections 4.2, 4.3 we observed that  
353 *operational* definitions for recourse are still unavailable. Second, Sections 4.4 and 4.8 underlined  
354 little consideration for people involved in recourse processes. Third, Sections 4.5, 4.6 highlighted the  
355 overwhelmingly technical approaches to recourse. Fourth, Section 4.7 stressed the lack of group-level  
356 analyses. Fifth, from Sections 4.8, and 4.9 we learned about the missing consideration of practitioners.

357 **1. Broadening the scope of research.** AR is generally seen as a service for affected individuals,  
358 but this formalization may be unnecessarily limiting. In fact, in many systems, these individuals may  
359 be unable to directly act on recommendations [see also 142]. Instead, we propose to operationalize  
360 the aim of AR as the provision of recommendations *aligned with the preferences of non-expert users*  
361 in an attempt to *help them improve outcomes* in an *ADM setting*, which emphasizes that providing  
362 *easy to understand* and *individually actionable* recommendations remains the key research problem.

363 **2. Engaging end users, affected individuals, and communities.** AR solutions are rarely evaluated  
364 with humans. Instead, they attempt to satisfy a variety of desiderata formulated by authors and  
365 assessed in an automated manner. Sparsity, proximity, or mutability of features are far from perfect  
366 proxies for individual actionability. For AR to be truly useful, it must be able to satisfy the preferences  
367 of its end users. Research is also necessary to learn about the needs of the affected individuals  
368 concerning recourse, and to validate its potential contributions and inherent limitations. Authors may  
369 also benefit from the rich literature on human-computer interaction [e.g., 11, 23] or psychology.



370 **3. Accepting a socio-technical perspective.** A pervasive assumption in the literature is that all  
371 challenges of AR require purely technical solutions. For instance, many authors emphasize the  
372 importance of causal modeling to guarantee recourse, but the models that aim to be explained are  
373 themselves *not* causal. Similarly, to improve the performance of CE generators many authors turn to  
374 deep generative models [35, 42, 61, 67, 81, 90, 99]. Not only do they explain the data rather than the  
375 model [10], but more importantly they shift the problem from improving the trust in non-interpretable  
376 models, to attempting to trust non-interpretable explainers. Although a socio-technical perspective  
377 on AR brings its own challenges, such as accounting for the roles of stakeholders involved in the  
378 provision of recourse, it creates important opportunities. For example, developing “recourse contracts”  
379 [34, 39] or designing feedback processes to account for imperfect robustness.

380 **4. Accounting for emergent effects.** Decision-making systems involve multiple individuals who  
381 may be interested in receiving recourse and may have competing interests. Research on AR should,  
382 from the onset, explore group-level effects such as external costs or fairness. While this may require  
383 expanding the boundaries of analysis, it is necessary to anticipate the emergent outcomes of recourse.  
384 These may even occur due to the multi-system dynamics of AR: recommendations implemented by  
385 an individual to improve their outcomes in one system will affect them in other contexts [see also 13].

386 **5. Attending to other operational aspects.** Finally, the artifacts of AR research should be  
387 practitioner-friendly. On the one hand, this requires being explicit about the position of the proposed  
388 methods in a broader system, for example, in the form of end-to-end case studies that allow practi-  
389 tioners to better understand the benefits of the proposed solutions. On the other hand, this suggests  
390 that authors should attempt to move away from merely providing scripts for experiments, and focus  
391 on developing well-documented frameworks that can be adapted to different ADM systems.

## 392 5.2 Limitations of our work

393 Our review is not without shortcomings. Most importantly, for each paper the extraction and coding  
394 of data was performed by a single author, which means that the quantitative results may be imperfect.  
395 We account for this by focusing the analysis on the *overarching themes* represented in existing  
396 publications, thus, even if another researcher would have carried out the coding in a somewhat  
397 different manner, they should arrive at similar results and our analysis remains valid. Additionally, as  
398 our review ultimately looks at the authors’ perception of recourse, we do not want to misconstrue  
399 their views. Thus, we do not infer any considerations unless they are provided explicitly. Our reading  
400 may be more strict than intended by the authors and the numbers reported in our results may be  
401 underestimated. At the same time, we believe that if certain considerations are deemed important  
402 by the researchers, they would choose to be explicit about them. Finally, although we followed a  
403 systematic process, we cannot claim that we collected AR literature in an exhaustive manner due to  
404 the specificities of computer science publishing. Thus, we acknowledge that there may exist some  
405 insightful publications addressing recourse that have not been covered in this literature review.

## 406 6 Conclusions

407 Algorithmic recourse concerns the provision of recommendations aligned with the preferences of  
408 non-expert users of algorithmic decision-making systems to help them achieve more desirable out-  
409 comes in the future. Existing research on the topic is predominantly theoretical, even though recourse,  
410 in expectation, is a real-world problem with strong practical implications. To that end, we conducted  
411 a systematized literature review of 127 publications that focus on algorithmic recourse, and more gen-  
412 erally on actionable counterfactual explanations. We evaluated the practical considerations provided  
413 by the authors. Our findings indicate that, indeed, AR tends to be perceived as a (predominantly)  
414 technical problem. Although we think highly of fundamental research, we note that for algorithmic  
415 recourse to leave computer science labs, it must be more strongly grounded and validated in the real  
416 world, and consider the requirements for systems that include not only technical but also social and  
417 institutional components. To help bridge this gap, we synthesize a list of five recommendations for  
418 other authors that aim to reinforce recourse as a practical problem. We believe that AR should not be  
419 seen as only a simple ad-hoc solution to improve the acceptance of black-box models in consequential  
420 domains, but rather as a full-fledged socio-technical mechanism that can benefit many systems and  
421 improve the agency of affected individuals and decision-makers across a variety of settings.

## 422 References

- 423 [1] Abubakar Abid, Mert Yuksekgonul, and James Zou. Meaningfully Debugging Model Mistakes  
424 using Conceptual Counterfactual Explanations. In *Proceedings of the 39th International*  
425 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,  
426 pages 66–88. PMLR, 17–23 Jul 2022. URL [https://proceedings.mlr.press/v162/a](https://proceedings.mlr.press/v162/abid22a.html)  
427 [bid22a.html](https://proceedings.mlr.press/v162/abid22a.html).
- 428 [2] Gediminas Adomavicius and Alexander Tuzhilin. Discovery of Actionable Patterns in  
429 Databases: The Action Hierarchy Approach. In *Proceedings of the Third International*  
430 *Conference on Knowledge Discovery and Data Mining*, KDD’97, page 111–114. AAAI Press,  
431 1997.
- 432 [3] Farzana Afrin, Margaret Hamilton, and Charles Thevathyan. Exploring Counterfactual Ex-  
433 planations for Predicting Student Success. In *Computational Science – ICCS 2023*, volume  
434 14074 LNCS, pages 413–420. Springer Nature Switzerland, 2023. doi: 10.1007/978-3-031-3  
435 6021-3\_44.
- 436 [4] Muhammad Afzaal, Jalal Nouri, Aayesha Zia, Panagiotis Papapetrou, Uno Fors, Xiu Wu,  
437 Yongchaoand Li, and Rebecka Weegar. Automatic and Intelligent Recommendations to  
438 Support Students’ Self-Regulation. In *2021 International Conference on Advanced Learning*  
439 *Technologies (ICALT)*, pages 336–338, July 2021. ISBN 2161-377X. doi: 10.1109/ICALT522  
440 72.2021.00107.
- 441 [5] Charu C. Aggarwal, Chen Chen, and Jiawei Han. The Inverse Classification Problem. *Journal*  
442 *of Computer Science and Technology*, 25:458–468, 2010.
- 443 [6] Emanuele Albini, Jason Long, Danial Dervovic, and Daniele Magazzeni. Counterfactual  
444 Shapley Additive Explanations. In *Proceedings of the 2022 ACM Conference on Fairness,*  
445 *Accountability, and Transparency*, FAccT ’22, page 1054–1070, New York, NY, USA, 2022.  
446 Association for Computing Machinery. ISBN 9781450393522. doi: 10.1145/3531146.3533168.  
447 URL <https://doi.org/10.1145/3531146.3533168>.
- 448 [7] Kars Alfrink, Ianus Keller, Gerd Kortuem, and Neelke Doorn. Contestable AI by design:  
449 Towards a framework. *Minds and Machines*, pages 1–27, 2022.
- 450 [8] Hissah Alotaibi and Ronal Singh. Metrics for Evaluating Actionability in Explainable AI. In  
451 *PRICAI 2023: Trends in Artificial Intelligence*, pages 481–487. Springer Nature Singapore,  
452 2023. ISBN 978-981-99-7022-3.
- 453 [9] Patrick Altmeyer, Giovan Angela, Aleksander Buszydlík, Karol Dobiczek, Arie van Deursen,  
454 and Cynthia C. S. Liem. Endogenous Macrodynamics in Algorithmic Recourse. In *2023 IEEE*  
455 *Conference on Secure and Trustworthy Machine Learning (SaTML)*, pages 418–431, 2023.  
456 doi: 10.1109/SaTML54575.2023.00036.
- 457 [10] Patrick Altmeyer, Mojtaba Farmanbar, Arie van Deursen, and Cynthia C. S. Liem. Faithful  
458 Model Explanations through Energy-Constrained Conformal Counterfactuals. In *Proceedings*  
459 *of the AAAI Conference on Artificial Intelligence*, volume 38, pages 10829–10837, 2024.
- 460 [11] Saleema Amershi, Maya Cakmak, William Bradley Knox, and Todd Kulesza. Power to the  
461 People: The Role of Humans in Interactive Machine Learning. *AI Magazine*, 35(4):105–120,  
462 2014.
- 463 [12] André Artelt, Valerie Vaquet, Riza Velioglu, Fabian Hinder, Johannes Brinkrolf, Malte  
464 Schilling, and Barbara Hammer. Evaluating Robustness of Counterfactual Explanations. In  
465 *2021 IEEE Symposium Series on Computational Intelligence (SSCI)*, pages 01–09, December  
466 2021. doi: 10.1109/SSCI50451.2021.9660058.
- 467 [13] Solon Barocas, Andrew D. Selbst, and Manish Raghavan. The Hidden Assumptions Behind  
468 Counterfactual Explanations and Principal Reasons. In *Proceedings of the 2020 Conference on*  
469 *Fairness, Accountability, and Transparency*, FAT\* ’20, page 80–89, New York, NY, USA, 2020.  
470 Association for Computing Machinery. ISBN 9781450369367. doi: 10.1145/3351095.3372830.  
471 URL <https://doi.org/10.1145/3351095.3372830>.

- 472 [14] Hosein Barzekar and Susan McRoy. Achievable Minimally-Contrastive Counterfactual  
473 Explanations. *Machine Learning and Knowledge Extraction*, 5(3):922–936, 2023. doi:  
474 10.3390/make5030048.
- 475 [15] Sander Beckers. Causal Explanations and XAI. In Bernhard Schölkopf, Caroline Uhler, and  
476 Kun Zhang, editors, *Proceedings of the First Conference on Causal Learning and Reasoning*,  
477 volume 177 of *Proceedings of Machine Learning Research*, pages 90–109. PMLR, 11–13 Apr  
478 2022. URL <https://proceedings.mlr.press/v177/beckers22a.html>.
- 479 [16] Alexander Berman, Ellen Breitholtz, Christine Howes, and Jean-Philippe Bernardy. Explaining  
480 Predictions with Enthymematic Counterfactuals. In *CEUR Workshop Proceedings*, volume  
481 3319, pages 95–100. CEUR-WS, 2022.
- 482 [17] Umang Bhatt, Alice Xiang, Shubham Sharma, Adrian Weller, Ankur Taly, Yunhan Jia, Joy-  
483 deep Ghosh, Ruchir Puri, José M. F. Moura, and Peter Eckersley. Explainable Machine  
484 Learning in Deployment. In *Proceedings of the 2020 Conference on Fairness, Accountability,  
485 and Transparency*, FAT\* ’20, page 648–657. New York, NY, USA, 2020. Association for  
486 Computing Machinery. ISBN 9781450369367. doi: 10.1145/3351095.3375624. URL  
487 <https://doi.org/10.1145/3351095.3375624>.
- 488 [18] Abeba Birhane, Pratyusha Kalluri, Dallas Card, William Agnew, Ravit Dotan, and Michelle  
489 Bao. The Values Encoded in Machine Learning Research. In *Proceedings of the 2022 ACM  
490 Conference on Fairness, Accountability, and Transparency*, FAccT ’22, page 173–184, New  
491 York, NY, USA, 2022. Association for Computing Machinery. ISBN 9781450393522. doi:  
492 10.1145/3531146.3533083. URL <https://doi.org/10.1145/3531146.3533083>.
- 493 [19] Miguel Á. Carreira-Perpiñán and Suryabhan Singh Hada. Counterfactual Explanations for  
494 Oblique Decision Trees: Exact, Efficient Algorithms. *Proceedings of the AAAI Conference  
495 on Artificial Intelligence*, 35:6903–6911, May 2021. doi: 10.1609/aaai.v35i8.16851. URL  
496 <https://ojs.aaai.org/index.php/AAAI/article/view/16851>.
- 497 [20] Yatong Chen, Jialu Wang, and Yang Liu. Strategic Recourse in Linear Classification. In  
498 *Workshop on Consequential Decision Making in Dynamic Environments*, 2020.
- 499 [21] Ziheng Chen, Fabrizio Silvestri, Gabriele Tolomei, Jia Wang, He Zhu, and Hongshik Ahn.  
500 Explain the Explainer: Interpreting Model-Agnostic Counterfactual Explanations of a Deep  
501 Reinforcement Learning Agent. *IEEE Transactions on Artificial Intelligence*, 5(4):1443–1457,  
502 2024. doi: 10.1109/TAI.2022.3223892.
- 503 [22] Furui Cheng, Yao Ming, and Huamin Qu. DECE: Decision Explorer with Counterfactual  
504 Explanations for Machine Learning Models. *IEEE Transactions on Visualization & Computer  
505 Graphics*, 27(02):1438–1447, feb 2021. ISSN 1941-0506. doi: 10.1109/TVCG.2020.3030342.
- 506 [23] Hao-Fei Cheng, Ruotong Wang, Zheng Zhang, Fiona O’Connell, Terrance Gray, F. Maxwell  
507 Harper, and Haiyi Zhu. Explaining Decision-Making Algorithms through UI: Strategies  
508 to Help Non-Expert Stakeholders. In *Proceedings of the 2019 CHI Conference on Human  
509 Factors in Computing Systems*, CHI ’19, page 1–12, New York, NY, USA, 2019. Association  
510 for Computing Machinery. ISBN 9781450359702. doi: 10.1145/3290605.3300789. URL  
511 <https://doi.org/10.1145/3290605.3300789>.
- 512 [24] Lea Cohausz. Towards Real Interpretability of Student Success Prediction Combining Methods  
513 of XAI and Social Science. In *Proceedings of the 15th International Conference on Educational  
514 Data Mining*, pages 361–367, Durham, United Kingdom, July 2022. International Educational  
515 Data Mining Society. ISBN 978-1-7336736-3-1. doi: 10.5281/zenodo.6853069.
- 516 [25] Riccardo Crupi, Alessandro Castelnovo, Daniele Regoli, and Beatriz San Miguel Gonzalez.  
517 Counterfactual Explanations as Interventions in Latent Space. *Data Mining and Knowledge  
518 Discovery*, 2022. doi: 10.1007/s10618-022-00889-2.
- 519 [26] Susanne Dandl, Christoph Molnar, Martin Binder, and Bernd Bischl. Multi-Objective Counter-  
520 factual Explanations. In *Parallel Problem Solving from Nature – PPSN XVI*, pages  
521 448–469, Cham, 2020. Springer International Publishing. ISBN 978-3-030-58112-1. doi:  
522 10.1007/978-3-030-58112-1\_3.

- 523 [27] Debanjan Datta, Feng Chen, and Naren Ramakrishnan. Framing Algorithmic Recourse for  
524 Anomaly Detection. In *Proceedings of the 28th ACM SIGKDD Conference on Knowledge*  
525 *Discovery and Data Mining*, KDD '22, page 283–293, New York, NY, USA, 2022. Association  
526 for Computing Machinery. ISBN 9781450393850. doi: 10.1145/3534678.3539344. URL  
527 <https://doi.org/10.1145/3534678.3539344>.
- 528 [28] Randall Davis, Andrew W. Lo, Sudhanshu Mishra, Arash Nourian, Manish Singh, Nicholas  
529 Wu, and Ruixun Zhang. Explainable Machine Learning Models of Consumer Credit Risk.  
530 *Journal of Financial Data Science*, 5(4):9–39, 2022. doi: 10.3905/jfds.2023.1.141.
- 531 [29] Marcelo de Sousa Balbino, Luis Enrique Zárata Gálvez, and Cristiane Neri Nobre. CSSE  
532 - An agnostic method of counterfactual, selected, and social explanations for classification  
533 models. *Expert Systems with Applications*, 228:120373, 2023. ISSN 0957-4174. doi: <https://doi.org/10.1016/j.eswa.2023.120373>. URL [https://www.sciencedirect.com/science](https://www.sciencedirect.com/science/article/pii/S0957417423008758)  
534 [article/pii/S0957417423008758](https://www.sciencedirect.com/science/article/pii/S0957417423008758).  
535
- 536 [30] Giovanni De Toni, Bruno Lepri, and Andrea Passerini. Synthesizing explainable counterfactual  
537 policies for algorithmic recourse with program synthesis. *Machine Learning*, 112(4):1389–  
538 1409, 2023. ISSN 0885-6125. doi: 10.1007/s10994-022-06293-7.
- 539 [31] Sarah Dean, Sarah Rich, and Benjamin Recht. Recommendations and User Agency: The  
540 Reachability of Collaboratively-Filtered Information. In *Proceedings of the 2020 Conference*  
541 *on Fairness, Accountability, and Transparency*, FAT\* '20, page 436–445, New York, NY, USA,  
542 2020. Association for Computing Machinery. ISBN 9781450369367. doi: 10.1145/3351095.  
543 3372866. URL <https://doi.org/10.1145/3351095.3372866>.
- 544 [32] Roel Dobbe and Anouk Wolters. Toward Sociotechnical AI: Mapping Vulnerabilities for  
545 Machine Learning in Context. *Minds and Machines*, 34(2):1–51, 2024.
- 546 [33] Roel Dobbe, Thomas Krendl Gilbert, and Yonatan Mintz. Hard choices in artificial intelligence.  
547 *Artificial Intelligence*, 300:103555, 2021. ISSN 0004-3702. doi: <https://doi.org/10.1016/j.arti>  
548 [nt.2021.103555](https://doi.org/10.1016/j.arti). URL [https://www.sciencedirect.com/science/article/pii/S0](https://www.sciencedirect.com/science/article/pii/S004370221001065)  
549 [004370221001065](https://www.sciencedirect.com/science/article/pii/S004370221001065).
- 550 [34] Ricardo Dominguez-Olmedo, Amir-Hossein Karimi, and Bernhard Schölkopf. On the Adver-  
551 sarial Robustness of Causal Algorithmic Recourse. In *Proceedings of the 39th International*  
552 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,  
553 pages 5324–5342. PMLR, 17–23 2022.
- 554 [35] Michael Downs, Jonathan L. Chu, Yaniv Yacoby, Finale Doshi-Velez, and Weiwei Pan.  
555 CRUDS: Counterfactual Recourse Using Disentangled Subspaces. *ICML Workshop on Human*  
556 *Interpretability in Machine Learning*, pages 1–23, 2020.
- 557 [36] Ahmad-Reza Ehyaei, Amir-Hossein Karimi, Bernhard Schoelkopf, and Setareh Maghsudi.  
558 Robustness Implies Fairness in Causal Algorithmic Recourse. In *Proceedings of the 2023*  
559 *ACM Conference on Fairness, Accountability, and Transparency*, FAccT '23, page 984–1001,  
560 New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701924.  
561 doi: 10.1145/3593013.3594057. URL <https://doi.org/10.1145/3593013.3594057>.
- 562 [37] Julia El Zini and Mariette Awad. Beyond Model Interpretability: On the Faithfulness and  
563 Adversarial Robustness of Contrastive Textual Explanations. In *Findings of the Association for*  
564 *Computational Linguistics: EMNLP 2022*, pages 1391–1402. Association for Computational  
565 Linguistics, 2022. doi: 10.18653/v1/2022.findings-emnlp.100.
- 566 [38] Andrew Estornell, Yatong Chen, Sanmay Das, Yang Liu, and Yevgeniy Vorobeychik. Incenti-  
567 vizing Recourse through Auditing in Strategic Classification. In *Proceedings of the Thirty-*  
568 *Second International Joint Conference on Artificial Intelligence, IJCAI-23*, pages 400–408.  
569 International Joint Conferences on Artificial Intelligence, 8 2023. doi: 10.24963/ijcai.2023/45.  
570 URL <https://doi.org/10.24963/ijcai.2023/45>.
- 571 [39] Andrea Ferrario and Michele Loi. The Robustness of Counterfactual Explanations Over Time.  
572 *IEEE Access*, 10:82736–82750, 2022. ISSN 2169-3536. doi: 10.1109/ACCESS.2022.3196917.

- 573 [40] Susanne Friese, Jacks Soratto, and Denise Pires de Pires. Carrying out a computer-aided  
574 thematic content analysis with ATLAS.ti. *IWMI Working Papers*, 18, 04 2018.
- 575 [41] Sainyam Galhotra, Romila Pradhan, and Babak Salimi. Explaining Black-Box Algorithms  
576 Using Probabilistic Contrastive Counterfactuals. In *Proceedings of the 2021 International  
577 Conference on Management of Data, SIGMOD '21*, pages 577–590, New York, NY, USA,  
578 2021. Association for Computing Machinery. ISBN 978-1-4503-8343-1. doi: 10.1145/344801  
579 6.3458455.
- 580 [42] Ruijiang Gao and Himabindu Lakkaraju. On the Impact of Algorithmic Recourse on Social  
581 Segregation. In *Proceedings of the 40th International Conference on Machine Learning,  
582 ICML'23*. JMLR.org, 2023.
- 583 [43] Azin Ghazimatin, Oana Balalau, Rishiraj Saha Roy, and Gerhard Weikum. PRINCE: Provider-  
584 Side Interpretability with Counterfactual Explanations in Recommender Systems. In *Pro-  
585 ceedings of the 13th International Conference on Web Search and Data Mining, WSDM '20*,  
586 pages 196–204, New York, NY, USA, 2020. Association for Computing Machinery. ISBN  
587 978-1-4503-6822-3. doi: 10.1145/3336191.3371824.
- 588 [44] Oscar Gomez, Steffen Holter, Jun Yuan, and Enrico Bertini. ViCE: Visual Counterfactual Ex-  
589 planations for Machine Learning Models. In *Proceedings of the 25th International Conference  
590 on Intelligent User Interfaces, IUI '20*, pages 531–535, New York, NY, USA, 2020. Association  
591 for Computing Machinery. ISBN 978-1-4503-7118-6. doi: 10.1145/3377325.3377536.
- 592 [45] Crystal Grant. Algorithms Are Making Decisions About Health Care, Which May Only  
593 Worsen Medical Racism, October 2022. URL [https://www.aclu.org/news/privacy-t  
594 echnology/algorithms-in-health-care-may-worsen-medical-racism](https://www.aclu.org/news/privacy-technology/algorithms-in-health-care-may-worsen-medical-racism). Accessed  
595 22.05.2024.
- 596 [46] Maria J. Grant and Andrew Booth. A typology of reviews: an analysis of 14 review types and  
597 associated methodologies. *Health Information & Libraries Journal*, 26(2):91–108, 2009. doi:  
598 <https://doi.org/10.1111/j.1471-1842.2009.00848.x>.
- 599 [47] Stephan Grimmelikhuijsen and Albert Meijer. Legitimacy of Algorithmic Decision-Making:  
600 Six Threats and the Need for a Calibrated Institutional Response. *Perspectives on Public  
601 Management and Governance*, 5(3):232–242, 03 2022. ISSN 2398-4910. doi: 10.1093/ppmg  
602 ov/gvac008. URL <https://doi.org/10.1093/ppmgov/gvac008>.
- 603 [48] Riccardo Guidotti. Counterfactual Explanations and How to Find Them: Literature Review  
604 and Benchmarking. *Data Mining and Knowledge Discovery*, 2022. doi: 10.1007/s10618-022  
605 -00831-6.
- 606 [49] Riccardo Guidotti and Salvatore Ruggieri. Ensemble of Counterfactual Explainers. In  
607 *Discovery Science: 24th International Conference, DS 2021, Halifax, NS, Canada, October  
608 11–13, 2021, Proceedings*, pages 358–368, Berlin, Heidelberg, 2021. Springer-Verlag. ISBN  
609 978-3-030-88941-8. doi: 10.1007/978-3-030-88942-5\_28.
- 610 [50] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Francesca Naretto, Franco Turini,  
611 Dino Pedreschi, and Fosca Giannotti. Stable and Actionable Explanations of Black-Box  
612 Models through Factual and Counterfactual Rules. *Data Mining and Knowledge Discovery*,  
613 2022. doi: 10.1007/s10618-022-00878-5.
- 614 [51] Hangzhi Guo, Feiran Jia, Jinghui Chen, Anna Squicciarini, and Amulya Yadav. RoCourseNet:  
615 Robust Training of a Prediction Aware Recourse Model. In *Proceedings of the 32nd ACM  
616 International Conference on Information and Knowledge Management, CIKM '23*, pages 619–  
617 628, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701245.  
618 doi: 10.1145/3583780.3615040.
- 619 [52] Vivek Gupta, Pegah Nokhiz, Chitradeep Dutta Roy, and Suresh Venkatasubramanian. Equaliz-  
620 ing Recourse Across Groups. *arXiv*, 2019.

- 621 [53] Victor Guyomard, Françoise Fessant, Tassadit Bouadi, and Thomas Guyet. Post-hoc Counterfactual Generation with Supervised Autoencoder. In *Communications in Computer and Information Science*, volume 1524 CCIS, pages 105–114. Springer Science and Business Media Deutschland GmbH, 2021. doi: 10.1007/978-3-030-93736-2\_10.
- 622  
623  
624
- 625 [54] Victor Guyomard, Françoise Fessant, Thomas Guyet, Tassadit Bouadi, and Alexandre Termier. Generating Robust Counterfactual Explanations. In *Machine Learning and Knowledge Discovery in Databases: Research Track. ECML PKDD 2023*, pages 394–409, Berlin, Heidelberg, 2023. Springer-Verlag. ISBN 978-3-031-43417-4. doi: 10.1007/978-3-031-43418-1\_24.
- 626  
627  
628
- 629 [55] Suryabhan Singh Hada and Miguel Á. Carreira-Perpiñán. Exploring Counterfactual Explanations for Classification and Regression Trees. In *Communications in Computer and Information Science*, volume 1524 CCIS, pages 489–504. Springer Science and Business Media Deutschland GmbH, 2021. doi: 10.1007/978-3-030-93736-2\_37.
- 630  
631  
632
- 633 [56] Aparajita Haldar, Teddy Cunningham, and Hakan Ferhatosmanoglu. RAGUEL: Recourse-Aware Group Unfairness Elimination. In *Proceedings of the 31st ACM International Conference on Information & Knowledge Management, CIKM '22*, pages 666–675, New York, NY, USA, 2022. Association for Computing Machinery. ISBN 978-1-4503-9236-5. doi: 10.1145/3511808.3557424.
- 634  
635  
636  
637
- 638 [57] Ian Hardy, Jayanth Yetukuri, and Yang Liu. Adaptive Adversarial Training Does Not Increase Recourse Costs. In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society, AIES '23*, pages 432–442, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400702310. doi: 10.1145/3600211.3604704.
- 639  
640  
641
- 642 [58] Zhian He and Eric Lo. Answering Why-not Questions on Top-k Queries. *2012 IEEE 28th International Conference on Data Engineering*, pages 750–761, 2012. doi: 10.1109/ICDE.2012.8.
- 643  
644
- 645 [59] Hans Hofmann. Statlog (German Credit Data). UCI Machine Learning Repository, 1994. DOI: <https://doi.org/10.24432/C5NC77>.
- 646
- 647 [60] Jacqueline Höllig, Aniek F. Markus, Jjef de Slegte, and Prachi Bagave. Semantic Meaningfulness: Evaluating Counterfactual Approaches for Real-World Plausibility and Feasibility. In *Communications in Computer and Information Science*, volume 1902 CCIS, pages 636–659. Springer Science and Business Media Deutschland GmbH, 2023. doi: 10.1007/978-3-031-44067-0\_32.
- 648  
649  
650  
651
- 652 [61] Shalmali Joshi, Oluwasanmi Koyejo, Warut Vijitbenjaronk, Been Kim, and Joydeep Ghosh. Towards Realistic Individual Recourse and Actionable Explanations in Black-Box Decision Making Systems. *arXiv*, 2019.
- 653  
654
- 655 [62] Sarathi K, Shania Mitra, Deepak P, and Sutanu Chakraborti. Counterfactuals as Explanations for Monotonic Classifiers. In *CEUR Workshop Proceedings*, volume 3389, pages 177–188. CEUR-WS, 2022.
- 656  
657
- 658 [63] Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Hiroki Arimura. Distribution-Aware Counterfactual Explanation by Mixed-Integer Linear Optimization. *Transactions of the Japanese Society for Artificial Intelligence*, 36(6), 2021. doi: 10.1527/TJSAI.36-6\_C-L44.
- 659  
660
- 661 [64] Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, Yuichi Ike, Kento Uemura, and Hiroki Arimura. Ordered Counterfactual Explanation by Mixed-Integer Linear Optimization. In *35th AAAI Conference on Artificial Intelligence, AAAI 2021*, volume 13A, pages 11564–11574. Association for the Advancement of Artificial Intelligence, 2021.
- 662  
663  
664
- 665 [65] Kentaro Kanamori, Takuya Takagi, Ken Kobayashi, and Yuichi Ike. Counterfactual Explanation Trees: Transparent and Consistent Actionable Recourse with Decision Trees. In *Proceedings of The 25th International Conference on Artificial Intelligence and Statistics*, volume 151, pages 1846–1870. PMLR, 2022.
- 666  
667  
668

- 669 [66] Amir-Hossein Karimi, Gilles Barthe, Borja Balle, and Isabel Valera. Model-Agnostic Counterfactual Explanations for Consequential Decisions. In *Proceedings of the Twenty Third International Conference on Artificial Intelligence and Statistics*, volume 108, pages 895–905. PMLR, 2020.
- 670  
671  
672
- 673 [67] Amir-Hossein Karimi, Julius Von Kügelgen, Bernhard Schölkopf, and Isabel Valera. Algorithmic recourse under imperfect causal knowledge: a probabilistic approach. *Advances in Neural Information Processing Systems*, 33:265–277, 2020.
- 674  
675
- 676 [68] Amir-Hossein Karimi, Julius von Kügelgen, Bernhard Schölkopf, and Isabel Valera. Towards Causal Algorithmic Recourse. In *xxAI - Beyond Explainable AI: International Workshop, Held in Conjunction with ICML 2020, July 18, 2020, Vienna, Austria, Revised and Extended Papers*, pages 139–166, Cham, 2020. Springer International Publishing. ISBN 978-3-031-04082-5. doi: 10.1007/978-3-031-04083-2\_8.
- 677  
678  
679  
680
- 681 [69] Amir-Hossein Karimi, Bernhard Schölkopf, and Isabel Valera. Algorithmic Recourse: From Counterfactual Explanations to Interventions. In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency, FAccT '21*, pages 353–362, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 978-1-4503-8309-7. doi: 10.1145/3442188.3445899.
- 682  
683  
684  
685
- 686 [70] Amir-Hossein Karimi, Gilles Barthe, Bernhard Schölkopf, and Isabel Valera. A Survey of Algorithmic Recourse: Contrastive Explanations and Consequential Recommendations. *ACM Computing Surveys*, 55(5), December 2022. ISSN 0360-0300. doi: 10.1145/3527848.
- 687  
688
- 689 [71] Mark T. Keane, Eoin M. Kenny, Eoin Delaney, and Barry Smyth. If Only We Had Better Counterfactual Explanations: Five Key Deficits to Rectify in the Evaluation of Counterfactual XAI Techniques. In Zhi-Hua Zhou, editor, *Proceedings of the Thirtieth International Joint Conference on Artificial Intelligence, IJCAI-21*, pages 4466–4474. International Joint Conferences on Artificial Intelligence Organization, 8 2021. doi: 10.24963/ijcai.2021/609. URL <https://doi.org/10.24963/ijcai.2021/609>. Survey Track.
- 690  
691  
692  
693  
694
- 695 [72] Nwaike Kelechi and Licheng Jiao. Quantifying Actionability: Evaluating Human-Recipient Models. *IEEE Access*, 11:119811–119823, 2023. ISSN 2169-3536. doi: 10.1109/ACCESS.2023.3324906.
- 696  
697
- 698 [73] Gunnar König, Timo Freiesleben, and Moritz Grosse-Wentrup. Causal Perspective on Meaningful and Robust Algorithmic Recourse. *ICML Workshop on Algorithmic Recourse*, 2021.
- 699
- 700 [74] Gunnar König, Timo Freiesleben, and Moritz Grosse-Wentrup. Improvement-Focused Causal Recourse (ICR). In *Proceedings of the Thirty-Seventh AAAI Conference on Artificial Intelligence and Thirty-Fifth Conference on Innovative Applications of Artificial Intelligence and Thirteenth Symposium on Educational Advances in Artificial Intelligence, AAAI'23/IAAI'23/EAAI'23*. AAAI Press, 2023. ISBN 978-1-57735-880-0. doi: 10.1609/aaai.v37i10.26398.
- 701  
702  
703  
704  
705
- 706 [75] Satyapriya Krishna, Jiaqi Ma, and Himabindu Lakkaraju. Towards Bridging the Gaps between the Right to Explanation and the Right to Be Forgotten. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- 707  
708
- 709 [76] Anisio Lacerda, Claudio Almeida, Leonardo Ferreira, Adriano Pereira, Gisele L. Pappa, Wagner Meira, Debora Miranda, Marco A. Romano-Silva, and Leandro Malloy Diniz. Algorithmic Recourse in Mental Healthcare. In *2023 International Joint Conference on Neural Networks (IJCNN)*, pages 1–8, June 2023. ISBN 2161-4407. doi: 10.1109/IJCNN54540.2023.10191158.
- 710  
711  
712
- 713 [77] Derek Leben. Explainable AI as evidence of fair decisions. *Frontiers in Psychology*, 14, 2023. doi: 10.3389/fpsyg.2023.1069426.
- 714
- 715 [78] Dan Ley, Saumitra Mishra, and Daniele Magazzeni. GLOBE-CE: A Translation Based Approach for Global Counterfactual Explanations. In *Proceedings of the 40th International Conference on Machine Learning, ICML'23*. JMLR.org, 2023.
- 716  
717

- 718 [79] Ana Lucic, Harrie Oosterhuis, Hinda Haned, and Maarten de Rijke. FOCUS: Flexible  
719 Optimizable Counterfactual Explanations for Tree Ensembles. In *Proceedings of the 36th*  
720 *AAAI Conference on Artificial Intelligence, AAAI 2022*, volume 36, pages 5313–5322, 2022.
- 721 [80] Shucen Ma, Jianqi Shi, Yanhong Huang, Shengchao Qin, and Zhe Hou. Minimal-unsatisfiable-  
722 core-driven Local Explainability Analysis for Random Forest. *International Journal of*  
723 *Software and Informatics*, 12(4):355–376, 2022. doi: 10.21655/ijsi.1673-7288.00280.
- 724 [81] Divyat Mahajan, Chenhao Tan, and Amit Sharma. Preserving Causal Constraints in Counter-  
725 factual Explanations for Machine Learning Classifiers. In *NeurIPS 2019 Workshop “Do the*  
726 *right thing”: machine learning and causal inference for improved decision making*, 2019.
- 727 [82] Raphael Mazzine, Sofie Goethals, Dieter Brughmans, and David Martens. Counterfactual  
728 Explanations for Employment Services. *International workshop on AI for Human Resources*  
729 *and Public Employment Services*, 2021.
- 730 [83] Md Golam Moula Mehedi Hasan and Douglas A. Talbert. Mitigating the Rashomon Effect in  
731 Counterfactual Explanation: A Game-theoretic Approach. In *Proceedings of the International*  
732 *Florida Artificial Intelligence Research Society Conference, FLAIRS*, volume 35. Florida  
733 Online Journals, University of Florida, 2022. doi: 10.32473/flairs.v35i.130711.
- 734 [84] Tim Miller. Explanation in artificial intelligence: Insights from the social sciences. *Artificial*  
735 *Intelligence*, 267:1–38, 2019. ISSN 0004-3702. doi: <https://doi.org/10.1016/j.artint.2018.07.007>. URL [https://www.sciencedirect.com/science/article/pii/S00043702183](https://www.sciencedirect.com/science/article/pii/S0004370218305988)  
736 [05988](https://www.sciencedirect.com/science/article/pii/S0004370218305988).  
737
- 738 [85] Jonathan Moore, Nils Hammerla, and Chris Watkins. Explaining Deep Learning Models with  
739 Constrained Adversarial Examples. In *PRICAI 2019: Trends in Artificial Intelligence: 16th*  
740 *Pacific Rim International Conference on Artificial Intelligence, Cuvu, Yanuca Island, Fiji,*  
741 *August 26–30, 2019, Proceedings, Part I*, pages 43–56, Berlin, Heidelberg, 2019. Springer-  
742 Verlag. ISBN 978-3-030-29907-1. doi: 10.1007/978-3-030-29908-8\_4.
- 743 [86] Ramaravind K. Mothilal, Amit Sharma, and Chenhao Tan. Explaining Machine Learning  
744 Classifiers through Diverse Counterfactual Explanations. In *Proceedings of the 2020 Confer-*  
745 *ence on Fairness, Accountability, and Transparency, FAT\* ’20*, pages 607–617, New York,  
746 NY, USA, 2020. Association for Computing Machinery. ISBN 978-1-4503-6936-7. doi:  
747 10.1145/3351095.3372850.
- 748 [87] Madhumita Murgia. Algorithms are deciding who gets organ transplants. Are their decisions  
749 fair?, November 2023. URL [https://www.ft.com/content/5125c83a-b82b-40c5-8](https://www.ft.com/content/5125c83a-b82b-40c5-8b35-99579e087951)  
750 [b35-99579e087951](https://www.ft.com/content/5125c83a-b82b-40c5-8b35-99579e087951). Accessed 22.05.2024.
- 751 [88] Philip Naumann and Eirini Ntoutsi. Consequence-Aware Sequential Counterfactual Generation.  
752 In *Machine Learning and Knowledge Discovery in Databases. Research Track: European*  
753 *Conference, ECML PKDD 2021, Bilbao, Spain, September 13–17, 2021, Proceedings, Part II*,  
754 pages 682–698, Berlin, Heidelberg, 2021. Springer-Verlag. ISBN 978-3-030-86519-1. doi:  
755 10.1007/978-3-030-86520-7\_42.
- 756 [89] Daniel Nemirovsky, Nicolas Thiebaut, Ye Xu, and Abhishek Gupta. Providing Actionable  
757 Feedback in Hiring Marketplaces using Generative Adversarial Networks. In *WSDM 2021 -*  
758 *Proceedings of the 14th ACM International Conference on Web Search and Data Mining*, pages  
759 1089–1092. Association for Computing Machinery, 2021. doi: 10.1145/3437963.3441705.
- 760 [90] Daniel Nemirovsky, Nicolas Thiebaut, Ye Xu, and Abhishek Gupta. CounterGAN: Gener-  
761 ating Counterfactuals for Real-Time Recourse and Interpretability using Residual GANs. In  
762 *Proceedings of the 38th Conference on Uncertainty in Artificial Intelligence, UAI 2022*, pages  
763 1488–1497. Association For Uncertainty in Artificial Intelligence (AUAI), 2022.
- 764 [91] Duy Nguyen, Ngoc Bui, and Viet Anh Nguyen. Distributionally Robust Recourse Action.  
765 *arXiv*, 2023.
- 766 [92] Duy Nguyen, Ngoc Bui, and Viet Anh Nguyen. Feasible Recourse Plan via Diverse Interpo-  
767 lation. In *Proceedings of The 26th International Conference on Artificial Intelligence and*  
768 *Statistics*, volume 206, pages 4679–4698. PMLR, 2023.



- 769 [93] Tuan-Duy H. Nguyen, Ngoc Bui, Duy Nguyen, Man-Chung Yue, and Viet Anh Nguyen.  
770 Robust Bayesian Recourse. In *Proceedings of the Thirty-Eighth Conference on Uncertainty in*  
771 *Artificial Intelligence*, volume 180, pages 1498–1508. PMLR, 2022.
- 772 [94] Andrew O’Brien and Edward Kim. Toward Multi-Agent Algorithmic Recourse: Challenges  
773 From a Game-Theoretic Perspective. In *Proceedings of the International Florida Artificial*  
774 *Intelligence Research Society Conference, FLAIRS*, volume 35. Florida Online Journals,  
775 University of Florida, 2022. doi: 10.32473/flairs.v35i.130614.
- 776 [95] Andrew O’Brien, Edward Kim, and Rosina Weber. Investigating Causally Augmented Sparse  
777 Learning as a Tool for Meaningful Classification. In *2023 IEEE Sixth International Conference*  
778 *on Artificial Intelligence and Knowledge Engineering (AIKE)*, pages 33–37, September 2023.  
779 ISBN 2831-7203. doi: 10.1109/AIKE59827.2023.00013.
- 780 [96] Ming Lun Ong, Anthony Li, and Mehul Motani. Explainable and Actionable Machine Learning  
781 Models for Electronic Health Record Data. In *IFMBE Proceedings*, volume 79, pages 91–99,  
782 Cham, 2021. Springer International Publishing. doi: 10.1007/978-3-030-62045-5\_9.
- 783 [97] Matthew J. Page, Joanne E. McKenzie, Patrick M. Bossuyt, Isabelle Boutron, Tammy C.  
784 Hoffmann, Cynthia D. Mulrow, Larissa Shamseer, Jennifer M. Tetzlaff, Elie A. Akl, Sue E.  
785 Brennan, et al. The PRISMA 2020 statement: an updated guideline for reporting systematic  
786 reviews. *Bmj*, 372, 2021.
- 787 [98] Axel Parmentier and Thibaut Vidal. Optimal Counterfactual Explanations in Tree Ensembles.  
788 In *Proceedings of the 38th International Conference on Machine Learning*, volume 139, pages  
789 8422–8431. PMLR, 2021.
- 790 [99] Martin Pawelczyk, Klaus Broelemann, and Gjergji Kasneci. Learning Model-Agnostic Counter-  
791 factual Explanations for Tabular Data. In *The Web Conference 2020 - Proceedings of the World*  
792 *Wide Web Conference, WWW 2020*, pages 3126–3132, 2020. doi: 10.1145/3366423.3380087.
- 793 [100] Martin Pawelczyk, Sascha Bielawski, Johannes van den Heuvel, Tobias Richter, and Gjergji  
794 Kasneci. CARLA: A Python Library to Benchmark Algorithmic Recourse and Counterfactual  
795 Explanation Algorithms. In *Proceedings of the Neural Information Processing Systems Track*  
796 *on Datasets and Benchmarks 1 (NeurIPS Datasets and Benchmarks 2021)*, 2021.
- 797 [101] Martin Pawelczyk, Chirag Agarwal, Shalmali Joshi, Sohini Upadhyay, and Himabindu  
798 Lakkaraju. Exploring Counterfactual Explanations Through the Lens of Adversarial Ex-  
799 amples: A Theoretical and Empirical Analysis. In *Proceedings of The 25th International*  
800 *Conference on Artificial Intelligence and Statistics*, volume 151, pages 4574–4594. PMLR,  
801 2022.
- 802 [102] Martin Pawelczyk, Himabindu Lakkaraju, and Seth Neel. On the Privacy Risks of Algorithmic  
803 Recourse. In *Proceedings of The 26th International Conference on Artificial Intelligence and*  
804 *Statistics*, volume 206, pages 9680–9696. PMLR, 2023.
- 805 [103] Judea Pearl. *Causality*. Cambridge University Press, 2 edition, 2009. ISBN 9780511803161.
- 806 [104] Rafael Poyiadzi, Kacper Sokol, Raul Santos-Rodriguez, Tijn De Bie, and Peter Flach. FACE:  
807 Feasible and actionable counterfactual explanations. In *Proceedings of the AAAI/ACM Con-*  
808 *ference on AI, Ethics, and Society, AIES ’20*, pages 344–350, New York, NY, USA, 2020.  
809 Association for Computing Machinery. doi: 10.1145/3375627.3375850.
- 810 [105] Wenting Qi and Charalampos Chelmis. Improving Algorithmic Decision-Making in the  
811 Presence of Untrustworthy Training Data. In *2021 IEEE International Conference on Big*  
812 *Data (Big Data)*, pages 1102–1108, 2021. doi: 10.1109/BigData52589.2021.9671677.
- 813 [106] Marcos M. Raimundo, Luis Gustavo Nonato, and Jorge Poco. Mining Pareto-optimal Counter-  
814 factual Antecedents with a Branch-and-Bound Model-Agnostic Algorithm. *Data Mining and*  
815 *Knowledge Discovery*, 2022. doi: 10.1007/s10618-022-00906-4.
- 816 [107] Goutham Ramakrishnan, Yun Chan Lee, and Aws Albarghouthi. Synthesizing Action Se-  
817 quences for Modifying Model Decisions. In *Proceedings of the AAAI Conference on Artificial*  
818 *Intelligence*, volume 34, pages 5462–5469, 2020.

- 819 [108] Natraj Raman, Daniele Magazzeni, and Sameena. Shah. Bayesian Hierarchical Models for  
820 Counterfactual Estimation. In *Proceedings of The 26th International Conference on Artificial  
821 Intelligence and Statistics*, volume 206, pages 1115–1128. PMLR, 2023.
- 822 [109] Gomathy Ramaswami, Teo Susnjak, and Anuradha Mathrani. Supporting Students’ Academic  
823 Performance Using Explainable Machine Learning with Automated Prescriptive Analytics.  
824 *Big Data and Cognitive Computing*, 6(4), 2022. doi: 10.3390/bdcc6040105.
- 825 [110] Zbigniew W. Ras and Alicja Wieczorkowska. Action-Rules: How to Increase Profit of a  
826 Company. In *Principles of Data Mining and Knowledge Discovery*, pages 587–592. Springer  
827 Berlin Heidelberg, 2000. ISBN 978-3-540-45372-7.
- 828 [111] Peyman Rasouli and Ingrid Chieh Yu. CARE: Coherent Actionable Recourse Based on Sound  
829 Counterfactual Explanations. *International Journal of Data Science and Analytics*, 17, 2022.  
830 doi: 10.1007/s41060-022-00365-6.
- 831 [112] Kaivalya Rawal and Himabindu Lakkaraju. Beyond Individualized Recourse: Interpretable  
832 and Interactive Summaries of Actionable Recourses. In *Proceedings of the 34th International  
833 Conference on Neural Information Processing Systems, NIPS’20*, Red Hook, NY, USA, 2020.  
834 Curran Associates Inc. ISBN 978-1-71382-954-6.
- 835 [113] Kaivalya Rawal, Ece Kamar, and Himabindu Lakkaraju. Algorithmic Recourse in the Wild:  
836 Understanding the Impact of Data and Model Shifts. *arXiv*, 2021.
- 837 [114] Annabelle Redelmeier, Martin Jullum, Kjersti Aas, and Anders Løland. MCCE: Monte Carlo  
838 sampling of realistic counterfactual explanations. In *Data Mining and Knowledge Discovery*,  
839 pages 421–437. Springer Nature, 2024. doi: 10.1007/s10618-024-01017-y.
- 840 [115] Alexis Ross, Himabindu Lakkaraju, and Osbert Bastani. Learning models for actionable  
841 recourse. In *Advances in Neural Information Processing Systems*, volume 34, pages 18734–  
842 18746, 2021.
- 843 [116] Pedram Salimi, Nirmalie Wiratunga, David Corsar, and Anjana Wijekoon. Towards Feasible  
844 Counterfactual Explanations: A Taxonomy Guided Template-Based NLG Method. In *Frontiers  
845 in Artificial Intelligence and Applications*, volume 372, pages 2057–2064. IOS Press BV, 2023.  
846 doi: 10.3233/FAIA230499.
- 847 [117] Maximilian Schleich, Zixuan Geng, Yihong Zhang, and Dan Suciu. GeCo: Quality Counter-  
848 factual Explanations in Real Time. *Proc. VLDB Endow.*, 14(9):1681–1693, may 2021. ISSN  
849 2150-8097. doi: 10.14778/3461535.3461555. URL <https://doi.org/10.14778/3461535.3461555>.
- 851 [118] Jakob Schoeffler, Niklas Kuehl, and Yvette Machowski. “There Is Not Enough Information”:  
852 On the Effects of Explanations on Perceptions of Informational Fairness and Trustworthiness  
853 in Automated Decision-Making. In *Proceedings of the 2022 ACM Conference on Fairness,  
854 Accountability, and Transparency, FAccT ’22*, pages 1616–1628, New York, NY, USA, 2022.  
855 Association for Computing Machinery. ISBN 978-1-4503-9352-2. doi: 10.1145/3531146.35  
856 33218.
- 857 [119] Andrew D. Selbst, danah boyd, Sorelle A. Friedler, Suresh Venkatasubramanian, and Janet  
858 Vertesi. Fairness and abstraction in sociotechnical systems. In *Proceedings of the Conference on  
859 Fairness, Accountability, and Transparency, FAT\* ’19*, page 59–68, New York, NY, USA, 2019.  
860 Association for Computing Machinery. ISBN 9781450361255. doi: 10.1145/3287560.3287598.  
861 URL <https://doi.org/10.1145/3287560.3287598>.
- 862 [120] Shubham Sharma, Jette Henderson, and Joydeep Ghosh. CERTIFAI: A Common Framework  
863 to Provide Explanations and Analyse the Fairness and Robustness of Black-Box Models.  
864 In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society, AIES ’20*, pages  
865 166–172, New York, NY, USA, 2020. Association for Computing Machinery. ISBN 978-1-  
866 4503-7110-0. doi: 10.1145/3375627.3375812.

- 867 [121] Shubham Sharma, Alan H. Gee, David Paydarfar, and Joydeep Ghosh. FaiR-N: Fair and Robust  
868 Neural Networks for Structured Data. In *Proceedings of the 2021 AAAI/ACM Conference on*  
869 *AI, Ethics, and Society*, AIES '21, pages 946–955, New York, NY, USA, 2021. Association for  
870 Computing Machinery. ISBN 978-1-4503-8473-5. doi: 10.1145/3461702.3462559.
- 871 [122] Sunny Shree, Jaganmohan Chandrasekaran, Yu Lei, Raghu N. Kacker, and D. Richard Kuhn.  
872 DeltaExplainer: A Software Debugging Approach to Generating Counterfactual Explanations.  
873 In *2022 IEEE International Conference On Artificial Intelligence Testing (AITest)*, pages  
874 103–110, 2022. doi: 10.1109/AITest55621.2022.00023.
- 875 [123] Manan Singh, Sai Srinivas Kancheti, Shivam Gupta, Ganesh Ghalme, Shweta Jain, and  
876 Narayanan C. Krishnan. Algorithmic Recourse Based on User’s Feature-Order Preference.  
877 In *Proceedings of the 6th Joint International Conference on Data Science & Management of*  
878 *Data (10th ACM IKDD CODS and 28th COMAD)*, CODS-COMAD '23, pages 293–294, New  
879 York, NY, USA, 2023. Association for Computing Machinery. ISBN 978-1-4503-9797-1. doi:  
880 10.1145/3570991.3571039.
- 881 [124] Ronal Singh, Tim Miller, Henrietta Lyons, Liz Sonenberg, Eduardo Velloso, Frank Vetere,  
882 Piers Howe, and Paul Dourish. Directive Explanations for Actionable Explainability in  
883 Machine Learning Applications. *ACM Trans. Interact. Intell. Syst.*, 13(4), December 2023.  
884 ISSN 2160-6455. doi: 10.1145/3579363.
- 885 [125] Dylan Slack, Sophie Hilgard, Himabindu Lakkaraju, and Sameer Singh. Counterfactual  
886 Explanations Can Be Manipulated. In *Advances in Neural Information Processing Systems*,  
887 volume 34, pages 62–75, 2021.
- 888 [126] Bevan I. Smith, Charles Chimedza, and Jacoba H. Bührmann. Individualized Help for At-Risk  
889 Students Using Model-Agnostic and Counterfactual Explanations. *Education and Information*  
890 *Technologies*, 27(2):1539–1558, March 2022. ISSN 1360-2357. doi: 10.1007/s10639-021-1  
891 0661-6.
- 892 [127] Jan-Tobias Sohns, Christoph Garth, and Heike Lette. Decision Boundary Visualization for  
893 Counterfactual Reasoning. *Computer Graphics Forum*, 42(1):7–20, 2023. doi: 10.1111/cgf.14  
894 650.
- 895 [128] Nina Spreitzer, Hinda Haned, and Ilse van der Linden. Evaluating the Practicality of Counter-  
896 factual Explanations. In *CEUR Workshop Proceedings*, volume 3277, pages 31–50. CEUR-WS,  
897 2022.
- 898 [129] Laura State, Salvatore Ruggieri, and Franco Turini. Reason to Explain: Interactive Contrastive  
899 Explanations (REASONX). In *Explainable Artificial Intelligence*, volume 1901 CCIS, pages  
900 421–437, Cham, 2023. Springer Nature Switzerland. ISBN 978-3-031-44064-9.
- 901 [130] Ilija Stepin, Jose M. Alonso, Alejandro Catala, and Martín Pereira-Fariña. A Survey of  
902 Contrastive and Counterfactual Explanation Generation Methods for Explainable Artificial  
903 Intelligence. *IEEE Access*, 9:11974–12001, 2021.
- 904 [131] Muhammad Suffian and Alessandro Bogliolo. Investigation and Mitigation of Bias in Ex-  
905 plainable AI. In *CEUR Workshop Proceedings*, volume 3319, pages 89–94. CEUR-WS,  
906 2022.
- 907 [132] Muhammad Suffian, Pierluigi Graziani, Jose M. Alonso, and Alessandro Bogliolo. FCE:  
908 Feedback Based Counterfactual Explanations for Explainable AI. *IEEE Access*, 10:72363–  
909 72372, 2022. ISSN 2169-3536. doi: 10.1109/ACCESS.2022.3189432.
- 910 [133] Emily Sullivan and Philippe Verreault-Julien. From Explanation to Recommendation: Ethical  
911 Standards for Algorithmic Recourse. In *Proceedings of the 2022 AAAI/ACM Conference on*  
912 *AI, Ethics, and Society*, AIES '22, pages 712–722, New York, NY, USA, 2022. Association for  
913 Computing Machinery. ISBN 978-1-4503-9247-1. doi: 10.1145/3514094.3534185.
- 914 [134] Gabriele Tolomei, Fabrizio Silvestri, Andrew Haines, and Mounia Lalmas. Interpretable  
915 Predictions of Tree-Based Ensembles via Actionable Feature Tweaking. In *Proceedings of*  
916 *the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*,  
917 pages 465–474, 2017. doi: 10.1145/3097983.3098039.

- 918 [135] Maria Tsiakmaki and Omiros Ragos. A Case Study of Interpretable Counterfactual Explanations for the Task of Predicting Student Academic Performance. In *2021 25th International Conference on Circuits, Systems, Communications and Computers (CSCC)*, pages 120–125, 919  
920 July 2021. doi: 10.1109/CSCC53858.2021.00029.  
921
- 922 [136] Stratis Tsirtsis and Manuel Gomez-Rodriguez. Decisions, Counterfactual Explanations and  
923 Strategic Behavior. In *Proceedings of the 34th International Conference on Neural Information  
924 Processing Systems*, NIPS’20, Red Hook, NY, USA, 2020. Curran Associates Inc. ISBN  
925 978-1-71382-954-6.
- 926 [137] Sohini Upadhyay, Shalmali Joshi, and Himabindu Lakkaraju. Towards Robust and Reliable  
927 Algorithmic Recourse. In *Advances in Neural Information Processing Systems*, volume 20,  
928 pages 16926–19937, 2021.
- 929 [138] Berk Ustun, Alexander Spangher, and Yang Liu. Actionable Recourse in Linear Classification.  
930 In *Proceedings of the Conference on Fairness, Accountability, and Transparency*, FAT\* ’19,  
931 pages 10–19, New York, NY, USA, 2019. Association for Computing Machinery. ISBN  
932 978-1-4503-6125-5. doi: 10.1145/3287560.3287566.
- 933 [139] Rens Van De Schoot, Jonathan De Bruin, Raoul Schram, Parisa Zahedi, Jan De Boer, Felix  
934 Weijdem, Bianca Kramer, Martijn Huijts, Maarten Hoogerwerf, Gerbrich Ferdinands,  
935 Albert Harkema, Joukje Willemsen, Yongchao Ma, Qixiang Fang, Sybren Hindriks, Lars  
936 Tummers, and Daniel L. Oberski. An open source machine learning framework for efficient  
937 and transparent systematic reviews. *Nature Machine Intelligence*, 3(2):125–133, 2021. doi:  
938 10.1038/s42256-020-00287-7.
- 939 [140] Mihaela van der Schaar and Andrew Rashbass. The case for Reality-centric AI, Feb 2023.  
940 URL <https://www.vanderschaar-lab.com/the-case-for-reality-centric-ai/>.  
941 Accessed 21.05.2024.
- 942 [141] Peter M. VanNostrand, Huayi Zhang, Dennis M. Hofmann, and Elke A. Rundensteiner. FACET:  
943 Robust Counterfactual Explanation Analytics. *Proc. ACM Manag. Data*, 1(4), December 2023.  
944 doi: 10.1145/3626729.
- 945 [142] Suresh Venkatasubramanian and Mark Alfano. The Philosophical Basis of Algorithmic Re-  
946 course. In *Proceedings of the 2020 Conference on Fairness, Accountability, and Transparency*,  
947 FAT\* ’20, pages 284–293, New York, NY, USA, 2020. Association for Computing Machinery.  
948 ISBN 978-1-4503-6936-7. doi: 10.1145/3351095.3372876.
- 949 [143] Sahil Verma, Varich Boonsanong, Minh Hoang, Keegan E. Hines, John P. Dickerson, and  
950 Chirag Shah. Counterfactual Explanations and Algorithmic Recourses for Machine Learning:  
951 A Review. *arXiv*, 2022.
- 952 [144] Sahil Verma, Keegan Hines, and John P. Dickerson. Amortized Generation of Sequential  
953 Algorithmic Recourses for Black-Box Models. In *Proceedings of the 36th AAAI Conference  
954 on Artificial Intelligence, AAAI 2022*, volume 36, pages 8512–8519. Association for the  
955 Advancement of Artificial Intelligence, 2022.
- 956 [145] Sahil Verma, Ashudeep Singh, Varich Boonsanong, John P. Dickerson, and Chirag Shah.  
957 RecRec: Algorithmic Recourse for Recommender Systems. In *Proceedings of the 32nd  
958 ACM International Conference on Information and Knowledge Management, CIKM ’23*,  
959 pages 4325–4329, New York, NY, USA, 2023. Association for Computing Machinery. ISBN  
960 9798400701245. doi: 10.1145/3583780.3615181.
- 961 [146] Kilian Vieth-Ditlmann. The algorithmic administration: automated decision-making in the  
962 public sector, May 2024. URL [https://algorithmwatch.org/en/algorithmic-admin-  
963 istration-explained/](https://algorithmwatch.org/en/algorithmic-administration-explained/). Accessed 22.05.2024.
- 964 [147] Marco Virgolin and Saverio Fracaros. On the Robustness of Sparse Counterfactual Explanations  
965 to Adverse Perturbations. *Artificial Intelligence*, 316(C), March 2023. ISSN 0004-3702.  
966 doi: 10.1016/j.artint.2022.103840.

- 967 [148] Vy Vo, Trung Le, Van Nguyen, He Zhao, Edwin V. Bonilla, Gholamreza Haffari, and Dinh  
968 Phung. Feature-Based Learning for Diverse and Privacy-Preserving Counterfactual Explanations. In *Proceedings of the 29th ACM SIGKDD Conference on Knowledge Discovery and Data Mining, KDD '23*, pages 2211–2222, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701030. doi: 10.1145/3580305.3599343.
- 972 [149] Julius Von Kugelgen, Amir-Hossein Karimi, Umang Bhatt, Isabel Valera, Adrian Weller, and  
973 Bernhard Scholkopf. On the Fairness of Causal Algorithmic Recourse. In *Proceedings of the 36th AAAI Conference on Artificial Intelligence, AAAI 2022*, volume 36, pages 9584–9594. Association for the Advancement of Artificial Intelligence, 2022.
- 976 [150] Julius von Kugelgen, Nikita Agarwal, Jakob Zeitler, Afsaneh Mastouri, and Bernhard  
977 Schölkopf. Algorithmic Recourse in Partially and Fully Confounded Settings Through Bounding  
978 Counterfactual Effects. *arXiv*, 2021.
- 979 [151] Sandra Wachter, Brent Mittelstadt, and Chris Russell. Counterfactual explanations without  
980 opening the black box: Automated decisions and the GDPR. *Harvard Journal of Law & Technology*, 31:841, 2017.
- 982 [152] Paul Y. Wang, Sainyam Galhotra, Romila Pradhan, and Babak Salimi. Demonstration of  
983 Generating Explanations for Black-Box Algorithms Using Lewis. *Proc. VLDB Endow.*, 14  
984 (12):2787–2790, July 2021. ISSN 2150-8097. doi: 10.14778/3476311.3476345.
- 985 [153] Yongjie Wang, Qinxu Ding, Ke Wang, Yue Liu, Xingyu Wu, Jinglong Wang, Yong Liu, and  
986 Chunyan Miao. The Skyline of Counterfactual Explanations for Machine Learning Decision  
987 Models. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management, CIKM '21*, pages 2030–2039, New York, NY, USA, 2021. Association for Computing Machinery. ISBN 978-1-4503-8446-9. doi: 10.1145/3459637.3482397.
- 990 [154] Yongjie Wang, Hangwei Qian, Yongjie Liu, Wei Guo, and Chunyan Miao. Flexible and Robust  
991 Counterfactual Explanations with Minimal Satisfiable Perturbations. In *Proceedings of the 32nd ACM International Conference on Information and Knowledge Management, CIKM '23*, pages 2596–2605, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 9798400701245. doi: 10.1145/3583780.3614885.
- 995 [155] Zhendong Wang, Isak Samsten, Vasiliki Kougia, and Panagiotis Papapetrou. Style-Transfer  
996 Counterfactual Explanations: An Application to Mortality Prevention of ICU Patients. *Artif. Intell. Med.*, 135(C), January 2023. ISSN 0933-3657. doi: 10.1016/j.artmed.2022.102457.
- 998 [156] Zijie J. Wang, Jennifer Wortman Vaughan, Rich Caruana, and Duen Horng Chau. GAM  
999 Coach: Towards Interactive and User-Centered Algorithmic Recourse. In *Proceedings of the 2023 CHI Conference on Human Factors in Computing Systems, CHI '23*, New York, NY, USA, 2023. Association for Computing Machinery. ISBN 978-1-4503-9421-5. doi: 10.1145/3544548.3580816.
- 1003 [157] Greta Warren, Mark T. Keane, and Ruth M. J. Byrne. Features of Explainability: How Users  
1004 Understand Counterfactual and Causal Explanations for Categorical and Continuous Features  
1005 in XAI. In *CEUR Workshop Proceedings*, volume 3251. CEUR-WS, 2022.
- 1006 [158] Greta Warren, Barry Smyth, and Mark T. Keane. “Better” Counterfactuals, Ones People  
1007 Can Understand: Psychologically-Plausible Case-Based Counterfactuals Using Categorical  
1008 Features for Explainable AI (XAI). In *Case-Based Reasoning Research and Development: 30th International Conference, ICCBR 2022, Nancy, France, September 12–15, 2022, Proceedings*, pages 63–78, Berlin, Heidelberg, 2022. Springer-Verlag. ISBN 978-3-031-14922-1. doi: 10.1007/978-3-031-14923-8\_5.
- 1012 [159] Daniel S. Weld and Gagan Bansal. The Challenge of Crafting Intelligible Intelligence. *Commun. ACM*, 62(6):70–79, may 2019. ISSN 0001-0782. doi: 10.1145/3282486. URL <https://doi.org/10.1145/3282486>.
- 1015 [160] Anjana Wijekoon, Nirmalie Wiratunga, Ikechukwu Nkisi-Orji, Kyle Martin, Chamath Pali-  
1016 hawadana, and David Corsar. Counterfactual Explanations for Student Outcome Prediction  
1017 with Moodle Footprints. In *CEUR Workshop Proceedings*, volume 2894, pages 1–8. CEUR-  
1018 WS, 2021.

- 1019 [161] Nirmalie Wiratunga, Anjana Wijekoon, Ikechukwu Nkisi-Orji, Kyle Martin, Chamath Pal-  
1020 ihawadana, and David Corsar. DisCERN: Discovering Counterfactual Explanations using  
1021 Relevance Features from Neighbourhoods. In *2021 IEEE 33rd International Conference on*  
1022 *Tools with Artificial Intelligence (ICTAI)*, pages 1466–1473, November 2021. ISBN 2375-0197.  
1023 doi: 10.1109/ICTAI52525.2021.00233.
- 1024 [162] Claes Wohlin. Guidelines for Snowballing in Systematic Literature Studies and a Replication  
1025 in Software Engineering. *EASE '14: Proceedings of the 18th International Conference on*  
1026 *Evaluation and Assessment in Software Engineering*, 2014. doi: 10.1145/2601248.2601268.  
1027 URL <https://doi.org/10.1145/2601248.2601268>.
- 1028 [163] Jingquan Yan and Hao Wang. Self-Interpretable Time Series Prediction with Counterfactual  
1029 Explanations. In *Proceedings of the 40th International Conference on Machine Learning*,  
1030 ICML'23. JMLR.org, 2023.
- 1031 [164] Jayanth Yetukuri, Ian Hardy, and Yang Liu. Towards User Guided Actionable Recourse.  
1032 In *Proceedings of the 2023 AAAI/ACM Conference on AI, Ethics, and Society*, AIES '23,  
1033 pages 742–751, New York, NY, USA, 2023. Association for Computing Machinery. ISBN  
1034 9798400702310. doi: 10.1145/3600211.3604708.
- 1035 [165] Songming Zhang, Xiaofeng Chen, Shiping Wen, and Zhongshan Li. Density-Based Reliable  
1036 and Robust Explainer for Counterfactual Explanation. *Expert Syst. Appl.*, 226(C), September  
1037 2023. ISSN 0957-4174. doi: 10.1016/j.eswa.2023.120214.

1038 **A Extended discussion of the search process**

1039 While our discussion of the search process in Section 3.1 in the main body of the document is  
1040 complete, we also provide an extended version of this discussion to allow for full reproducibility.

1041 We make use of 3 search engines to collect the initial set of studies: ACM Digital Library, IEEE  
1042 Xplore, and SCOPUS. Given the blurry distinction between AR and CEs, we consider the papers  
1043 discussing either problem. In a small scoping review, we identify several keywords common to  
1044 publications on recourse, as well as several equivalent terms to build the query shown below.

```
("Machine Learning" OR "Artificial Intelligence"  
OR "Algorithmic Decision*" OR "Consequential Decision*"  
OR Classif* OR Predict* OR "Explainable AI" OR AI OR XAI)  
AND (((Counterfactual OR Contrastive OR Actionable) AND Explanation*)  
OR ((Algorithmic OR Individual* OR Actionable) AND Recourse)  
OR Counterfactual?)
```

1045 We modify this query to account for the semantic differences between the search engines.

1046 For ACM Digital Library:

```
Title:(("Machine Learning" OR "Artificial Intelligence"  
OR "Algorithmic Decision*" OR "Consequential Decision*"  
OR classif* OR predict* OR "Explainable AI" OR ai OR xai )  
AND ( ( ( counterfactual OR contrastive OR actionable )  
AND explanation* ) OR ( ( algorithmic OR individual* OR actionable )  
AND recourse ) OR counterfactual? ) )  
OR Abstract:(("Machine Learning" OR "Artificial Intelligence"  
OR "Algorithmic Decision*" OR "Consequential Decision*"  
OR classif* OR predict* OR "Explainable AI" OR ai OR xai )  
AND ( ( ( counterfactual OR contrastive OR actionable )  
AND explanation* ) OR ( ( algorithmic OR individual* OR actionable )  
AND recourse ) OR counterfactual? ) )  
OR Keyword:(("Machine Learning" OR "Artificial Intelligence"  
OR "Algorithmic Decision*" OR "Consequential Decision*"  
OR classif* OR predict* OR "Explainable AI" OR ai OR xai )  
AND ( ( ( counterfactual OR contrastive OR actionable )  
AND explanation* ) OR ( ( algorithmic OR individual* OR actionable )  
AND recourse ) OR counterfactual? ) )
```

1047 For IEEE Xplore:

```
((("All Metadata":"Machine Learning"  
OR "All Metadata":"Artificial Intelligence"  
OR "All Metadata":"Algorithmic Decision*"  
OR "All Metadata":"Consequential Decision*"  
OR "All Metadata":classif* OR "All Metadata":predict*  
OR "All Metadata":"Explainable AI" OR "All Metadata":ai  
OR "All Metadata":xai )  
AND (((("All Metadata":counterfactual OR "All Metadata":contrastive  
OR "All Metadata":actionable ) AND "All Metadata":explanation* )  
OR ( ("All Metadata":algorithmic OR "All Metadata":individual*  
OR "All Metadata":actionable )  
AND "All Metadata":recourse )  
OR "All Metadata":counterfactual? )))
```

1048 For SCOPUS:

```
TITLE-ABS-KEY ( ( "Machine Learning" OR "Artificial Intelligence"  
OR "Algorithmic Decision*" OR "Consequential Decision*"  
OR classif* OR predict* OR "Explainable AI" OR ai OR xai )  
AND ( ( ( counterfactual OR contrastive OR actionable ) AND explanation* )  
OR ( ( algorithmic OR individual* OR actionable ) AND recourse )  
OR counterfactual? ) )
```

1049 The search is carried out on January 12<sup>th</sup> 2024 in titles, abstracts, and keywords, with 1267 results  
1050 from ACM Digital Library (The ACM Guide to Computing Literature), 513 results from IEEE Xplore,  
1051 and 2139 results from SCOPUS. This leads to a total of 3919 results, which are imported to the  
1052 Zotero reference management software for de-duplication. After removing the duplicates, we are left  
1053 with 3136 results, 44 of which are the meta-data of conference proceedings that we also remove.

1054 To facilitate the screening process, we employ the open-source ASReview tool, which makes use of  
1055 an active learning approach to re-order the set of publications, such that the most relevant ones are  
1056 always “at the top of the stack” [139]. We run ASReview on the default settings, i.e.:

```
Feature extraction technique: TF-IDF  
Classifier: Naive Bayes  
Query strategy: Maximum  
Balance strategy: Dynamic resampling (Double)
```

1057 The researchers behind the tool suggest employing a stopping rule measured in the number of  
1058 consecutive irrelevant records, which we set to 30, or 1% of the entire dataset. We accept all papers  
1059 that focus on algorithmic recourse and counterfactual explanations, completing the screening after  
1060 evaluating 1040 abstracts (33.67% of the dataset), leading to 504 (16.30%) records among which we  
1061 identify further 4 duplicates to remove. This results in the reported number of 499 relevant records.

1062 We observe that some important publications may be missing from our results. For instance, [151]  
1063 was published in the Harvard Journal of Law & Technology that is not indexed by computer science  
1064 search engines. Thus, we decide to augment the set of records by applying snowballing, which has  
1065 been shown as a good alternative to databases in systematic reviews in software engineering [162].

1066 We decide to make use of citation counts as a proxy for impact. Due to the lack of a suitable tool that  
1067 would provide unbiased citation counts for *all* papers in our dataset, we collect them from Google  
1068 Scholar. Unfortunately, citation counts on Google Scholar tend to be inflated, but as we make use of  
1069 snowballing purely to enrich the dataset, these does not impact the validity of our study. We manually  
1070 collect Google Scholar citation counts for all 499 results from the first screening on January 27<sup>th</sup>  
1071 and 28<sup>th</sup>, order them descendingly, and collect references for the top 50 (10%) “most impactful”  
1072 publications. Snowballing results in a total of 1519 new records. Indeed, we observe that [151]  
1073 (mentioned above) is referenced by 39 of the 50 publications used for snowballing.

1074 While this strategy introduces several pre-prints into our result set [52, 61, 91, 113, 143, 150], we  
1075 decide not to exclude them. Our review remains primarily concerned with peer-reviewed work. Here,  
1076 we also note that [114], which we collected as a pre-print has been published between the search and  
1077 appraisal. As such we decided to evaluate its published version and refer to it in this paper.

1078 After adding the snowballed references into our dataset, we are left with 2018 records for the second  
1079 screening with ASReview, again on the default settings. This time, we look for publications that  
1080 specifically refer to the problem of AR, “actionable” CEs, or modifying outcomes of automated  
1081 decision-making systems. We employ a stricter stopping rule to minimize the risk of false neg-  
1082 atives, completing the screening after 60 consecutive irrelevant records. We evaluate 538 results  
1083 (26.71% of the dataset), with 203 (10.06%) relevant results that are considered for full-text appraisal.  
1084 This concludes the extended discussion of the search process.



**B Evaluation of contributions**

Table 1: Evaluation of the collected publications on the types of contributions, 2017-2021.

Year	Reference	Propose methods	Theoretical frameworks	Analyses	Apply	Benchmark	Review
2017	[151]	✓	✓				
2019	[52]	✓					
	[61]	✓					
	[81]	✓					
	[85]	✓					
	[138]	✓					
2020	[35]	✓					
	[86]	✓					
	[136]	✓					
	[20]	✓					
	[26]	✓					
	[44]	✓					
	[67]	✓					
	[66]	✓					
	[99]	✓					
	[104]	✓					
	[107]	✓					
	[120]	✓					
	[112]	✓					
	[13]		✓				
	[142]		✓				
2021	[69]	✓	✓				
	[137]	✓			✓		
	[41]	✓					
	[49]	✓					
	[53]	✓					
	[73]	✓					
	[150]	✓					
	[105]	✓					
	[19]	✓					
	[22]	✓					
	[63]	✓					
	[64]	✓					
	[88]	✓					
	[98]	✓					
	[115]	✓					
	[117]	✓					
	[153]	✓					
	[161]	✓					
	[121]	✓					
	[55]		✓				
	[12]				✓		
	[113]				✓		
	[125]				✓		
	[4]					✓	
	[82]					✓	
	[89]					✓	
	[96]					✓	
	[135]					✓	
	[152]					✓	
	[160]					✓	
	[100]						✓

Table 2: Evaluation of the collected publications on the types of contributions, 2022.

Year	Reference	Propose methods	Theoretical frameworks	Analyses	Apply	Benchmark	Review	
2022	[39]	✓		✓				
	[34]	✓		✓				
	[6]	✓						
	[25]	✓						
	[50]	✓						
	[62]	✓						
	[158]	✓						
	[83]	✓						
	[56]	✓						
	[79]	✓						
	[80]	✓						
	[90]	✓						
	[93]	✓						
	[106]	✓						
	[111]	✓						
	[132]	✓						
	[131]	✓						
	[144]	✓						
	[65]	✓						
	[101]			✓	✓			
	[24]			✓		✓		
	[70]			✓				✓
	[15]			✓				
	[16]			✓				
	[94]			✓				
	[118]			✓				
	[133]			✓				
	[157]			✓				
	[128]			✓				
	[149]				✓			
[28]					✓			
[109]					✓			
[126]					✓			
[48]						✓	✓	
[143]						✓	✓	

Table 3: Evaluation of the collected publications on the types of contributions, 2023-2024.

Year	Reference	Propose methods	Theoretical frameworks	Analyses	Apply	Benchmark	Review	
2023	[36]	✓	✓					
	[29]	✓	✓					
	[116]	✓	✓					
	[9]	✓		✓				
	[42]	✓		✓				
	[75]	✓		✓				
	[147]	✓						
	[156]	✓				✓		
	[155]	✓				✓		
	[54]	✓						
	[123]	✓						
	[14]	✓						
	[72]	✓						
	[30]	✓						
	[51]	✓						
	[91]	✓						
	[92]	✓						
	[95]	✓						
	[108]	✓						
	[127]	✓						
	[129]	✓						
	[141]	✓						
	[154]	✓						
	[163]	✓						
	[164]	✓						
	[165]	✓						
	[78]	✓						
	[148]	✓						
	[74]	✓						
	[77]			✓				
	[124]			✓				
	[38]				✓			
	[57]				✓			
[102]				✓				
[3]					✓			
[76]					✓			
[8]						✓		
[60]						✓		
2024	[21]	✓						
	[114]	✓						

1086 **NeurIPS Paper Checklist**

1087 **1. Claims**

1088 Question: Do the main claims made in the abstract and introduction accurately reflect the  
1089 paper's contributions and scope?

1090 Answer: [Yes]

1091 Justification: Our main claim is that existing research on recourse is disconnected from the  
1092 practical requirements of systems where it would be applied (see Section 4 and Section 5.1).  
1093 Our claim is supported by a systematized literature review which is the contribution of this  
1094 work (Section 3). These are reflected in the abstract and the introduction.

1095 Guidelines:

- 1096 • The answer NA means that the abstract and introduction do not include the claims  
1097 made in the paper.
- 1098 • The abstract and/or introduction should clearly state the claims made, including the  
1099 contributions made in the paper and important assumptions and limitations. A No or  
1100 NA answer to this question will not be perceived well by the reviewers.
- 1101 • The claims made should match theoretical and experimental results, and reflect how  
1102 much the results can be expected to generalize to other settings.
- 1103 • It is fine to include aspirational goals as motivation as long as it is clear that these goals  
1104 are not attained by the paper.

1105 **2. Limitations**

1106 Question: Does the paper discuss the limitations of the work performed by the authors?

1107 Answer: [Yes]

1108 Justification: We highlight and discuss the three main limitations of our work in Section 5.2.

1109 Guidelines:

- 1110 • The answer NA means that the paper has no limitation while the answer No means that  
1111 the paper has limitations, but those are not discussed in the paper.
- 1112 • The authors are encouraged to create a separate "Limitations" section in their paper.
- 1113 • The paper should point out any strong assumptions and how robust the results are to  
1114 violations of these assumptions (e.g., independence assumptions, noiseless settings,  
1115 model well-specification, asymptotic approximations only holding locally). The authors  
1116 should reflect on how these assumptions might be violated in practice and what the  
1117 implications would be.
- 1118 • The authors should reflect on the scope of the claims made, e.g., if the approach was  
1119 only tested on a few datasets or with a few runs. In general, empirical results often  
1120 depend on implicit assumptions, which should be articulated.
- 1121 • The authors should reflect on the factors that influence the performance of the approach.  
1122 For example, a facial recognition algorithm may perform poorly when image resolution  
1123 is low or images are taken in low lighting. Or a speech-to-text system might not be  
1124 used reliably to provide closed captions for online lectures because it fails to handle  
1125 technical jargon.
- 1126 • The authors should discuss the computational efficiency of the proposed algorithms  
1127 and how they scale with dataset size.
- 1128 • If applicable, the authors should discuss possible limitations of their approach to  
1129 address problems of privacy and fairness.
- 1130 • While the authors might fear that complete honesty about limitations might be used by  
1131 reviewers as grounds for rejection, a worse outcome might be that reviewers discover  
1132 limitations that aren't acknowledged in the paper. The authors should use their best  
1133 judgment and recognize that individual actions in favor of transparency play an impor-  
1134 tant role in developing norms that preserve the integrity of the community. Reviewers  
1135 will be specifically instructed to not penalize honesty concerning limitations.

1136 **3. Theory Assumptions and Proofs**

1137 Question: For each theoretical result, does the paper provide the full set of assumptions and  
1138 a complete (and correct) proof?

1139  
1140  
1141  
1142  
1143  
1144  
1145  
1146  
1147  
1148  
1149  
1150  
1151  
1152  
1153  
1154  
1155  
1156  
1157  
1158  
1159  
1160  
1161  
1162  
1163  
1164  
1165  
1166  
1167  
1168  
1169  
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180  
1181  
1182  
1183  
1184  
1185  
1186  
1187  
1188  
1189  
1190  
1191  
1192  
1193

Answer: [NA]

Justification: Our work, as a literature review, does not rely on theoretical results or proofs. Nonetheless, we are explicit about the “assumptions” in that we discuss our approach to the collection and analysis of results in depth in Section 3.

Guidelines:

- The answer NA means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

**4. Experimental Result Reproducibility**

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [NA]

Justification: Our work does not rely on any experiments, so this question is not applicable. Nonetheless, we believe that we provide sufficiently in-depth characterization of the review process where other authors should be able to reproduce it (Section 3 and Appendix A).

Guidelines:

- The answer NA means that the paper does not include experiments.
- If the paper includes experiments, a No answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
  - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
  - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
  - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
  - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

1194  
1195  
1196  
1197  
1198  
1199  
1200  
1201  
1202  
1203  
1204  
1205  
1206  
1207  
1208  
1209  
1210  
1211  
1212  
1213  
1214  
1215  
1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223  
1224  
1225  
1226  
1227  
1228  
1229  
1230  
1231  
1232  
1233  
1234  
1235  
1236  
1237  
1238  
1239  
1240  
1241  
1242  
1243

## 5. Open access to data and code

Question: Does the paper provide open access to the data and code, with sufficient instructions to faithfully reproduce the main experimental results, as described in supplemental material?

Answer: [NA]

Justification: Our work does not rely on any experiments, so this question is not applicable. Nonetheless, we provide the complete list of publications covered in this review. We will also release the review data upon acceptance.

Guidelines:

- The answer NA means that paper does not include experiments requiring code.
- Please see the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- While we encourage the release of code and data, we understand that this might not be possible, so “No” is an acceptable answer. Papers cannot be rejected simply for not including code, unless this is central to the contribution (e.g., for a new open-source benchmark).
- The instructions should contain the exact command and environment needed to run to reproduce the results. See the NeurIPS code and data submission guidelines (<https://nips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- The authors should provide instructions on data access and preparation, including how to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- The authors should provide scripts to reproduce all experimental results for the new proposed method and baselines. If only a subset of experiments are reproducible, they should state which ones are omitted from the script and why.
- At submission time, to preserve anonymity, the authors should release anonymized versions (if applicable).
- Providing as much information as possible in supplemental material (appended to the paper) is recommended, but including URLs to data and code is permitted.

## 6. Experimental Setting/Details

Question: Does the paper specify all the training and test details (e.g., data splits, hyper-parameters, how they were chosen, type of optimizer, etc.) necessary to understand the results?

Answer: [NA]

Justification: Our work does not rely on any experiments, so this question is not applicable.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The experimental setting should be presented in the core of the paper to a level of detail that is necessary to appreciate the results and make sense of them.
- The full details can be provided either with the code, in appendix, or as supplemental material.

## 7. Experiment Statistical Significance

Question: Does the paper report error bars suitably and correctly defined or other appropriate information about the statistical significance of the experiments?

Answer: [NA]

Justification: Our work does not rely on any experiments, so this question is not applicable.

Guidelines:

- The answer NA means that the paper does not include experiments.
- The authors should answer "Yes" if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.

- 1244 • The factors of variability that the error bars are capturing should be clearly stated (for  
1245 example, train/test split, initialization, random drawing of some parameter, or overall  
1246 run with given experimental conditions).
- 1247 • The method for calculating the error bars should be explained (closed form formula,  
1248 call to a library function, bootstrap, etc.)
- 1249 • The assumptions made should be given (e.g., Normally distributed errors).
- 1250 • It should be clear whether the error bar is the standard deviation or the standard error  
1251 of the mean.
- 1252 • It is OK to report 1-sigma error bars, but one should state it. The authors should  
1253 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis  
1254 of Normality of errors is not verified.
- 1255 • For asymmetric distributions, the authors should be careful not to show in tables or  
1256 figures symmetric error bars that would yield results that are out of range (e.g. negative  
1257 error rates).
- 1258 • If error bars are reported in tables or plots, The authors should explain in the text how  
1259 they were calculated and reference the corresponding figures or tables in the text.

## 1260 8. Experiments Compute Resources

1261 Question: For each experiment, does the paper provide sufficient information on the com-  
1262 puter resources (type of compute workers, memory, time of execution) needed to reproduce  
1263 the experiments?

1264 Answer: [NA]

1265 Justification: Our work does not rely on any experiments, so this question is not applicable.

1266 Guidelines:

- 1267 • The answer NA means that the paper does not include experiments.
- 1268 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,  
1269 or cloud provider, including relevant memory and storage.
- 1270 • The paper should provide the amount of compute required for each of the individual  
1271 experimental runs as well as estimate the total compute.
- 1272 • The paper should disclose whether the full research project required more compute  
1273 than the experiments reported in the paper (e.g., preliminary or failed experiments that  
1274 didn't make it into the paper).

## 1275 9. Code Of Ethics

1276 Question: Does the research conducted in the paper conform, in every respect, with the  
1277 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

1278 Answer: [Yes]

1279 Justification: We have reviewed the NeurIPS Code of Ethics and we confirm that our work  
1280 conforms to it in every respect.

1281 Guidelines:

- 1282 • The answer NA means that the authors have not reviewed the NeurIPS Code of Ethics.
- 1283 • If the authors answer No, they should explain the special circumstances that require a  
1284 deviation from the Code of Ethics.
- 1285 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-  
1286 eration due to laws or regulations in their jurisdiction).

## 1287 10. Broader Impacts

1288 Question: Does the paper discuss both potential positive societal impacts and negative  
1289 societal impacts of the work performed?

1290 Answer: [Yes]

1291 Justification: Although this is not covered in a separate section, positive and negative societal  
1292 impacts of our work (and algorithmic recourse in general) are a key consideration throughout  
1293 this paper. See for instance Section 1 or Section 6.

1294 Guidelines:

- 1295 • The answer NA means that there is no societal impact of the work performed.
- 1296 • If the authors answer NA or No, they should explain why their work has no societal
- 1297 impact or why the paper does not address societal impact.
- 1298 • Examples of negative societal impacts include potential malicious or unintended uses
- 1299 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
- 1300 (e.g., deployment of technologies that could make decisions that unfairly impact specific
- 1301 groups), privacy considerations, and security considerations.
- 1302 • The conference expects that many papers will be foundational research and not tied
- 1303 to particular applications, let alone deployments. However, if there is a direct path to
- 1304 any negative applications, the authors should point it out. For example, it is legitimate
- 1305 to point out that an improvement in the quality of generative models could be used to
- 1306 generate deepfakes for disinformation. On the other hand, it is not needed to point out
- 1307 that a generic algorithm for optimizing neural networks could enable people to train
- 1308 models that generate Deepfakes faster.
- 1309 • The authors should consider possible harms that could arise when the technology is
- 1310 being used as intended and functioning correctly, harms that could arise when the
- 1311 technology is being used as intended but gives incorrect results, and harms following
- 1312 from (intentional or unintentional) misuse of the technology.
- 1313 • If there are negative societal impacts, the authors could also discuss possible mitigation
- 1314 strategies (e.g., gated release of models, providing defenses in addition to attacks,
- 1315 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
- 1316 feedback over time, improving the efficiency and accessibility of ML).

## 1317 11. Safeguards

1318 Question: Does the paper describe safeguards that have been put in place for responsible  
 1319 release of data or models that have a high risk for misuse (e.g., pretrained language models,  
 1320 image generators, or scraped datasets)?

1321 Answer: [NA]

1322 Justification: Our work poses no such risks, so this question is not applicable. We do not  
 1323 introduce any data or models.

1324 Guidelines:

- 1325 • The answer NA means that the paper poses no such risks.
- 1326 • Released models that have a high risk for misuse or dual-use should be released with
- 1327 necessary safeguards to allow for controlled use of the model, for example by requiring
- 1328 that users adhere to usage guidelines or restrictions to access the model or implementing
- 1329 safety filters.
- 1330 • Datasets that have been scraped from the Internet could pose safety risks. The authors
- 1331 should describe how they avoided releasing unsafe images.
- 1332 • We recognize that providing effective safeguards is challenging, and many papers do
- 1333 not require this, but we encourage authors to take this into account and make a best
- 1334 faith effort.

## 1335 12. Licenses for existing assets

1336 Question: Are the creators or original owners of assets (e.g., code, data, models), used in  
 1337 the paper, properly credited and are the license and terms of use explicitly mentioned and  
 1338 properly respected?

1339 Answer: [NA]

1340 Justification: Our work does not use existing assets (other than the referenced papers), so  
 1341 this question is not applicable. All papers covered in the review are referenced in sufficient  
 1342 detail, so that the readers can access them.

1343 Guidelines:

- 1344 • The answer NA means that the paper does not use existing assets.
- 1345 • The authors should cite the original paper that produced the code package or dataset.
- 1346 • The authors should state which version of the asset is used and, if possible, include a
- 1347 URL.



- 1348 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 1349 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 1350 service of that source should be provided.
- 1351 • If assets are released, the license, copyright information, and terms of use in the package
- 1352 should be provided. For popular datasets, [paperswithcode.com/datasets](https://paperswithcode.com/datasets) has
- 1353 curated licenses for some datasets. Their licensing guide can help determine the license
- 1354 of a dataset.
- 1355 • For existing datasets that are re-packaged, both the original license and the license of
- 1356 the derived asset (if it has changed) should be provided.
- 1357 • If this information is not available online, the authors are encouraged to reach out to
- 1358 the asset's creators.

### 1359 13. **New Assets**

1360 Question: Are new assets introduced in the paper well documented and is the documentation

1361 provided alongside the assets?

1362 Answer: [NA]

1363 Justification: Our work does not release any new assets, so this question is not applicable.

1364 We release the paper with the most permissible license available for NeurIPS submissions.

1365 Finally, we will release the review data upon acceptance.

1366 Guidelines:

- 1367 • The answer NA means that the paper does not release new assets.
- 1368 • Researchers should communicate the details of the dataset/code/model as part of their
- 1369 submissions via structured templates. This includes details about training, license,
- 1370 limitations, etc.
- 1371 • The paper should discuss whether and how consent was obtained from people whose
- 1372 asset is used.
- 1373 • At submission time, remember to anonymize your assets (if applicable). You can either
- 1374 create an anonymized URL or include an anonymized zip file.

### 1375 14. **Crowdsourcing and Research with Human Subjects**

1376 Question: For crowdsourcing experiments and research with human subjects, does the paper

1377 include the full text of instructions given to participants and screenshots, if applicable, as

1378 well as details about compensation (if any)?

1379 Answer: [NA]

1380 Justification: Our paper does not involve crowdsourcing or research with human subjects, so

1381 this question is not applicable. The work was in its entirety carried out by the authors.

1382 Guidelines:

- 1383 • The answer NA means that the paper does not involve crowdsourcing nor research with
- 1384 human subjects.
- 1385 • Including this information in the supplemental material is fine, but if the main contribu-
- 1386 tion of the paper involves human subjects, then as much detail as possible should be
- 1387 included in the main paper.
- 1388 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 1389 or other labor should be paid at least the minimum wage in the country of the data
- 1390 collector.

### 1391 15. **Institutional Review Board (IRB) Approvals or Equivalent for Research with Human**

1392 **Subjects**

1393 Question: Does the paper describe potential risks incurred by study participants, whether

1394 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

1395 approvals (or an equivalent approval/review based on the requirements of your country or

1396 institution) were obtained?

1397 Answer: [NA]

1398 Justification: Our work does not involve crowdsourcing or research with human subjects, so

1399 this question is not applicable. We did not require an IRB approval or equivalent to carry

1400 out this work.

1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411

**Guidelines:**

- The answer NA means that the paper does not involve crowdsourcing nor research with human subjects.
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
- We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
- For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.