

Opinion: Leveraging AI to improve evidence synthesis in conservation

Authors

- Oded Berger-Tal*, Mitrani Department of Desert Ecology, Jacob Blaustein Institutes for Desert Research, Ben Gurion University, Israel.
<https://orcid.org/0000-0002-7396-456X>
- Bob B.M. Wong*, School of Biological Sciences, Monash University, Victoria 3800, Australia.
<https://orcid.org/0000-0001-9352-6500>
- Carrie Ann Adams, Department of Fish, Wildlife, and Conservation Biology, Colorado State University, 1474 Campus Delivery, Colorado State University, Fort Collins, CO, USA.
<https://orcid.org/0000-0001-7381-0555>
- Daniel T. Blumstein, Department of Ecology & Evolutionary Biology, University of California, 621 Young Drive South, Los Angeles, CA 90095-1606, USA.
<https://orcid.org/0000-0001-5793-9244>
- Ulrika Candolin, Organismal and Evolutionary Biology Research Programme, University of Helsinki, 00014 Helsinki, Finland.
<https://orcid.org/0000-0001-8736-7793>
- Matthew J. Gibson, Evolution & Ecology Research Centre, Centre for Ecosystem Science, and School of Biological, Earth and Environmental Sciences, University of New South Wales, Sydney, NSW 2052, Australia.
<https://orcid.org/0000-0003-1856-8959>
- Alison L. Greggor, Conservation Science and Wildlife Health, San Diego Zoo Wildlife Alliance, Escondido, CA, USA
<https://orcid.org/0000-0003-0998-618X>
- Malgorzata Lagisz, Evolution & Ecology Research Centre, Centre for Ecosystem Science, and School of Biological, Earth and Environmental Sciences, University of New South Wales, Sydney, NSW 2052, Australia.
<https://orcid.org/0000-0002-3993-6127>
- Biljana Macura, Stockholm Environment Institute (HQ), Box 24218, Stockholm, 10451, Sweden.
<https://orcid.org/0000-0002-4253-1390>
- Catherine J. Price, School of Life & Environmental Sciences, University of Sydney, NSW 2006, Australia.
<https://orcid.org/0000-0001-9631-2479>
- Breanna J. Putman, Department of Biology, California State University, San Bernardino, CA, United States
<https://orcid.org/0000-0002-1079-5519>

- 40 ● Lysanne Snijders, Behavioural Ecology Group, Wageningen University &
41 Research, De Elst 1, 6708 WD, Wageningen, The Netherlands.
42 <https://orcid.org/0000-0003-0911-3418>
43 ● Shinichi Nakagawa, Evolution & Ecology Research Centre, Centre for Ecosystem
44 Science, and School of Biological, Earth and Environmental Sciences, University
45 of New South Wales, Sydney, NSW 2052, Australia.
46 <https://orcid.org/0000-0002-7765-5182>

47

48 * Equal contribution

49

50 Correspondence: Shinichi Nakagawa (s.nakagawa@unsw.edu.au), Oded Berger-Tal
51 (bergerod@bgu.ac.il), Bob B.M. Wong (bob.wong@monash.edu)

52

53 Keywords: artificial intelligence, biodiversity conservation, evidence synthesis, large
54 language models, systematic reviews

55

56

57 **Abstract**

58 Systematic evidence syntheses (systematic reviews and maps) summarize knowledge
59 and are used to support decisions and policies in a variety of applied fields, from
60 medicine and public health to biodiversity conservation. However, conducting these
61 exercises in conservation is often expensive and slow, which can impede their use and
62 hamper progress in addressing the biodiversity crisis. With the explosive growth of large
63 language models (LLM) and other forms of artificial intelligence (AI), we discuss the
64 promise and perils associated with their use. We conclude that, when judiciously used,
65 AI has the potential to speed up and hopefully improve the process of evidence
66 synthesis, which can be particularly useful for underfunded applied fields such as
67 conservation science.

68

69

70 Main text

71 Biodiversity conservation needs evidence synthesis

72 Biodiversity conservation requires rapid decisions that, ideally, are made with the best
73 available scientific evidence. Systematic **evidence syntheses (systematic reviews**
74 **and systematic maps**; see Glossary)—through rigorous, transparent and repeatable
75 methods—are recognized as the gold standard for cataloging, collating and
76 synthesizing the available evidence to support decision making from public health to
77 environmental management and conservation ([1], Box 1). However, conducting
78 systematic evidence syntheses can often be expensive and slow [2]. With the
79 conservation literature growing exponentially, the endeavor can rapidly become
80 unmanageable for human reviewers and irrelevant for managers and policy advisors
81 that look for timely scientific evidence to support their decisions [3]. Several solutions
82 have been proposed for more rapid forms of evidence synthesis (e.g., [1,4,5]), which
83 raises the challenge of potentially having to trade speed of the review process with
84 comprehensiveness or exhaustiveness, thus reducing the reliability of the review
85 findings.

86
87 **Artificial intelligence (AI) and machine learning (especially deep learning)** tools are
88 currently revolutionizing how evidence is synthesized in biomedical sciences [6]. While
89 there are key differences between biomedicine and conservation research, in this
90 opinion piece, we make the case that AI tools can also dramatically improve evidence
91 syntheses and decision-making for biodiversity conservation. We do so by first
92 highlighting the potential role of AI in biodiversity conservation, and then discussing the
93 benefits and challenges of using AI, especially **large language models (LLMs)** in this
94 field. Because these tools are still in their infancy [7,8], we clarify their role in
95 synthesizing text-based scientific evidence for conservation decision-making, and
96 propose suggestions for responsible and ethical use of AI in conservation science.

97

98 Artificial intelligence is revolutionizing conservation science

99 Artificial Intelligence, initially the realm of science fiction, is now firmly entrenched in our
100 daily lives, and continues to revolutionize the way we interact with each other, our world
101 and even the universe. In conservation science, AI technologies are already extensively
102 and creatively deployed in a myriad of ways for research and management purposes —
103 from AI tools to expose online wildlife trafficking [9] and drones with machine and deep
104 learning capabilities to identify, track and monitor wild animals [10], to the use of
105 interactive robots to understand and control the spread of invasive species [11]. By
106 contrast, using rapidly emerging AI tools, such as LLMs, to allow for more efficient

107 evidence synthesis to support conservation decision-making, holds great potential but is
108 still relatively new.

109
110 Machine learning algorithms employ **artificial neural networks** that are trained by large
111 amounts of data (referred to as a corpus). Whereas simple machine learning is an
112 approach to classify and facilitate discrimination between two or more entities, LLMs are
113 able to recognize, summarize, translate, predict and generate text without any training
114 or only a few instructions as a form of **prompts** (known as **zero-shot or few-shot**
115 **learning**). In the medical sciences, where evidence synthesis methods are well
116 developed and widely used, recent studies demonstrate the promising role that AI tools
117 can play in carrying out rapid and extensive literature reviews [8,12]. At the same time,
118 there is also discourse around potential challenges and limitations regarding the
119 usefulness of these platforms [7,13–15].

120 **Benefits and challenges of using AI for evidence synthesis**

121 *Speed*

122 Conservation science is a race against time. Employing AI and LLM tools can reduce
123 the time required to perform systematic evidence syntheses by assisting in various
124 stages of the work [6], including communicating the results to relevant stakeholders [3].
125 Researchers have shown that the use of LLM tools can substantially shorten, by as
126 much as six-fold, the time spent screening relevant research ([8,12,13], Box 2). LLMs
127 could also be applied to (meta)data extraction from relevant studies and summarize a
128 collection of articles more efficiently [8,16,17]. At present, different AI tools have
129 different limits to the amount of data that can be inputted into them or processed by
130 them. Some free versions of AI tools may be swamped by large screening tasks [17],
131 which could limit their use by funding-restricted conservation agencies. Speed is
132 desirable, but without expert oversight there are likely to be issues with accuracy and
133 reliability by increasing the pace of evidence syntheses (i.e., a human-in-the-loop, HITL
134 process is necessary).

135 136 *Comprehensiveness, accuracy and reliability*

137 Systematic evidence syntheses aim to reduce human bias in the assessment of
138 scientific evidence, but human biases (e.g., selection and language biases; [18]) and
139 inconsistencies among human reviewers in study selection and data extraction, are
140 known issues in these syntheses [19]. Using LLM tools can assist in reducing these
141 human biases. For example, by improving prompts, Spillias et al. ([13]) were able to
142 increase the accuracy of screening with ChatGPT (reducing type II errors to < 1%). By
143 helping locate potentially useful gray literature sources, which can be a critical source of
144 biodiversity conservation evidence [20,21], LLMs can help further reduce the effects of

145 publication bias on review comprehensiveness, and can act as a second or third non-
146 human reviewer to tackle screening inconsistencies [13].

147

148 While AI tools may reduce some human biases, they can introduce errors. LLMs can
149 miss important and relevant articles during screening [8] and, more broadly, the
150 reliability of different AI tools can vary greatly throughout the synthesis process [22,
151 Table 1]. Missing relevant information may be especially problematic in conservation
152 research where the best solutions are often context-dependent [23], which can lead to
153 incorrect management guidance. AI tools may also generate overconfident and
154 potentially erroneous conclusions and create harm in real-world applications [17].
155 Misinterpretation errors, where text is improperly summarized, creates an improper
156 understanding of the content. Fabrication errors, where a summary includes information
157 not in the original text, refer to a broad class of 'hallucinations' that are well-known
158 outputs from LLMs. Attribute errors relate to any non-key elements in the review
159 question (e.g., the mis-evaluation of the number of interventions or treatments). Thus,
160 substantial human validation of LLM outputs is essential at each stage of review
161 construction (i.e., HITL; [8]).

162

163 *Complexity*

164 Compounding the problem of reliability, conservation research is characterized by some
165 unique complexities. Specifically, the field is highly heterogeneous, and includes studies
166 that span a variety of ecosystems and species applying a panoply of study designs and
167 dependent variables that can be measured in various ways (c.f., [24]). The field often
168 draws on evidence from many different disciplines, from psychology and physiology to
169 biochemistry and animal behavior. In addition, the language and terminology used in
170 conservation can be highly inconsistent, with many synonyms for similar terms [25]. For
171 example, the terms invasive, introduced, exotic, alien or non-native species, weed, and
172 pest can all have the same meaning, depending on context. Finally, the majority of
173 published conservation research does not test practical, real-world interventions [26].
174 Evidence producers must therefore make fine-grained decisions about where academic
175 studies are sufficiently solution-oriented or relevant, while trudging through disparate
176 and highly variable gray literature. Such complexities and nuances need to be taken into
177 account in developing search prompts, screening and oversight of results, and when
178 models are updated to ensure reliability and accuracy of results generated by LLMs
179 [27]. However, robust methods for dealing with such complexities are yet to be
180 developed.

181

182 *Relevance over time*

183 The evidence base for conservation is rapidly accumulating and evidence syntheses
184 can quickly become outdated. In a rapidly changing world, the effectiveness of
185 interventions might also change with time. Thus, systematic reviews that are not

186 regularly updated may lead to significant inaccuracies over time [28,29]. **Living**
187 **systematic reviews** have been developed to provide high quality, up-to-date online
188 summaries that incorporate relevant new evidence as it becomes available [28,30].
189 Such reviews require continuous work and a level of commitment that is often hard to
190 achieve. Here, LLMs can be used to support living reviews and ensure that the
191 evidence-base remains up to date with minimal human effort [30,31]. However, as the
192 outputs of LLMs may change over time (because the algorithms and training sets
193 change), their performance will require human evaluation.

194

195 *Inclusivity*

196 In our view, one of the major benefits of using LLMs in synthesis is their ability to find
197 conservation evidence from across the globe, particularly in languages other than
198 English [32]. Most of the world's remaining biodiversity is found in the Global South, yet
199 most scientific evidence to inform decision-making comes from authors in the Global
200 North and is published in English [33]. Local studies from the Global South are often
201 missed or discarded from reviews if they are not written in English.

202

203 By translating languages, AI tools can make all stages of the review process more
204 inclusive (Box 1). For example, a review on community-based fisheries management
205 focusing on the Pacific Islands [13] benefited from AI rapidly providing a list of non-
206 English relevant terms to be integrated into the search string and yielding additional
207 articles not previously identified by the original search. AI-suggested terms should,
208 however, be checked by proficient speakers of the language in question before
209 inclusion in the search string.

210

211 Nevertheless, it is important to emphasize that AI tools require accessible digitized
212 information. Moreover, the original training to create LLMs requires sufficiently large
213 data sets that currently exclude most of the world's languages [34,35]. Therefore,
214 exclusively relying on AI for information means that some traditional and local
215 knowledge may be ignored. This process could reduce the effectiveness of
216 conservation interventions at the local scale and widen the divide between conservation
217 agencies and local communities [36]. In this respect, we emphasize that effective
218 conservation work relies just as heavily on building strong relationships with the relevant
219 stakeholders as using the most accurate scientific evidence (e.g., [37]). The use of AI
220 may alienate local collaborators if not conveyed and properly communicated to all
221 stakeholders and rightsholders.

222 **Ethical considerations**

223 The question, in our view, is not whether AI tools will/should be used in conservation
224 science (the singularity is nigh!), but rather how they are used. Issues of data privacy

225 and informed consent created by emerging AI technologies can be exacerbated through
226 their use in systematic evidence syntheses. People may not wish that their published
227 data are used for AI training, or repurposed and applied to new problems. In this regard,
228 continuous effort to actively engage various stakeholders in the synthesis process is
229 even more crucial in the context of AI application to evidence synthesis.

230

231 A well-recognized concern with using AI is the presence of (algorithmic) biases that
232 result from factors such as the unknown data quality and representativeness in training
233 corpus [38,39]. As previously discussed, it is likely that documents written in English
234 and from developed countries form the bulk of the training corpus — this may limit the
235 nature of responses to specific queries and enhance existing biases. Therefore, there is
236 an urgent need for culturally sensitive multi-lingual LLMs [40]. Moreover, in the current
237 LLM landscape there is a lack of transparency around algorithm development and
238 reporting related to decisions algorithms make during the review process. Lack of
239 transparency leads to limited peer scrutiny and accountability in AI-supported evidence
240 syntheses and prevents equitable and responsible development of AI.

241

242 Hence, the best practice moving forward is to be explicitly clear about how AI is being
243 used in evidence syntheses, which may include detailed reporting of the prompts and
244 instructions given to an LLM and how it was tested for replicability and reliability. This
245 ensures transparency and reproducibility to some extent. Repeatability can be limited
246 because models are probabilistic and constantly updated with new data. Thus, multiple
247 runs of the same model over time may produce different responses. This is a challenge
248 that requires future research to fully understand its impact on evidence synthesis and,
249 ultimately, on conservation management decisions.

250 **Concluding remarks**

251 AI is not a silver bullet and conducting a reliable evidence synthesis requires a lot of
252 work and will always be time-consuming and require attention to detail (Box 3).
253 However, AI tools can help improve the location and consideration of gray literature and
254 evidence in a variety of languages that were not traditionally included in syntheses. AI
255 may make evidence synthesis faster, more accessible, and inclusive to a greater
256 number of researchers. Although decision-making in conservation involves more than
257 just scientific evidence, expanding the availability of the information base will increase
258 opportunities for developing informed policies and management actions (see
259 Outstanding Questions).

260

261 More broadly, while we have focused on how AI tools can be used to synthesize
262 evidence for biodiversity conservation, we suggest that ecologists and evolutionary

263 biologists, more broadly, can also benefit from using these tools to efficiently identify the
264 state of knowledge in their respective disciplines.

265
266

267 **Acknowledgements:** This paper emerged from a workshop conducted at Monash
268 University's Prato Centre, partially supported by Monash University and Ben Gurion
269 University of the Negev (to OB-T and BMW) through a grant from the Pratt
270 Foundation. We thank the staff at the Prato Centre for their wonderful hospitality and
271 support. In addition, we acknowledge the following funding agencies for financial
272 support: the Australian Research Council (FT190100014 and DP220100245 to BMW,
273 DP210100812 and DP230101248 to ML and SN, and DE220101316 to CP), a NASA
274 Biodiversity grant (#80NSSC21K114 to CA), the Swedish Cultural Foundation in Finland
275 (Nr 179446 to UC), and the Netherlands Organisation for Scientific Research
276 (VI.Veni.192.018 to LS).

277

278 **AI disclosure:** We used Chat GPT 4.0, accessed through Microsoft Edge, to develop a
279 list of benefits of AI for decision making for conservation. We largely ignored the specific
280 list and wrote this paper collectively using purely human-synthesized knowledge. Chat
281 GPT 4.0 was also used as a resource to help understand key ideas and tools used in
282 this rapidly growing field. The authors are fully responsible for the content contained in
283 this manuscript.

284

285 **Declaration of interests**

286 No interests are declared.

287

288

289 References

- 290
- 291 1. Collaboration for Environmental Evidence (2022) Guidelines and Standards for Evidence
 292 Synthesis in Environmental Management. Version 5.1. (Pullin, A.S. *et al.*, eds). URL:
 293 www.environmentalevidence.org/information-for-authors
- 294 2. Haddaway, N.R. and Westgate, M.J. (2019) Predicting the time needed for environmental
 295 systematic reviews and systematic maps. *Conserv. Biol.* 33, 434–443
- 296 3. Tyler, C. *et al.* (2023) AI tools as science policy advisers? The potential and the pitfalls.
 297 *Nature* 622, 27–30
- 298 4. Haby, M.M. *et al.* (2016) What are the best methodologies for rapid reviews of the research
 299 evidence for evidence-informed decision making in health policy and practice: a rapid review.
 300 *Health Res. Policy and Syst.* 14, 83
- 301 5. Sutherland, W.J. and Wordley, C.F.R. (2018) A fresh approach to evidence synthesis. *Nature*
 302 558, 364–366
- 303 6. Jimenez, R.C. *et al.* (2022) Machine learning computational tools to assist the performance of
 304 systematic reviews: A mapping review. *BMC Med. Res. Methodol.* 22, 322
- 305 7. Qureshi, R. *et al.* (2023) Are ChatGPT and large language models “the answer” to bringing
 306 us closer to systematic review automation? *Syst. Rev.* 12, 72
- 307 8. Blaizot, A. *et al.* (2022) Using artificial intelligence methods for systematic review in health
 308 sciences: A systematic review. *Res. Synth. Methods* 13, 353–362
- 309 9. Cardoso, A.S. *et al.* (2023) Detecting wildlife trafficking in images from online platforms: A
 310 test case using deep learning with pangolin images. *Biol. Conserv.* 279, 109905
- 311 10. Couzin, I.D. and Heins, C. (2023) Emerging technologies for behavioral research in
 312 changing environments. *Trends Ecol. Evol.* 38, 346–354
- 313 11. Polverino, G. *et al.* (2022) Ecology of fear in highly invasive fish revealed by robots. *iScience*
 314 25, 103529
- 315 12. van Dijk, S.H.B. *et al.* (2023) Artificial intelligence in systematic reviews: promising when
 316 appropriately used. *BMJ Open* 13, e072254
- 317 13. Spillias, S. *et al.* (2023) Human-AI Collaboration to Identify Literature for Evidence
 318 Synthesis. *Research Square Preprint*. DOI: 10.21203/rs.3.rs-3099291/v1
- 319 14. Zhu, J.-J. *et al.* (2023) ChatGPT and Environmental Research. *Environ. Sci. Technol.* 57,
 320 17667–17670
- 321 15. Demszky, D. *et al.* (2023) Using large language models in psychology. *Nat. Rev. Psychol.* 2,
 322 688–701
- 323 16. Shaib, C. *et al.* (2023) Summarizing, Simplifying, and Synthesizing Medical Evidence using
 324 GPT-3 (with Varying Success). *Proc. Conf. Assoc. Comput. Linguist. Meet.* 2, 1387–1407
 325 DOI: 10.18653/v1/2023.acl-short.119
- 326 17. Tang, L. *et al.* (2023) Evaluating large language models on medical evidence
 327 summarization. *npj Digit. Med.* 6, 1–8
- 328 18. Frampton, G. *et al.* (2022) Principles and framework for assessing the risk of bias for studies
 329 included in comparative quantitative environmental systematic reviews. *Environ Evid* 11, 12
- 330 19. Felson, D.T. (1992) Bias in meta-analytic research. *J. Clin. Epidemiol.* 45, 885–892
- 331 20. Haddaway, N.R. and Bayliss, H.R. (2015) Shades of grey: Two forms of grey literature
 332 important for reviews in conservation. *Biol. Conserv.* 191, 827–829
- 333 21. Amano, T. *et al.* (2023) The role of non-English-language science in informing national
 334 biodiversity assessments. *Nat. Sustain.* 6, 845–854
- 335 22. Zhao, Z. *et al.* (2021) Calibrate Before Use: Improving Few-shot Performance of Language
 336 Models. *Proc. Intern. Conf. Mach. Learn.* 139, 12697–12706
- 337 23. Christie, A.P. *et al.* (2020) Poor availability of context-specific evidence hampers decision-
 338 making in conservation. *Biol. Conserv.* 248, 108666

- 339 24. Christie, A.P. *et al.* (2020) Quantifying and addressing the prevalence and bias of study
340 designs in the environmental and social sciences. *Nat. Commun.* 11, 6377
- 341 25. Cheng, S.H. *et al.* (2018) Using machine learning to advance synthesis and use of
342 conservation and environmental evidence. *Conserv. Biol.* 32, 762–764
- 343 26. Williams, D.R. *et al.* (2020) The past and future role of conservation science in saving
344 biodiversity. *Conserv. Lett.* 13, e12720–e12720
- 345 27. Clusmann, J. *et al.* (2023) The future landscape of large language models in medicine.
346 *Commun. Med.* 3, 141
- 347 28. Brooker, J. *et al.* (2019) Guidance for the production and publication of Cochrane living
348 systematic reviews: Cochrane Reviews in living mode [Online].
349 [https://community.cochrane.org/sites/default/files/uploads/inline-](https://community.cochrane.org/sites/default/files/uploads/inline-files/Transform/201912_LSR_Revised_Guidance.pdf)
350 [files/Transform/201912_LSR_Revised_Guidance.pdf](https://community.cochrane.org/sites/default/files/uploads/inline-files/Transform/201912_LSR_Revised_Guidance.pdf)
- 351 29. Shojania, K.G. *et al.* (2007) *Updating Systematic Reviews*, Agency for Healthcare Research
352 and Quality (US), Rockville, MD. AHRQ Publication No. 07-0087.
- 353 30. Elliott, J.H. *et al.* (2014) Living Systematic Reviews: An Emerging Opportunity to Narrow the
354 Evidence-Practice Gap. *PLOS Medicine* 11, e1001603
- 355 31. Shackelford, G.E. *et al.* (2021) Dynamic meta-analysis: a method of using global evidence
356 for local decision making. *BMC Biol.* 19, 33
- 357 32. Amano, T. *et al.* (2021) Tapping into non-English-language science for the conservation of
358 global biodiversity. *PLoS Biol.* 19, e3001296
- 359 33. Christie, A.P. *et al.* (2021) The challenge of biased evidence in conservation. *Conservation*
360 *Biology* 35, 249–262
- 361 34. Joshi, P. *et al.* (2021) The State and Fate of Linguistic Diversity and Inclusion in the NLP
362 World. *arXiv Preprint* DOI:10.48550/arXiv.2004.09095
- 363 35. Ranathunga, S. and de Silva, N. (2022) Some Languages are More Equal than Others:
364 Probing Deeper into the Linguistic Disparity in the NLP World. *Proc. Conf. Asia-Pacific*
365 *Assoc. Comput. Linguist. Joint Conf. Nat. Lang. Process. (Volume 1: Long Papers)*, pp. 823–
366 848 URL: <https://aclanthology.org/2022.aacl-main.62>
- 367 36. Droz, L. *et al.* (2023) Multilingualism for pluralising knowledge and decision making about
368 people and nature relationships. *People Nat.* 5, 874–884
- 369 37. Chaplin-Kramer, R. *et al.* (2023) Transformation for inclusive conservation: evidence on
370 values, decisions, and impacts in protected areas. *Curr. Opin. Environ. Sustain.* 64, 101347
- 371 38. Hovy, D. and Prabhumoye, S. (2021) Five sources of bias in natural language processing.
372 *Lang. Linguist. Compass* 15, e12432
- 373 39. Chen, Y. *et al.* (2023) Human-Centered Design to Address Biases in Artificial Intelligence. *J.*
374 *Med. Internet Res.* 25, e43251
- 375 40. Ramesh, K. *et al.* (2023) Fairness in Language Models Beyond English: Gaps and
376 Challenges. *arXiv Preprint* DOI: 10.48550/arXiv.2302.12578
- 377 41. Fan, W. *et al.* (2023) Recommender Systems in the Era of Large Language Models (LLMs).
378 *arXiv Preprint* DOI: 10.48550/arXiv.2307.02046
- 379 42. O'Donoghue, O. *et al.* (2023) BioPlanner: Automatic Evaluation of LLMs on Protocol
380 Planning in Biology. *arXiv Preprint* DOI: 0.48550/arXiv.2310.10632
- 381 43. Khraisha, Q. *et al.* (2023) Can large language models replace humans in the systematic
382 review process? Evaluating GPT-4's efficacy in screening and extracting data from peer-
383 reviewed and grey literature in multiple languages. *arXiv Preprint*. DOI:
384 10.48550/arXiv.2310.17526
- 385 44. Michelson, M. *et al.* (2020) Artificial Intelligence for Rapid Meta-Analysis: Case Study on
386 Ocular Toxicity of Hydroxychloroquine. *J. Med. Internet Res.* 22, e20007
- 387 45. Valizadeh, A. *et al.* (2022) Abstract screening using the automated tool Rayyan: results of
388 effectiveness in three diagnostic test accuracy systematic reviews. *BMC Med. Res.*
389 *Methodol.* 22, 160

- 390 46. Chen, L. *et al.* (2023) How is ChatGPT's behavior changing over time? *arXiv Preprint* DOI:
391 10.48550/arXiv.2307.09009
- 392 47. Koehler, M. and Sauermann, H. (2023) Algorithmic Management in Scientific Research.
393 *SSRN Journal* DOI: 10.2139/ssrn.4497871
- 394 48. Bannach-Brown, A. *et al.* (2019) Machine learning algorithms for systematic review:
395 reducing workload in a preclinical review of animal studies and reducing human screening
396 error. *Syst. Rev.* 8, 23
- 397 49. Hill, J.E. *et al.* (2023) Methods for using Bing's AI-powered search engine for data extraction
398 for a systematic review. *Res. Synth. Methods* DOI: 10.1002/jrsm.1689
- 399 50. Waffenschmidt, S. *et al.* (2023) Increasing the efficiency of study selection for systematic
400 reviews using prioritization tools and a single-screening approach. *Syst. Rev.* 12, 161
- 401 51. Syriani, E. *et al.* (2023) Assessing the Ability of ChatGPT to Screen Articles for Systematic
402 Reviews. *arXiv Preprint* DOI: 10.48550/arXiv.2307.06464
- 403 52. Jardim, P.S.J. *et al.* (2022) Automating risk of bias assessment in systematic reviews: a
404 real-time mixed methods comparison of human researchers to a machine learning system.
405 *BMC Med. Res. Methodol.* 22, 167
- 406 53. Marshall, I.J. *et al.* (2017) Automating Biomedical Evidence Synthesis: RobotReviewer. *Proc*
407 *Conf. Assoc. Comput. Linguist. Meet.* 2017, 7–12
- 408 54. Gates, A. *et al.* (2018) Technology-assisted risk of bias assessment in systematic reviews: a
409 prospective cross-sectional evaluation of the RobotReviewer machine learning tool. *J. Clin.*
410 *Epidemiol.* 96, 54–62
- 411 55. Tsafnat, G. *et al.* (2014) Systematic review automation technologies. *Syst. Rev.* 3, 74
- 412 56. Gates, A. *et al.* (2021) Creating efficiencies in the extraction of data from randomized trials:
413 a prospective evaluation of a machine learning and text mining tool. *BMC Med. Res.*
414 *Methodol.* 21, 169
- 415 57. Mutinda, F.W. *et al.* (2022) Automatic data extraction to support meta-analysis statistical
416 analysis: a case study on breast cancer. *BMC Med. Inform. Decis. Mak.* 22, 158
- 417 58. West, R. *et al.* (2023) Using machine learning to extract information and predict outcomes
418 from reports of randomised trials of smoking cessation interventions in the Human Behaviour-
419 Change Project. *Wellcome Open Res.* 8, 452 DOI: 10.12688/wellcomeopenres.20000.1
- 420 59. Ho, I.M.K. *et al.* (2023) Using Machine Learning Algorithms to Pool Data from Meta-Analysis
421 for the Prediction of Countermovement Jump Improvement. *Int. J. Environ. Res. Public*
422 *Health* 20, 5881
- 423 60. Xu, J.-L. *et al.* (2022) Combining machine learning with meta-analysis for predicting
424 cytotoxicity of micro- and nanoplastics. *J. Hazard. Mater. Adv.* 8, 100175
- 425 61. Marshall, I.J. and Wallace, B.C. (2019) Toward systematic review automation: a practical
426 guide to using machine learning tools in research synthesis. *Syst. Rev.* 8, 163
- 427 62. Huang, J. and Tan, M. (2023) The role of ChatGPT in scientific communication: writing
428 better scientific review articles. *Am. J. Cancer Res.* 13, 1148–1154
- 429 63. van de Schoot, R. *et al.* (2021) An open source machine learning framework for efficient and
430 transparent systematic reviews. *Nat. Mach. Intell.* 3, 125–133
- 431 64. Lombaers, P. *et al.* (2023) Reproducibility and Data storage Checklist for Active Learning-
432 Aided Systematic Reviews. *PsyArXiv Preprint* DOI: 10.31234/osf.io/g93zf
- 433 65. Dicks, L.V. *et al.* (2014) Organising evidence for environmental management decisions: a
434 '4S' hierarchy. *Trends Ecol. Evol.* 29, 607–613
- 435 66. Orgeolet, L. *et al.* (2020) Can artificial intelligence replace manual search for systematic
436 literature? Review on cutaneous manifestations in primary Sjögren's syndrome.
437 *Rheumatology* 59, 811–81967.
- 438 67. Ouzzani, M. *et al.* (2016) Rayyan—a web and mobile app for systematic reviews. *Syst Rev*
439 5, 210

- 440 68. Adams, C.A. *et al.* (2021) Effects of artificial light on bird movement and distribution: a
 441 systematic map. *Environ. Evid.* 10, 37
- 442 69. Ali, S. *et al.* (2023) Explainable Artificial Intelligence (XAI): What we know and what is left to
 443 attain Trustworthy Artificial Intelligence. *Inf. Fusion* 99, 101805
- 444 70. Turpin, M. *et al.* (2023) Language Models Don't Always Say What They Think: Unfaithful
 445 Explanations in Chain-of-Thought Prompting. *arXiv Preprint*. DOI:10.48550/arXiv.2305.04388
- 446 71. Shi, F. *et al.* (2023) Large Language Models Can Be Easily Distracted by Irrelevant Context.
 447 *Proc. Intern. Conf. Mach. Learn.* 202, 31210-31227
- 448 72. O'Dea, R.E. *et al.* (2021) Preferred reporting items for systematic reviews and meta-
 449 analyses in ecology and evolutionary biology: a PRISMA extension. *Biol. Rev. Camb. Phil.*
 450 *Soc.* 96, 1695–1722
- 451 73. Haddaway, N.R. *et al.* (2018) ROSES RepOrting standards for Systematic Evidence
 452 Syntheses: pro forma, flow-diagram and descriptive summary of the plan and conduct of
 453 environmental systematic reviews and systematic maps. *Environ. Evid.* 7, 7
- 454 74. Susnjak, T. (2023) PRISMA-DFLLM: An Extension of PRISMA for Systematic Literature
 455 Reviews using Domain-specific Finetuned Large Language Models. *arXiv Preprint* DOI:
 456 10.48550/arXiv.2306.14905
 457
 458

459 Glossary

460

461 **Artificial intelligence (AI):** a machine or model which can perform what appears to
 462 require human intelligence. It also refers to a branch of computer science dedicated to
 463 creating these models. Recently, generative AI has gained much attention with its ability
 464 to create text, images, audio and other media.

465

466 **Artificial neural networks:** a method used in machine learning whereby the
 467 connections and strength of connections between a set of nodes (which are modeled
 468 after neurons in the brain) is iteratively modified to maximize some desired output (e.g.,
 469 a discrimination). Originally, these networks had several layers of nodes between input
 470 and output but deep learning models have many layers of nodes.

471

472 **Evidence syntheses:** involve a process of combining information from multiple studies
 473 on a specific topic and to inform decision making. The term is also used as an umbrella
 474 term for the family of reviews that include systematic reviews, systematic maps, rapid
 475 reviews, and reviews of reviews.

476

477 **Deep learning:** a type of machine learning that relies on multiple layers of connected
 478 nodes whose connections and weights are iteratively modified so as to maximize their
 479 ability to make discriminations or identifications. It requires a huge amount of training.

480

481 **Human-in-the-loop (HITL):** a process of model development in machine learning
482 where humans play an interactive and iterative role.

483

484 **Large language model (LLM):** a type of generative artificial intelligence created by a
485 deep learning, neural network trained on a large written corpus that can “understand”
486 human language and generate responses to specific queries.

487

488 **Living systematic reviews:** systematic reviews that are continuously updated that
489 incorporate new evidence as it is produced.

490

491 **Machine learning:** a process by which data are fed into neural network models which
492 are iteratively modified without specific instructions that permit the identification of
493 patterns in data.

494

495 **Prompts:** specific inputs or instructions to a LLM designed to elicit an answer. The
496 growing field of prompt engineering studies the characteristics of effective prompts
497 which in general should be specific, and constrained. Creating a role (‘you are a
498 fastidious researcher conducting a systematic review...’) can help improve output
499 accuracy.

500

501 **Systematic review:** a formal and highly structured process to comprehensively,
502 rigorously and transparently collate and synthesize evidence, including the academic
503 and gray literature sources. Can be used to support policy formation and biodiversity
504 management decisions.

505

506 **Systematic map:** comprehensive catalogues of the literature on a broad topic of
507 interest. Systematic maps follow the same step-wise process as systematic reviews, but
508 they tackle broader questions, and their final output is a narrative report and a
509 searchable catalogue of the literature that can be used to identify areas where evidence
510 is lacking or is under-represented (knowledge gaps), or areas with sufficient evidence to
511 conduct full synthesis (knowledge clusters) **Zero-shot or few-shot learning:** a direct
512 query to an existing LLM is referred to as a zero-shot query where the results of zero-
513 shot queries are based entirely on the information already contained in the LLM. By
514 contrast, few-shot learning requires some additional data, for instance, where the LLM
515 is provided a list of papers that, based on their title and abstract, that should be included
516 or excluded from a systematic review.

517

518 Table 1: AI tools and platforms for evidence synthesis ^a

Stage of synthesis	Example tools and Platforms ^b	Opportunities	Potential challenges and considerations
Identify and formulate review questions	<ul style="list-style-type: none"> • Gemini (Google DeepMind; https://gemini.google.com/) • Scite (scite; https://scite.ai/) 	Facilitate question formulation through assistance with brainstorming and refinement [7]	Some stakeholders might feel disengaged or excluded by the process, potentially hampering innovation and even reinforcing existing biases [7,41]
Draft review protocol	<ul style="list-style-type: none"> • Gemini (Google DeepMind; https://gemini.google.com/) • ChatGPT (OpenAI, https://chat.openai.com) 	Assist in creating a good initial outline and, hence, speeding up the process for protocol writing [7,42]	Risk of 'hallucinations' may cast doubt on protocol accuracy [16,17]; Protocol may lack details and/or correct references [16]
Search for evidence	<ul style="list-style-type: none"> • Elicit (Elicit; https://elicit.com/) • Scite (scite; https://scite.ai/) • Consensus (Consensus; https://consensus.app/) • Scispace (PubGenius Inc; https://typeset.io/) • ConnectedPapers (Connected Papers; https://www.connectedpapers.com/) • Inciteful (Weishun, M. 2024; https://inciteful.xyz/) • Litmaps (Litmaps Ltd; https://www.litmaps.com) • Gemini (Google DeepMind; https://gemini.google.com/) • ChatGPT (OpenAI, https://chat.openai.com) 	Help with suggesting and finding a variety of gray literature sources, including in different languages [43]; Suggest alternative terms for the search [7]; Help to incorporate evidence as it becomes available [44]	Inconsistent and incomplete search terms that can reduce search efficiency and increase the potential for selection bias [45]; Changes to the algorithm may change search results [7,46]; Search results may be probabilistic, erroneous, and not repeatable [7]; Can only make use of digitized knowledge [47]

<p>Include relevant studies</p>	<ul style="list-style-type: none"> • Rayyan (Ouzanni et al. 2016; https://www.rayyan.ai) • Abstrackr (Brown University; http://abstrackr.cebm.brown.edu/account/login) • DistillerSR DistillerSR Inc; https://www.distillersr.com/) • EPPI-Reviewer (EPPI Centre; eppi.ioe.ac.uk/EPPIReviewer-Web) • SWIFT-Active Screener (Sciome; https://www.sciome.com/swift-activescreener/) • ASReview (ASReview Lab; https://asreview.nl/) • Silivi (A-Evidence ApS; https://www.silvi.ai/) 	<p>Substantially reduce screening time [Box 2]; In the case of double screening, act as the second reviewer to tackle screening inconsistencies [48,49]</p>	<p>May inadvertently pass on relevant studies [50,51]; Changes to the algorithm may change screening results [7,46]; Lack of transparency around algorithm development and decision-making [52]; Screening decisions may be probabilistic and not repeatable [7]</p>
<p>Critically appraise studies</p>	<ul style="list-style-type: none"> • RobotReviewer [53] (https://www.robotreviewer.net/) • Elicit (Elicit; https://elicit.com/) 	<p>Speed up an otherwise very time-consuming process [53,54]</p>	<p>Difficulties in dealing with more complex and diverse study designs and different reporting styles [55]; Interpretation and extraction errors [16,56]; Lack of transparency around algorithm development and decision-making [52]</p>

Extract data	<ul style="list-style-type: none"> • Scispace (PubGenius Inc; https://typeset.io/) • RobotReviewer [53] (https://www.robotreviewer.net/) • SWIFT-Review (Sciome; https://www.sciome.com/swift-review/) • Silivi (A-Evidence ApS; https://www.silvi.ai/) • ExaCT (https://exact.cluster.gctools.nrc.ca/ExactDemo/intro.php) • Elicit (Elicit; https://elicit.com/) 	Efficient at extracting data and metadata (e.g. moderators and study descriptors) [53,57]	Difficulties in dealing with more complex and diverse study designs and different reporting styles [53,55,57]; Interpretation and extraction errors [16,56]; Lack of transparency around algorithm development and decision-making [52]; May not be reliable in obtaining effect sizes [58]
Synthesize data/study findings	<ul style="list-style-type: none"> • ChatGPT (OpenAI, https://chat.openai.com) • Gemini (Google DeepMind; https://gemini.google.com/) 	Potentially efficient at running simple quantitative syntheses (meta-analysis) of evidence as well as narratively synthesizing study findings [59,60]	Sophisticated quantitative (e.g. meta-regression) synthesis is still difficult to conduct [59,61]
Report findings	<ul style="list-style-type: none"> • ChatGPT (OpenAI, https://chat.openai.com) • Scispace (PubGenius Inc; https://typeset.io/) 	Efficient at scientific communication as it can assist scientists in improving their writing style by analyzing text and provide suggestions for improvements [14,62]	Lack of transparency around algorithm development and decision-making [63,64]

519 ^a We highlight both opportunities, as well as potential challenges and considerations. In regard to the latter, many of the challenges
520 we have identified can be resolved by having humans-in-the-loop and greater procedural transparency. Stages of synthesis mirror
521 those outlined in Figure I in Box 1.

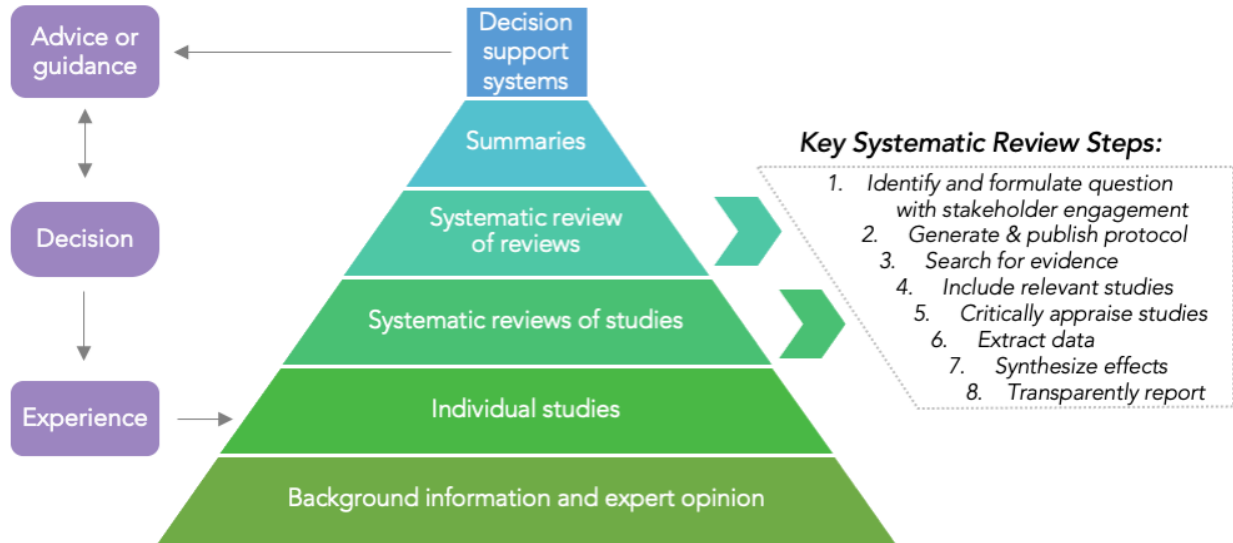
522 ^b A non-exhaustive list with an emphasis on new and popular platforms

523 **Box 1: Evidence hierarchy for decision support in conservation with**
524 **AI**

525
526 For scientific evidence to be useful or usable, information must be distilled,
527 amalgamated and translated from a large collection of individual studies to an output
528 that can inform decision-making. Figure I illustrates how different types of knowledge,
529 information and expert opinions, primary (individual studies) and secondary research
530 (e.g., systematic reviews and review of reviews) feed into decision support systems
531 (i.e., tools that provide different scenarios and logical sets of steps to assist with
532 decision making; [65]). Outputs from these systems help create evidence-informed
533 advice and guides. The pyramid demonstrates how, at each step, the scientific evidence
534 gradually becomes more “condensed” and hence more accessible to conservation
535 decision makers.

536
537 Each step of evidence synthesis could potentially be supported and expedited by AI and
538 LLMs, including: 1) question formulation, 2) protocol generation, 3) literature search, 4)
539 screening to select relevant papers (including deduplication), 5) critical appraisal of
540 included studies, 6) data extraction, 7) synthesizing information and 8) transparent
541 reporting (Figure I). Recently, Jimenez and colleagues ([6]) identified 63 machine-
542 learning tools for systematic evidence syntheses. They showed that most of the
543 currently available tools primarily support the three review stages: searching, screening
544 and data extraction. For example, *BIBOT* uses keywords to search and retrieve relevant
545 papers from PubMed [66], while *Rayyan* facilitates screening by reordering papers in
546 the order of relevance, learning from included and excluded papers [67] (also Box 2).
547 None of the tools in their review used LLMs, but LLMs can immediately be used in these
548 three stages and more. A generative AI platform, *Elicit* (elicit.com), for instance, can
549 extract information and summarize pdf documents.

550
551 In addition, LLMs can facilitate “Summaries” turning long academic documents (such as
552 systematic reviews) into distilled key messages for policy and practice. Furthermore,
553 LLMs can help create algorithms and software for decision support systems [3].



554

555

556

557

558

Figure I: Hierarchy of scientific evidence used in conservation decision-making (Modified and redrawn from Dicks et al. [65]).

559 **Box 2:** Speeding up screening with AI: a case study

560
561 There are a number of AI-assisted article screening tools, most of which use re-ordering
562 algorithms that learn from included/excluded articles as researchers screen based on
563 title and abstract. More recently, large language models (LLMs) have been suggested to
564 be used for such screening [13]. We tested both types: Rayyan.ai (re-ordering
565 algorithm) and GPT 3.5 (LLM) to screen 11,270 article search records from the Web of
566 Science for relevance to the question: *how does artificial light affect bird movement and*
567 *distribution?* These articles were manually screened by Adams *et al.* [68] (Figure 1).

568
569 Rayyan.ai's relevance ratings could have reduced the manual screening burden at the
570 title/abstract level by over 80%, with accuracy comparable to a human-alone screening.
571 We provided initial training data by classifying 46 articles we knew to be relevant as
572 "include" and classified 46 additional articles as "exclude". Rayyan computed relevance
573 ratings for the remaining articles, and we sorted them by relevance and screened the
574 first 100. We then recomputed the ratings, re-sorted the records, and screened the next
575 100 articles. We repeated the process until no additional relevant articles were found,
576 which occurred at ~ 2,200 articles. This method identified 169 (97%) out of 174 relevant
577 articles in the screening dataset after screening less than 20% of the articles. Notably,
578 this process yielded 5 articles missed by a human screener during the original
579 screening process, meaning that the human-alone and this AI-assisted method
580 (Rayyan.ai) had equivalent false negative rates in this case (2.9%).

581
582 For GPT 3.5, we used the following prompt "*Classify the given research paper as*
583 *worthy of inclusion or exclusion... The paper should be classified as "include" or*
584 *"exclude". You are a careful and thorough researcher conducting a systematic review of*
585 *the effect of artificial light on bird movement and distribution. Given a title and an*
586 *abstract of a research paper, your task is to determine whether the paper meets the*
587 *criteria for inclusion in a review study."* Following this message, this prompt also
588 included the published abstract along with screening criteria. For the initial run (i.e. zero-
589 shot learning) it retrieved 66 of 215 relevant articles (30%). For the second run, we
590 provided 46 included and excluded articles, and GPT 3.5 was able to retrieve 200 out of
591 215 (93%) articles. It took 2.5 hours for each run to screen 11,270 articles.

592



593

594 Figure I: Many studies have investigated the relationship between artificial light at night
595 and bird movements (credit: JoshuaWoroniecki)

596

597

598

599 **Box 3: Guiding principles for responsible AI use in evidence**
600 **syntheses for conservation**

601
602 Acceptable practices of using AI are evolving rapidly. For example, AI has been used to
603 improve writing for years (many already use Grammarly or Microsoft Grammar Checker)
604 but some publishers currently limit or prohibit LLM-produced text from being used in
605 papers. With this state of flux in mind, we make the following recommendations (Figure
606 I).

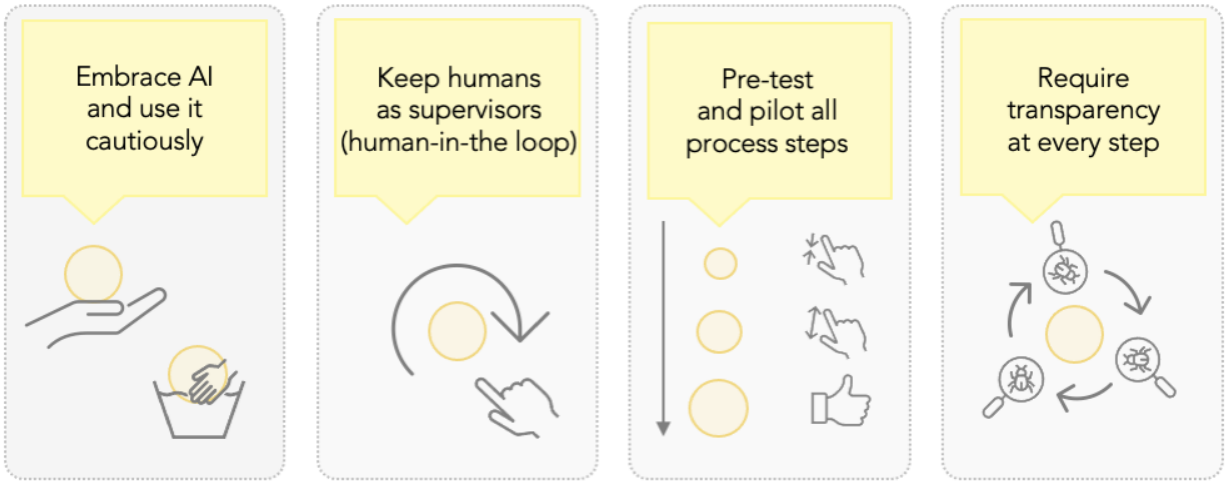
607
608 First, while AI tools offer considerable promise, use them cautiously. We do not
609 currently understand, in various contexts, its precision, accuracy, specificity, or reliability
610 and the developers themselves are unclear about how some AI tools and models work
611 [69]. As these tools are applied to specific conservation issues, effort will have to be
612 allocated to estimate these sources of error and optimize algorithms [70,71].

613
614 Second, view AI tools as a research assistant—it is essential to keep humans as
615 supervisors of AI decision-making (i.e., human-in-the-loop). In the context of systematic
616 evidence syntheses, validate AI decisions against established evidence synthesis
617 standards and guidelines for conduct and reporting (e.g., [1,72,73]).

618 Third, at the moment, AI is more reliable in some evidence synthesis steps (such as title
619 and abstract screening, and to some extent search strategy design and full-text
620 screening) than others (such as data extraction and critical appraisal). To prevent
621 relevant omissions for search strategy and screening supported by AI, there is a need
622 for detailed scoping exercise that will test all phases of the review before it is conducted.

623 Finally, we urge AI developers to provide decision files that facilitate the scrutiny of AI
624 algorithms, because transparency is crucial (e.g., see ASReview AI software, [63]), and
625 we should make decision data files accessible [12]. The evidence synthesis community
626 urgently needs a guide for reporting of AI-supported reviews (e.g., PRISMA extension
627 PRISMA-DfLLM for LLM; [74]). Such transparency will help with trust building between
628 evidence producers and evidence users.

Recommendations



629
630
631
632
633
634

Figure I: Recommendations for responsible AI use for evidence synthesis in conservation.