

1 Title: The views of health guideline developers on the use of automation  
2 in health evidence synthesis

3 Authors:

4 Anneliese Arno<sup>1\*</sup>: anneliese.arno.17@ucl.ac.uk

5 Julian Elliott<sup>2</sup>: julian.elliott@monash.edu

6 Byron Wallace<sup>3</sup>: b.wallace@northeastern.edu

7 Tari Turner<sup>2</sup>: tari.turner@monash.edu

8 James Thomas<sup>1</sup>: james.thomas@ucl.ac.uk

9 <sup>1</sup> Eppi-Centre, University College London, United Kingdom

10 <sup>2</sup> School of Public Health and Preventive Medicine, Monash University, Australia

11 <sup>3</sup> College of Computer and Information Science, Northeastern University, United States

12 \*corresponding author

13

14

15

## 16 Abstract

## 17 Background

18 The increasingly rapid rate of evidence publication has made it difficult for evidence synthesis –  
19 systematic reviews and health guidelines -- to be continually kept up to date. One proposed solution for  
20 this is the use of automation in health evidence synthesis. Guideline developers are key gatekeepers in  
21 the acceptance and use of evidence, and therefore their opinions on the potential use of automation are  
22 crucial.

## 23 Methods

24 The objective of this study was to analyze the attitudes of guideline developers towards the use of  
25 machine learning and crowd-sourcing in evidence. The Diffusion of Innovations framework was chosen  
26 as an initial analytical framework because it encapsulates some of the core issues which are thought to  
27 affect the adoption of new innovations in practice. This well-established theory posits five dimensions  
28 which affect the adoption of novel technologies: *Relative Advantage*, *Compatibility*, *Complexity*,  
29 *Trialability*, and *Observability*.

30 Eighteen interviews were conducted with individuals who were currently working, or had previously  
31 worked, in guideline development. After transcription, a multiphase mixed deductive and grounded  
32 approach was used to analyze the data. First, transcripts were coded with a deductive approach using  
33 Rogers' Diffusion of Innovation as the top-level themes. Second, sub-themes within the framework were  
34 identified using a grounded approach.

## 35 Results

36 Participants were consistently most concerned with the extent to which an innovation is in line with  
37 current values and practices (i.e. *Compatibility* in the Diffusion of Innovations framework). Participants

38 were also concerned with *Relative Advantage* and *Observability*, which were discussed in approximately  
39 equal amounts. For the latter, participants expressed a desire for transparency in methodology of  
40 automation software. Participants were noticeably less interested in *Complexity* and *Triability*, which  
41 were discussed infrequently. These results were reasonably consistent across all participants.

## 42 Conclusions

43 If machine learning and other automation technologies are to be used more widely and to their full  
44 potential in systematic reviews and guideline development, it is crucial to ensure new technologies are  
45 in line with current values and practice. It will also be important to maximize the transparency of the  
46 methods of these technologies to address the concerns of guideline developers.

## 47 Keywords

48 Automation; systematic reviews; machine learning; guideline development; diffusion of innovation;  
49 evidence synthesis

50

51

52

53

54

55

56

57

## 58 Background

### 59 Evidence-based guidelines are overwhelmed by the rate of research publication

60 As guidelines increasingly incorporate an evidence-based medicine approach, the systematic reviews  
61 which are a crucial component of this evidence have become overwhelmed by the rate of publication of  
62 new evidence [1]. With nearly 4000 health research articles published daily, systematic reviews cannot  
63 keep up with the deluge of data [2]. Research is at risk of being wasted, leading to out-of-date  
64 healthcare and guidelines, and consequently impacting on population health outcomes.

### 65 Limited literature addresses adoption of automation

66 Given this, there is increasing interest in the use of automation in the completion of systematic reviews  
67 [3-6]. While some herald the use of automation, others are hesitant to adopt these novel methods.  
68 Automation technologies include machine learning (ML), natural language processing (NLP), text mining,  
69 among other technologies.

70 The literature on the topic of automation in health evidence synthesis is somewhat lacking. Previous  
71 publications largely focus on two main areas: potential applications of automation, and the validity or  
72 accuracy of automation tools [5]. Uptake of automation has been notably slow [6, 7], despite the broad  
73 availability of various tools, particularly for study screening, and the availability of peer-reviewed  
74 literature addressing accuracy and potential integration opportunities. This leads to the question: what  
75 factors are inhibiting uptake of automation into systematic reviews and into guidelines?

76 There has not yet been significant primary research into barriers and facilitators to uptake of  
77 automation in health evidence synthesis contexts [7, 8]. Considering the slow adoption rate,  
78 understanding the barriers and facilitators to uptake is important, as is exploring the perceptions of key  
79 stakeholders in evidence production towards the uptake of automation. As outlined above, guidelines

80 are a key component in the translation of knowledge to practice, therefore evidence from guideline  
81 developers detailing their views of the use of automation would be helpful in elucidating the reasons for  
82 the slow adoption of automation.

### 83 Diffusion of Innovations

84 Rogers' Diffusion of Innovations is a highly applicable framework for analysis of these views and the  
85 adoption of automation in health evidence production more broadly [9]. This theory describes how,  
86 why, and at what rate an innovation spreads, the characteristics of the innovation that play a role in this  
87 process, and the typical categories of innovators. Its insights have been repeatedly supported by  
88 empirical data in a broad range of contexts.

89 The characteristics of an innovation that impact its diffusion are described as follows:

#### 90 *Complexity*

91 *Complexity* is how easy an innovation is, or is perceived to be, to comprehend and to put into use.

#### 92 *Compatibility*

93 *Compatibility* refers to an innovation being in line with existing values and practices, and the needs of  
94 potential future users.

#### 95 *Trialability*

96 *Trialability* refers to the ability of users or potential users to experiment with an innovation prior to  
97 adopting it.

#### 98 *Observability*

99 *Observability* is the degree to which potential users may examine the results of an innovation.

## 100 *Relative Advantage*

101 *Relative advantage* refers to how much better an innovation is, or is perceived to be, than the system it  
102 is replacing.

103 These five elements collectively influence potential adopters' decisions toward adoption of an  
104 innovation. In distinct contexts, and with distinct populations, some characteristics may play a greater  
105 role than some others. Understanding the comparative role of these characteristics in the context of  
106 systematic reviews' and health guidelines' potential use of automation should prove useful in describing  
107 the current state of adoption, as well as considering future research foci and organizational norm  
108 setting.

## 109 *Research questions*

110 The goal of this research was to gather data from guideline developers regarding their attitudes and  
111 perceptions of automation, and specifically automation applied to health evidence production. Research  
112 questions were:

- 113 1) How do the opinions of guideline developers of automation of health evidence synthesis fit into  
114 the Diffusion of Innovations framework?
- 115 2) Within the Diffusion of Innovations themes, what important concepts were identified by  
116 participants?

## 117 *Methods*

### 118 *Participants*

119 Participants were recruited via existing personal and professional networks for guideline developers and  
120 systematic review researchers. These networks included Guidelines International Network (GIN) and  
121 National Institute for Health Care and Excellence (NICE). Participants were required to be developers of

122 health policy or clinical practice guidelines, and/or to be familiar with current practices of guideline  
123 development.

124 Potential participants were invited to participate in a semi-structured one-on-one interview conducted  
125 via phone or via Skype with Ms Arno, with the session to be audio recorded with participant consent.

126 None personally knew Ms Arno prior to completion of the interviews. Participants were provided  
127 information on her background as a PhD student at University College London studying the adoption of  
128 automation in health evidence synthesis. An interview instrument was developed by the research team  
129 prior to the interviews and was applied in all interview sessions, with variation in follow-up questions as  
130 relevant according to participants' responses.

### 131 [Data Collection and Analysis](#)

132 Following verbatim transcription of the interviews, they were provided to participants for validation.  
133 They were then analyzed using QSR NVivo 12 [10].

134 A thematic approach, as outlined by Braun and Clarke [11], was applied. This method was adapted to  
135 incorporate both deductive and inductive analyses. This allowed for framework analysis in addition to  
136 reflexive and iterative insights from the resulting data [12].

137 Analysis took place in five stages, described in more detail below.

#### 138 [Stage 1: Assignment within predefined frameworks](#)

139 First, Rogers' Diffusion of Innovation framework was used as the top-level deductive codes, using a line-  
140 by-line verbatim assignment of transcripts to one or more of the five themes (i.e. *Complexity*,  
141 *Compatibility*, *Trialability*, *Observability*, and *Relative Advantage*).

142 Stage 2: Open coding within Diffusion of Innovations framework

143 Once each transcript was coded according to the top-level frameworks (i.e. Diffusion of Innovations), a  
144 codebook -- a document containing all data belonging to a code or theme -- was generated for each of  
145 the five themes. These codebooks were then examined with an open coding method.

146 Stage 3: Generation of themes

147 The codebook of each Diffusion of Innovation theme was reviewed across all transcripts together to  
148 identify shared patterns among the grounded open codes. Each individual verbatim code was grouped  
149 with others with similar meaning and content, forming preliminary explanatory themes.

150 Following formation of these themes, a further review process was undertaken to reconsider how the  
151 themes fitted together. In addition, outlying codes were identified as those that either had not been  
152 grouped with codes from other transcripts, or those that had relatively few grounded codes grouped  
153 together.

154 Stage 4: Generation of matrices

155 A matrix was generated comparing each of the top-level framework themes against the data-driven  
156 themes. This approach not only allowed for an overview of the relative significance of each overall  
157 theme -- thus addressing research question 1 -- but also to examine this significance through different  
158 lenses. For instance, isolating the framework to selected grounded codes might give a different  
159 impression of the results of the framework analysis.

160 Stage 5: Identifying patterns and outliers

161 These matrices were finally used to analyze the data in relation to the first research question (i.e. how  
162 do guideline developers' opinions on automation relate to the Diffusion of Innovations framework?),  
163 and to expand upon these data in relation to the second research question (i.e. what important  
164 concepts were identified by participants?).

## 165 Results

### 166 Participants

167 Eighteen interviews were conducted and varied in length from approximately 30 minutes to  
168 approximately 80. Five participants were male, and 13 were female. Half of the participants had  
169 between five and ten years of experience in evidence synthesis; five participants had between ten and  
170 20 years of experience; two participants had more than 20 years of experience, and two had less than 5  
171 years of experience. No participants withdrew from the study, and no repeat interviews were required.

### 172 Overview

173 Interview transcripts demonstrated high consistency in distribution within the Diffusion of Innovations  
174 thematic framework. Following initial coding (Stage 1), *Compatibility* was the most prominently  
175 discussed theme across all participants. *Relative Advantage* and *Observability* were also given  
176 substantial attention from participants' discussions, though to a lesser extent than *Compatibility*.  
177 *Trialability* and *Complexity* were the least discussed themes among all participants.

### 178 Compatibility

179 All participants discussed their values as guideline developers at length, both within the context of  
180 potential use of automation and independent from it. Emphasis on the *Compatibility* theme was  
181 consistently far stronger than any of the other four themes in the deductive framework being applied in  
182 this analysis. Some examples of values were a "rigorous" approach to evaluation of evidence and careful  
183 construction of questions. Relating to automation more specifically, participants highlighted a need for  
184 human and organizational involvement.

185 *"how you synthesize it, how you pull it together is kind of key" Participant 3*

186 *"I think it would be a shame if humans weren't involved in [synthesis]." Participant 9*

187 Two sub-themes were identified within the *Compatibility* theme which further detailed participants'  
188 desire to match new practices with the values which underpin current practices: *ability to double check*,  
189 and *transparency as accountability*.

### 190 *Ability to double check*

191 Most participants indicated the importance of the *ability to double check* the output of automation with  
192 human researcher input. These discussions often cited the rationale that current practices usually  
193 involve a human double checking the work of another human and posited that newer workflows should  
194 therefore maintain this pattern.

195 Some, but not all, participants indicated that reproducibility was the underlying reason for the double-  
196 checking status quo. It is possible that views of participants would be different should rigorous research  
197 alter overall perceptions of the reproducibility of automated screening and extraction; this is further  
198 discussed as a contextual factor in subsequent sections.

199 *"I can see it could be done. But surely it would need to be checked by someone*  
200 *anyway. Because even if it's done by a human with vast experience, it's always*  
201 *important to have a second person to check it." Participant 5*

202 *"At the minute the standard is for two operators. So you'd want it to have been*  
203 *checked by a second method, if not person. So that would be my only thing –*  
204 *reproducibility." Participant 7*

## 205 Transparency as accountability

206 Several participants wanted to ensure that any automation methods used in synthesizing evidence were  
207 freely accessible and transparent to examination. Many emphasized that they are accountable to  
208 stakeholders who need to be sure they have not missed any information, and therefore require the  
209 ability to freely examine methods used, including any automation.

210 Trustworthiness of evidence in general is integral to the professional culture of guideline development  
211 and was emphasized by the participants. Trustworthiness and methods to verify it therefore extend to  
212 new tools that use automation, in the view of participants, in the form of transparency and validation.

213 *“A group of experts can apply judgement to that body of evidence, and needs to*  
214 *know they can trust the evidence that you’d found.” Participant 12*

215 *“The key part of working with a face to face committee ... Is you have they have to*  
216 *have total confidence in what the technical team has done” Participant 16*

## 217 Relative Advantage

218 When discussing *Relative Advantage*, participants focused on the *freeing up of human resources*, and to  
219 a lesser extent on *time and cost saving*. When prompted to discuss ML directly (in contrast to general  
220 views of evidence synthesis and guideline development approaches), participants tended to more  
221 frequently discuss ideas relating to the *Relative Advantage* of automation. Participants were interested  
222 in freeing time and money, but contingent upon the automation perfectly matching perceived human  
223 quality.

## 224 Freeing up human resources

225 The primary advantage specified in discussion with participants was the potential to free up human  
226 resources for rededication to additional tasks within the health evidence ecosystem.

227 *“In research time is always limited and you know there’s never enough grant money*  
228 *to help employ staff ... by having a machine do it, it would be cost-effective, and*  
229 *spare the researchers’ time to do other research-related tasks.” Participant 17*

### 230 Time and cost saving

231 Some participants also identified that automation might potentially save time and/or save money.  
232 Strikingly, no participants indicated an openness to any trade-off between accuracy and time.

233 *“No matter how quickly a guideline’s done, everybody always wants it faster and to*  
234 *be of high quality. So anything that can improve on that would be welcome, I think.”*  
235 *Participant 11*

### 236 Observability

237 Participants communicated that they would like to see evidence prior to implementing new practices, as  
238 well as a sustained ability to cross-examine the behavior of the technologies.

### 239 Need for evidence

240 The need for rigorously produced, disseminated, and easily accessed evidence was clear in the data.  
241 Several participants expressed an openness to automation being integrated into evidence synthesis, on  
242 the condition that accuracy has been demonstrated.

243 *“I think at the moment it has a potentially high level of risk of being incorrect. But I*  
 244 *don’t really know enough about it. I’d need to be convinced about it I think to*  
 245 *consider it.” Participant 9*

246 *“If the whole process were done by some machine or machine learning application, I*  
 247 *think it would need to be properly trialed.” Participant 5*

248 *“As long as there was clear data to support that ... machine-learning is a reliable*  
 249 *method, but you know, better than or equal to humans doing it.” Participant 17*

250 One notable outlier indicated they were already convinced of automation’s abilities within the specific  
 251 context of screening. This unusual case raises the possibility that these results would be different given  
 252 further evidence production and dissemination.

253 *“I do think it’s been well demonstrated for the screening aspects, for the hit rates of*  
 254 *what gets included and what doesn’t, and how correct it is.” Participant 11*

### 255 Personal need for double-checking

256 Participants often wanted an established and ongoing method of observing the inner workings of the ML  
 257 processes, frequently described as a desire to “check” what ML had done. While similar to the previous  
 258 theme of *Compatibility: ability to double check*, the latter discussed that guideline developers believe  
 259 the ability to check methods should be available as a matter of principle, while *Observability: personal*  
 260 *need for double-checking* discusses that guideline developers want to do such checking themselves.

261 This need to be able to continually check how the machine learning has processed information could be  
 262 interpreted as a desire to maintain control over the evidence synthesis process. As previously discussed,  
 263 guideline developers must convince other stakeholders of their recommendations’ integrity, so personal  
 264 quality control fits in with the cultural expectations of guideline development.

265 *“The thing that’s sort of a little bit distressing from a novice point of view with*  
 266 *machine-learning is not feeling like I have a way to check it... I’d need some way to be*  
 267 *confident .... [I’d need] a way to check the algorithms” Participant 3*

## 268 Complexity and Trialability

269 Selected participants identified that the learning process would need to be simple if researchers were to  
 270 adopt ML. They also expressed a preference for familiarity over the novel.

271 *“Whenever you try and really change things, I think there’s a degree of skepticism*  
 272 *anyway...I think that might just be the nature of human beings.” Participant 9*

273 *“If they have to learn the process, and if it’s hard, then that sort of discourages*  
 274 *them.” Participant 18*

275 *“So unless the technology offers a value add that’s substantial enough to overcome*  
 276 *the learning curve...however much time it takes to do that has to not be more time*  
 277 *than you’re gonna save.” Participant 3*

## 278 Contextual themes

279 Upon re-examination of how the data informed the deductive framework (described in Step 5 of the  
 280 Methods section), several contextual factors were identified.

### 281 Participant familiarity with automation

282 Participants nearly always offered disclaimers prior to commenting, indicating they felt they did not  
 283 have sufficient experience with automation technologies to be able to comment at their desired level of

284 expertise. These data were of significant interest as they demonstrated a current lack of robust  
 285 knowledge of the capabilities of automation within the target population.

286 *“I’ve done a very little bit with machine-learning.” Participant 3*

287 *“It’s just my concern would be that I’ve not had any experience with it.” Participant 7*

288 *“I haven’t had much to do with machine-learning. Like I’ve kind of heard about it”*

289 *Participant 17*

290 *“I think that’s something I have no personal experience with” Participant 11*

291 *“To be honest I actually haven’t had much experience with it” Participant 8*

292 *“Yeah, I don’t know, I don’t really understand that [OBJ] process.” Participant 5*

### 293 Overall skepticism towards Machine Learning

294 Overall skepticism or mistrust towards automation, both towards current technologies and anticipated  
 295 future ones, was clear in the contributions from participants. They particularly expressed doubt over the  
 296 ability of a machine to mimic human judgement calls they felt are currently essential to well-formulated  
 297 health guidelines.

298 *“It would be very difficult to train a machine to make the sort of value decisions that*

299 *we have to make” Participant 10*

300 *“I’m still a bit nervous about some of the interpretation of that...it just might be a*

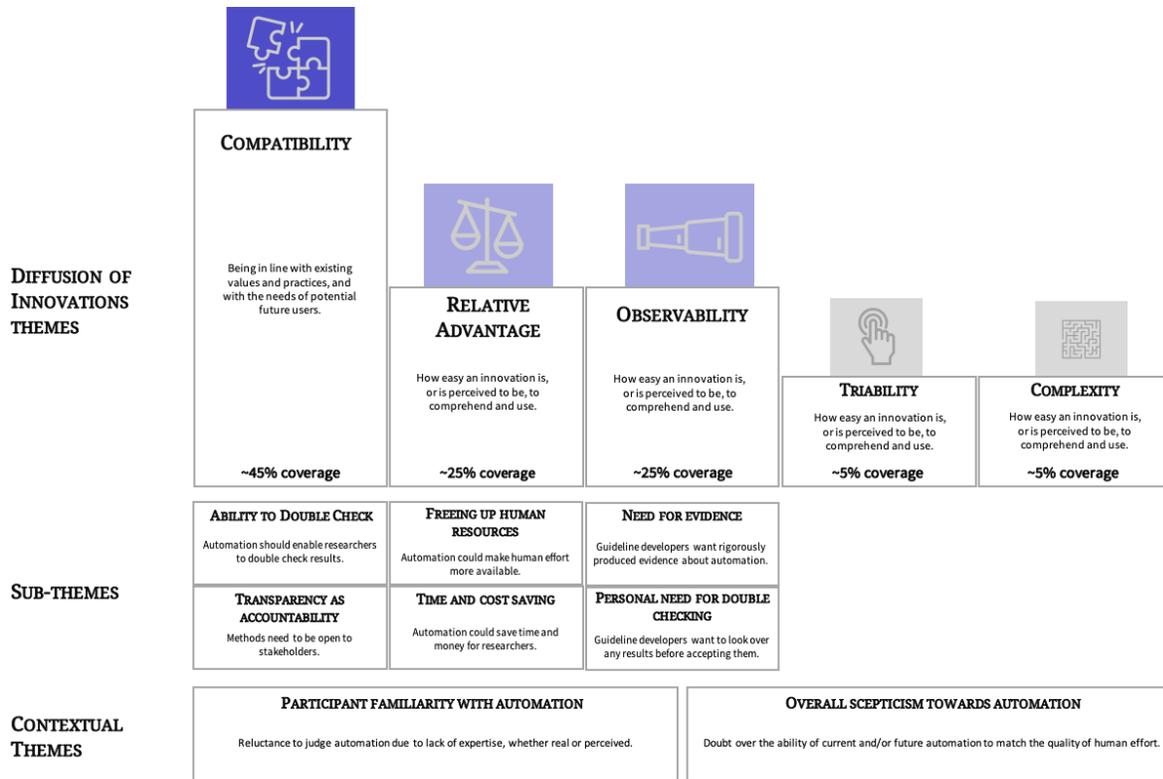
301 *distrust about it, I think?” Participant 13*

302                    *“How can a computer apply judgement? ...There’s judgement required when it*  
303                    *comes to things like quality or – they are not things I expect to be evidence that could*  
304                    *be accurate.” Participant 12*

305                    *“I don’t think it could fully replace a human ... I think there can be subtleties between*  
306                    *how things can interact... I think there’s always going to be some sort of human*  
307                    *element.” Participant 9*

308                    *“I don’t know if we’re there yet. Maybe we’ll get to the point where we can do that,*  
309                    *but to do that, like quality rating, or to do a level of evidence, or strength of*  
310                    *evidence... I mean there’s still a lot of value judgements in that. And I don’t know how*  
311                    *much machine learning could help with that at this point.” Participant 3*

312    Figure 1 presents a summary of the themes, sub-themes, and contextual themes.



313

314 Figure 1. Summary of results

315 Discussion

316 Cultural standards of practice greatly influence decision-making

317 Guideline developers demonstrated deeply held core beliefs about evidence synthesis methodology and  
 318 perceived quality. These will be central in the potential adoption of automation to health evidence  
 319 synthesis.

320 A 2013 paper commenting on the reasons for slow uptake of automation posed the broad question:

321 “why is [automation] not yet widely used?” [13]. At the time, the authors concluded that “further  
 322 technical and empirical work is needed ... [to] develop solutions which have a demonstrative relative  
 323 advantage, and which are clearly compatible with the needs of systematic reviewers and their users.”

324 That is, they considered *Relative Advantage* and *Compatibility* the most important themes, playing the

325 most significant roles in the adoption of automation. Considering the data presented in this study in  
326 relation to Thomas et al's [13] question, the prior conclusions should be adjusted slightly.

327 The most significant reason appears to be that key stakeholders of EBM are not persuaded that  
328 automation is compatible with their guiding values, principles, or standards. While *Relative Advantage*  
329 was important, it was secondary to the far more prominent discussion of *Compatibility*. Further, the  
330 identified sub-themes of *Compatibility* focused more on automation's fitting in with the values behind  
331 current practices, rather than on fitting in with existing infrastructure

332 The preceding points not only represent a shift from the hypotheses presented in the literature, but also  
333 from the focus of previous discussions at relevant conferences. The International Collaboration for the  
334 Automation of Systematic Reviews (ICASR), formed in 2015, is a global network endeavoring to  
335 successfully integrate all the parts of automation of systematic review production together. In the notes  
336 of the third ICASR meeting in 2017, the group concluded that the "most pressing needs at present are to  
337 develop approaches for validating" automation and integration with existing system architecture [14].  
338 Stated another way, ICASR believed *Observability* to be critical to uptake, as well as *Compatibility*  
339 specifically in reference to fitting into existing practice.

340 As in the previous case discussed, this research has provided some evidence to support this assertion  
341 but suggests a slight redirection in which priorities would be best suited to promotion of automation  
342 adoption. While evidence gathered in this study reinforces that *Compatibility* plays a significant role, it  
343 also demonstrates that alignment with values is more highly prioritized than alignment with current  
344 practice and system architecture. In addition, the "most pressing need" may not be validation  
345 (*Observability*), but instead the demonstration and communication of methodological standards and  
346 cultural coordination (*Compatibility*).

347 Researcher effort will be redirected rather than replaced

348 Guideline developers anticipate that automation will be most useful in redirecting person-time rather  
349 than replacing it. They highlighted that a critical (and in their view, irreplaceably human) part of their  
350 professional contribution is the nuanced judgements applied to the presented evidence, often derived  
351 from lived experience. Automation could contribute to an improvement in guideline quality by providing  
352 additional resources (namely, time) to more difficult aspects of guideline development, and not simply  
353 by cutting costs and workload.

354 Contributions from participants in this study relating to reluctance to relinquish human judgement align  
355 with notes from the previously mentioned ICASR meeting [14]. They stated:

356 *“For example, external stakeholders might believe the current vision is automated*  
357 *reviews devoid of valuable human control and input, that is, a general autonomous*  
358 *artificial intelligence system. That view, however, was neither represented nor*  
359 *sanctioned at the meeting. Therefore, improving the terminology associated with*  
360 *systematic review automation to reflect the goal more accurately is likely valuable.”*

361 This study provides evidence in support of this proposition: participants were wary of automation in part  
362 due to the idea that it might remove crucial human judgement in the process of guideline development;  
363 notably, encouraging complete and total replacement is “neither represented nor sanctioned.” Given  
364 that participants in this study echoed this fear, it raises the question of why guideline developers hold  
365 this view, and how to best communicate a more accurate representation of the goals of promoters of  
366 automation in systematic reviews and guidelines.

367 This observed anticipation of automation allowing for refocusing of effort is what should be expected if  
368 the results of this study are situated in the historical evidence and context. From the late nineteenth  
369 century onwards, there have been repeated waves of automation of production and consequent

370 population-level job panic [15]. With each wave, however, human effort has not been removed, but  
371 rather redirected. In some cases, job opportunities have expanded rather than contracted as a by-  
372 product of widespread automation. Therefore, in addition to enabling the valuable skills of EBM  
373 researchers to be better spent, it is possible the field will see an expansion of opportunities.

374 As in the previous section, perhaps proponents of automation have a choice to steer the general  
375 conversation to clarify that expert opinion will not be superseded but will instead be made less costly  
376 and more available by freeing up person-time and other resources. An enabling environment for the  
377 promotion and adoption of automation in a manner that redirects rather than replaces research effort  
378 could be an effective strategy in building consensus among guideline developers, as key stakeholders of  
379 the evidence synthesis process, in accepting, implementing, and promoting automation practices.

## 380 [Study limitations](#)

### 381 [Participant sample](#)

382 One limitation of this study is the use of purposive sampling; an inherent weakness of this technique is a  
383 tendency towards similar respondents [16]. Purposive sampling was used to ensure that participants  
384 could be recruited within the time available. While it was necessary to take a purposive sample in order  
385 to gather this data, it must be acknowledged that this method introduces potential for bias to the  
386 results.

387 It might be expected that purposive sampling used in this manner (i.e. using personal direct contacts and  
388 networks) would result in respondents with similar views to the investigators. This was not observed,  
389 however. The generally low awareness of the capabilities of ML, and moreover the aims of integration of  
390 ML into EBM, indicates that participants most likely did not hold similar opinions to the researchers.

391 One clear skew in the resulting sample was the inclusion of only 28% male participants. Despite this, the  
392 views of participants did not appear to vary according to gender. It is possible that this balance is

393 representative of the current balance in the field of guideline development; for example, according to  
394 data available from NICE, they employ 68.63% women and 31.37% men [17].

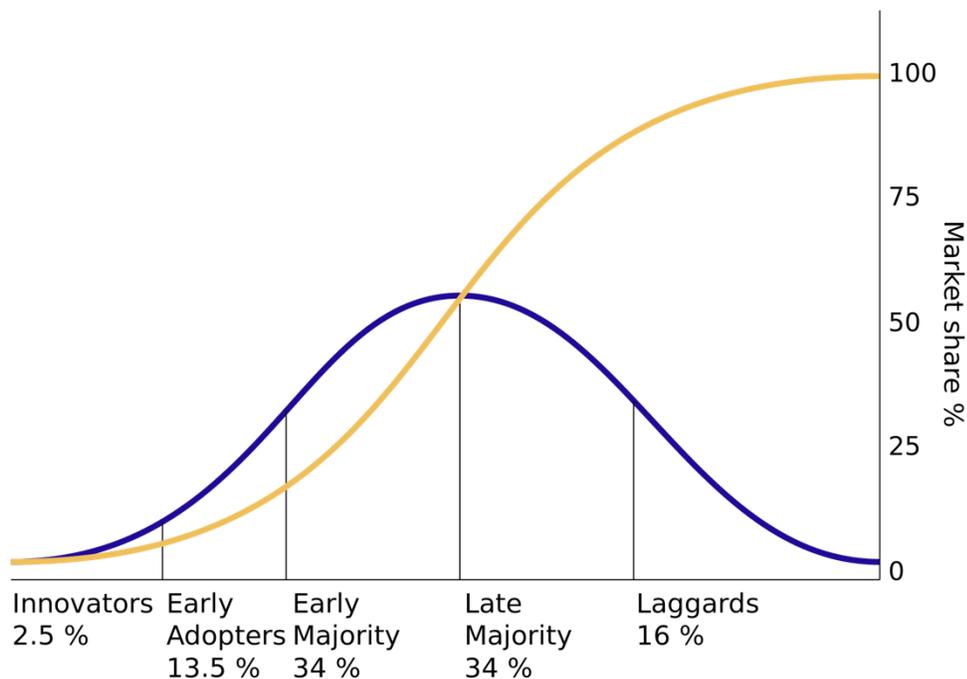
395 Current state on the adoption curve

396 Finally, as previously outlined in the Background section, *Diffusion of Innovations* describes the typical  
397 categories of innovators (personas) and provides an approximation of the expected proportions of each  
398 category. These are shown in Figure 2. As time progresses, successive groups will adopt a given  
399 innovation, until a critical mass of the market share is reached. In addition to the five categories shown,  
400 non-adopters are sometimes added as a sixth category.

401 The finding that many of the participants perceived themselves to be inexperienced with automation in  
402 the context of evidence synthesis, raises the question of where the field currently resides within this  
403 adoption curve. The evidence of this study suggests that the field is in very early stages of broader  
404 adoption (i.e. innovators) with only a small minority taking on use of this new technology.

405 Once a later stage has been reached, subsequent studies may well return different results. For example,  
406 the case from Participant 11 in *Observability: Need for evidence*, while certainly an outlier in the context  
407 of this study, may fall under an innovator or early adopter persona while other participants may fall  
408 under late majority or laggards. Additional analysis and/or data collection using persona categories as  
409 the deductive framework could build upon the results of this study. However, until a later stage of

410 diffusion is reached, it may be difficult to find sufficient contributors within each category.



411

412 *Figure 2. Diffusion of Innovations adoption curve and market share*

413 [Suggestions for future research](#)

414 While guideline developers are a crucial group within the field of evidence-based medicine, they are far  
 415 from the only one. Therefore, it would be logical to repeat this study with different population groups.

416 One important group to consult is systematic reviewers -- the individuals who will be using automation  
 417 software directly, as opposed to those who are gatekeepers of the evidence produced by such software  
 418 (i.e. guideline developers).

419 Patient stakeholders are also an important group to consult. Patients are often involved on guideline  
 420 panels, and there have been recent pushes to include more consumers and patients in development of  
 421 health guidance [18]. Health guidelines should ultimately aim to benefit patients and the community,  
 422 and organizational mission statements often (and rightly) include statements about patient  
 423 transparency and empowerment. Finally, policy makers should also be examined, as they were

424 identified by some of the participants in this study was fellow stakeholders in the process of creating  
425 guidelines.

426 Future research should be prioritized and proceed in parallel to the forms of validation highlighted by  
427 participants as crucial to their decision making. Select examples of automation have long been available  
428 for evidence synthesis, and several prominent organizations are encouraging automation uptake.  
429 Despite this reality, they are not being integrated into workflows at large or even medium scale. The  
430 data from this study show hesitation from a key stakeholder group, and additional data relating to other  
431 user stakeholders will be helpful in identifying barriers and facilitators for these groups.

## 432 Conclusions

433 Analyzed via the lens of the Diffusion of Innovations framework, the results of this study strongly  
434 demonstrate that *Compatibility* with professional cultural values is the most significant consideration for  
435 guideline developers in the potential adoption of automation. Participating guideline developers  
436 identified increased availability of person-time as a primary *Relative Advantage*, and desired rigorous  
437 validation (*Observability*) to occur both prior to adoption and on an ongoing basis. Participants' lack  
438 knowledge of ML is a contributing contextual factor to slow uptake of automation, along with a  
439 generalized anxiety towards relinquishing human control to a computer. This contextual factor means  
440 that future studies may return different results if and when the evidence synthesis field reaches a later  
441 stage in the adoption curve. The data demonstrated an inaccurate perception that nuanced human  
442 judgement is likely to be removed from evidence synthesis by automation.

443 The creation and dissemination of empirical evidence that systematically demonstrates automation's  
444 alignment with the values and standards of guideline development and EBM should therefore be  
445 prioritized. In addition, disseminated evidence and communications around automation tools may  
446 benefit from focusing on the combination of human and ML effort, rather than the replacement of

447 human insight. Finally, proponents should prioritize communication of transparency in automation  
448 methods and on strengthening automation competency and familiarity among EBM professionals.

## 449 [List of Abbreviations](#)

450 ML: Machine Learning

451 NICE: National Institute for Health Care and Excellence (UK)

452 NHMRC: National Health and Medical Research Council (Australia)

453 EBM: Evidence-based Medicine

454 ICASR: International Collaboration for the Automation of Systematic Reviews

## 455 [Declarations](#)

### 456 [Ethics approval and consent to participate](#)

457 This study was prospectively approved according to University College London's ethical standards.

458 Participants all consented to have their interviews recorded and were free to withdraw at any time.

### 459 [Consent for publication](#)

460 Not applicable (all data is anonymized).

### 461 [Availability of data and materials](#)

462 A COREQ checklist is available as supplementary material.

463 Anonymized data will be securely stored according to UCL Institute of Education guidance and may be

464 made available upon reasonable request.

## 465 Competing interests

466 The authors declare that they have no competing interests.

## 467 Funding

468 This research was jointly funded by a UCL and Monash PhD Studentship.

## 469 Authors' contributions

470 All members of the research team contributed to personal networking for recruitment. The interview  
 471 instrument was developed by Ms Arno with input from Dr Turner and Dr Elliott. Interviews,  
 472 transcription, and coding were completed by Ms Arno. Coding was checked over by Dr Thomas and by  
 473 Dr Elliott. The manuscript was primarily prepared by Ms Arno, with secondary input from all other  
 474 authors.

## 475 Acknowledgements

## 476 References

- 477 1. Bastian, H., P. Glasziou, and I. Chalmers, *Seventy-five trials and eleven systematic reviews a day: how will we ever keep up?* PLoS medicine, 2010. **7**(9): p. e1000326.
- 478 2. Shojania, K.G., et al., *How quickly do systematic reviews go out of date? A survival analysis.* Annals of internal medicine, 2007. **147**(4): p. 224-233.
- 479 3. Elliott, J.H., et al., *Living systematic reviews: an emerging opportunity to narrow the evidence-practice gap.* PLoS medicine, 2014. **11**(2): p. e1001603.
- 480 4. Marshall, I.J. and B.C. Wallace, *Toward systematic review automation: a practical guide to using machine learning tools in research synthesis.* Systematic reviews, 2019. **8**(1): p. 163.
- 481 5. Tsafnat, G., et al., *The automation of systematic reviews.* BMJ: British Medical Journal (Online), 2013. **346**.
- 482 6. Thomas, J., et al., *Living systematic reviews: 2. Combining human and machine effort.* Journal of Clinical Epidemiology, 2017. **91**: p. 31-37.
- 483 7. van Altena, A.J., R. Spijker, and S.D. Olabarriga, *Usage of automation tools in systematic reviews.* Research Synthesis Methods, 2019. **10**(1): p. 72-82.
- 484 8. Cleo, G., et al., *Usability and acceptability of four systematic review automation software packages: a mixed method design.* Systematic Reviews, 2019. **8**(1): p. 145.
- 485 9. Rogers, E.M., *Diffusion of innovations.* 2010: Simon and Schuster.
- 486 10. NVivo qualitative data analysis software. 2018, QSR International Pty Ltd.

- 495 11. Braun, V., V. Clarke, and G. Terry, *Thematic analysis*. Qual Res Clin Health Psychol, 2014. **24**: p.  
 496 95-114.
- 497 12. Gale, N.K., et al., *Using the framework method for the analysis of qualitative data in multi-*  
 498 *disciplinary health research*. BMC Medical Research Methodology, 2013. **13**(1): p. 117.
- 499 13. Thomas, J., *Diffusion of innovation in systematic review methodology: why is study selection not*  
 500 *yet assisted by automation*. OA Evidence-Based Medicine, 2013. **1**(2): p. 1-6.
- 501 14. O'Connor, A.M., et al., *Still moving toward automation of the systematic review process: a*  
 502 *summary of discussions at the third meeting of the International Collaboration for Automation of*  
 503 *Systematic Reviews (ICASR)*. Systematic reviews, 2019. **8**(1): p. 57.
- 504 15. David, H., *Why are there still so many jobs? The history and future of workplace automation*.  
 505 Journal of economic perspectives, 2015. **29**(3): p. 3-30.
- 506 16. Acharya, A.S., et al., *Sampling: Why and how of it*. Indian Journal of Medical Specialties, 2013.  
 507 **4**(2): p. 330-333.
- 508 17. NICE. *Gender pay gap report*. 2020; Available from: [https://www.nice.org.uk/about/who-we-](https://www.nice.org.uk/about/who-we-are/corporate-publications/gender-pay-gap-report)  
 509 [are/corporate-publications/gender-pay-gap-report](https://www.nice.org.uk/about/who-we-are/corporate-publications/gender-pay-gap-report).
- 510 18. Rashid, A., et al., *Patient and public involvement in the development of healthcare guidance: an*  
 511 *overview of current methods and future challenges*. The Patient-Patient-Centered Outcomes  
 512 Research, 2017. **10**(3): p. 277-282.

513