

BMJ Open

BMJ Open is committed to open peer review. As part of this commitment we make the peer review history of every article we publish publicly available.

When an article is published we post the peer reviewers' comments and the authors' responses online. We also post the versions of the paper that were used during peer review. These are the versions that the peer review comments apply to.

The versions of the paper that follow are the versions that were submitted during the peer review process. They are not the versions of record or the final published versions. They should not be cited or distributed as the published version of this manuscript.

BMJ Open is an open access journal and the full, final, typeset and author-corrected version of record of the manuscript is available on our site with no access controls, subscription charges or pay-per-view fees (<http://bmjopen.bmj.com>).

If you have any questions on BMJ Open's open peer review process please email info.bmjopen@bmj.com

BMJ Open

A modern era systematic review using artificial intelligence: considerations to ensure methodological quality

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2023-072254.R1
Article Type:	Communication
Date Submitted by the Author:	22-Feb-2023
Complete List of Authors:	van Dijk, Sanne H B; University of Twente Technical Medical Centre, Health Technology & Services Research; Medisch Spectrum Twente, Pulmonary Medicine Brusse-Keizer, Marjolein; Medisch Spectrum Twente, Medical School Twente; University of Twente Technical Medical Centre, Health Technology & Services Research Bucsán, Charlotte; Medisch Spectrum Twente, Pulmonary Medicine; University of Twente Faculty of Behavioural Sciences, Cognition, Data & Education van der Palen, Job; Medisch Spectrum Twente, Medical School Twente; University of Twente Faculty of Behavioural Sciences, Cognition, Data & Education Doggen, Carine J.M.; University of Twente Technical Medical Centre, Health Technology & Services Research; Rijnstate Hospital, Clinical Research Centre Lenferink, Anke; University of Twente Technical Medical Centre, Health Technology & Services Research; Medisch Spectrum Twente, Pulmonary Medicine
Primary Subject Heading:	Communication
Secondary Subject Heading:	Research methods
Keywords:	Systematic Review, STATISTICS & RESEARCH METHODS, Information technology < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in BMJ Open and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

TITLE PAGE***Title:***

A modern era systematic review using artificial intelligence: considerations to ensure methodological quality

Authors and affiliations:

Sanne H B van Dijk^{1,2}, Marjolein G J Brusse-Keizer^{1,3}, Charlotte C Bucsán^{2,4}, Job van der Palen^{3,4}, Carine J M Doggen^{1,5}, Anke Lenferink^{1,2,5}

¹ Health Technology & Services Research, Technical Medical Centre, University of Twente, Enschede, the Netherlands

² Department of Pulmonary Medicine, Medisch Spectrum Twente, Enschede, the Netherlands

³ Medical School Twente, Medisch Spectrum Twente, Enschede, the Netherlands

⁴ Cognition, Data & Education, Faculty of Behavioural, Management & Social Sciences, University of Twente, Enschede, the Netherlands

⁵ Clinical Research Centre, Rijnstate Hospital, Arnhem, the Netherlands

Corresponding author:

Anke Lenferink, a.lenferink@utwente.nl

ABSTRACT

Systematic reviews provide a structured overview of the available evidence in medical-scientific research. However, due to the increasing medical-scientific research output, it is a time-consuming task to conduct systematic reviews. To accelerate this process, artificial intelligence (AI) can be used in the screening of titles and abstracts. In this communication paper, we suggest how to conduct a transparent and reliable systematic review using the AI tool 'ASReview' in the title and abstract screening. Using the tool in our review resulted in much time saved: 23% of the articles were screened with the AI tool, leaving 77% of the articles unseen. Considerations to ensure methodological quality when using AI in systematic reviews included: the choice of when to use AI, the need of both deduplication and inter-reviewer agreement, how to choose a stopping criterion, and the quality of reporting. The AI tool is an important innovation for current systematic reviewing practice, as long as it is appropriately used and methodological quality can be assured.

STRENGTHS AND LIMITATIONS OF THIS STUDY

- Potential pitfalls regarding the use of artificial intelligence in systematic reviewing were identified.
- Remedies for each pitfall were provided to ensure methodological quality.
- A time-efficient approach is suggested on how to conduct a transparent and reliable systematic review using an artificial intelligence tool.
- The artificial intelligence tool described in the paper was not evaluated for its accuracy.

BACKGROUND

Medical-scientific research output has grown exponentially since the very first medical papers were published (1–3). The output in the field of clinical medicine increased and keeps doing so (4). To illustrate, a quick PubMed search for “cardiology” shows a fivefold increase in annual publications from 10,420 (2007) to 52,537 (2021). Although the medical-scientific output growth rate is not higher when compared to other scientific fields (1–3), this field creates the largest output (3). Staying updated by reading all published articles is therefore not feasible. However, systematic reviews facilitate up-to-date and accessible summaries of evidence, as they synthesise previously published results in a transparent and reproducible manner (5,6). Hence, conclusions can be drawn that provide the highest considered level of evidence in medical research (5,7). Therefore, systematic reviews are not only crucial in science, but they have a large impact on clinical practice and policy-making as well (6). They are, however, highly labour-intensive to conduct due to the necessity of screening a large amount of research. Thus, efficient and innovative reviewing methods are desired (8).

An open-source artificial intelligence (AI) tool ‘ASReview’ (9) was published in 2021 to facilitate the title and abstract screening process in systematic reviews. Applying this tool facilitates researchers to conduct systematic reviews: simulations already showed its time-saving potential (9–11). We used the tool in the study selection of our own systematic review and came across scenarios that needed consideration to prevent loss of methodological quality. In this communication paper, we provide a reliable and transparent AI-supported systematic reviewing approach.

METHODS

We first describe how the AI tool was used in a systematic review conducted by our research group. For more detailed information regarding searches and eligibility criteria of the review, we refer to the protocol (PROSPERO registry: CRD42022283952) (12). Subsequently, when deciding on the AI screening-related methodology, we applied appropriate remedies against foreseen scenarios and their pitfalls to maintain a reliable and transparent approach. These potential scenarios, pitfalls and remedies will be discussed in the result section.

In our systematic review, the AI tool ‘ASReview’ (version 0.17.1) (9) was used for the screening of titles and abstracts by the first reviewer (SvD). The tool uses an active researcher-in-the-loop machine learning algorithm to rank the articles from high to low probability of eligibility for inclusion by text mining. The AI tool offers several classifier models by which the relevancy of the included articles can be determined (9). In a simulation study using six large systematic review datasets on various topics, a Naïve Bayes (NB) and a term frequency-inverse document frequency (TF-IDF) outperformed other model settings (10). The NB classifier estimates the probability of an article being relevant, based on TF-IDF measurements. TF-IDF measures the originality of a certain word within the article relative to the total number of articles the word appears in (13). This combination of NB and TF-IDF were chosen for our systematic review.

Before the AI tool can be used for the screening of relevant articles, its algorithm needs training with at least one relevant and one irrelevant article (i.e., prior knowledge). It is assumed that the more prior knowledge, the better the algorithm is trained at the start of the screening process, and the faster it will identify relevant articles (9). In our review, the prior knowledge consisted of three relevant articles (14–16) selected from a systematic review on the topic (17) and three randomly picked irrelevant articles .

1
2
3 After training with the prior knowledge, the AI tool made a first ranking of all
4 unlabelled articles (i.e., articles not yet decided on eligibility). Unlabelled articles were ranked
5 from highest to lowest probability of being relevant by the classifier, based on the prior
6 knowledge (i.e., articles already decided on). The first reviewer read the title and abstract of
7 the highest ranked article and made a decision ('relevant' or 'irrelevant') following the
8 eligibility criteria. This decision was added to the prior knowledge and the AI tool made a new
9 ranking. Again, the next top ranked article was proposed to the reviewer, who made a decision
10 regarding eligibility. The process of AI taking additional prior knowledge into account for each
11 ranking and a reviewer making decisions was repeated until the predefined stopping criterion
12 of – in our case - 100 subsequent irrelevant articles was reached.
13
14
15
16
17
18
19
20
21
22
23
24
25

26 The articles that were labelled relevant during the title and abstract screening were each
27 screened on full text independently by two reviewers (SvD & MBK, AL, JvdP, CD, CB) to
28 minimise the influence of subjectivity on inclusion. Disagreements regarding inclusion were
29 solved by a third independent reviewer.
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

RESULTS

How to maintain reliability and transparency when using AI in systematic reviewing

A summary of the potential scenarios, and their pitfalls and remedies, when using the AI tool in a systematic review is given in Table 1. These potential scenarios should not be ignored, but acted upon to maintain reliability and transparency. Figure 1 shows when and where to act upon during the screening process, from literature search results to publishing the review.

Table 1 Per-scenario overview of potential pitfalls and how to prevent these when using ASReview in a systematic review

Potential scenario		Pitfall	Remedy
①	Only a small (i.e., manually feasible) number of articles (with possibly a high proportion relevant) available for screening	Time wasted by considering AI-related choices, and no time saved by using AI	Do not use AI: conduct manual screening
②	Presence of duplicate articles in ASReview	Unequal weighing of labelled articles in AI-supported screening	Apply deduplication methods before using AI
③	Reviewer's own opinion, expertise or mistakes influence(s) AI algorithm on article selection	Not all relevant articles are included, potentially introducing selection bias	Reviewer training in title and abstract screening Perform (partial) double screening and check inter-reviewer agreement
④	AI-supported screening is stopped before or a long time after all relevant articles are found	Not all relevant articles are included, potentially introducing selection bias, or time is wasted	Formulate a data-driven stopping criterion (i.e., number of consecutive irrelevant articles)
⑤	AI-related choices not (completely) described	Irreproducible results, leading to a low-quality systematic review	Describe and substantiate the choices that are made
⑥	Study selection is not transparent	Irreproducible results (black box algorithm), leading to a low-quality systematic review	Publish open data (i.e., extracted file with all decisions)

In our systematic review, by means of broad literature searches in several scientific databases, a first set of potentially relevant articles was identified. This yielded 8,456 articles, enough to apply the AI tool in the title and abstract screening (scenario ①) was avoided, see

1
2
3 Table 1). Subsequently, this complete set of articles was uploaded in reference manager
4 EndNote X9 (18) and Covidence (19), where 3,761 duplicate articles were removed.
5
6 Deduplication is usually applied in systematic reviewing (20), but is increasingly important
7
8 prior to the use of AI. Since multiple decisions regarding a duplicate article weigh more than
9
10 one, this will disproportionately influence classification and possibly the results (Table 1,
11
12 scenario ②). In our review, a deduplicated set of articles was uploaded in the AI tool. Prior to
13
14 the actual AI-supported title and abstract screening, the reviewers (SvD & AL, MBK) trained
15
16 themselves with a small selection of 74 articles. The first reviewer became familiar with the
17
18 ASReview software, and all three reviewers learnt how to apply the eligibility criteria, to
19
20 minimise personal influence on the article selection (Table 1, scenario ③).
21
22
23
24
25

26
27 Defining the stopping criterion used in the screening process is left to the reviewer (9).
28
29 An optimal stopping criterion in active learning is considered a perfectly balanced trade-off
30
31 between a certain cost (in terms of time spent) of screening one more article versus the
32
33 predictive performance (in terms of identifying a new relevant article) that could be increased
34
35 by adding one more decision (21). The optimal stopping criterion in systematic reviewing
36
37 would be the moment that screening additional articles will not result in more relevant articles
38
39 being identified (22). Therefore, in our review, we predetermined a data-driven stopping
40
41 criterion for the title and abstract screening as ‘100 consecutive irrelevant articles’ in order to
42
43 prevent the screening from being stopped before or a long time after all relevant articles were
44
45 identified (Table 1, scenario ④).
46
47
48
49

50
51 Due to the fact that the stopping criterion was reached after 1,063 of the 4,695 articles,
52
53 only a part of the total number of articles was seen. Therefore, this approach might be sensitive
54
55 to possible mistakes when articles are screened by only one reviewer, influencing the
56
57 algorithm, possibly resulting in an incomplete selection of articles (Table 1, scenario ③) (23).
58
59 As a remedy, second reviewers (AL, MBK) checked 20% of the titles and abstracts seen by the
60

1
2
3 first reviewer. This 20% had a comparable ratio regarding relevant versus irrelevant articles
4 over all articles seen. The percentual agreement and Cohen's Kappa (), a measure for the
5 inter-reviewer agreement above chance, were calculated to express the reliability of the
6 decisions taken (24). The decisions were agreed in 96% and was 0.83. A equal of at least
7 0.6 is generally considered high (24), and thus it was assumed that the algorithm was reliably
8 trained by the first reviewer.
9

10
11
12 The reporting of the use of the AI tool should be transparent. If the choices made
13 regarding the use of the AI tool are not entirely reported (Table 1, scenario ⑤), the reader will
14 not be able to properly assess the methodology of the review, and review results may even be
15 graded as low-quality due to the lack of transparent reporting. The ASReview tool offers the
16 possibility to extract a data file providing insight into all decisions made during the screening
17 process, in contrast to various other "black box" AI-reviewing tools (9). This file will be
18 published alongside our systematic review to provide full transparency of our AI-supported
19 screening. This way, the screening with AI is reproducible (remedy to scenario ⑥, Table 1).
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37

38 ***Results of AI-supported study selection in a systematic review***

39
40 We experienced an efficient process of title and abstract screening in our systematic review.
41
42 Whereas the screening was performed with a database of 4,695 articles, the stopping criterion
43 was reached after 1,063 articles, so 23% were seen. Figure 2 shows the proportion of articles
44 identified as being relevant at any point during the AI-supported screening process. It can be
45 observed that the articles are indeed prioritised by the active learning algorithm: in the
46 beginning, relatively many relevant articles were found, but this decreased as the stopping
47 criterion (vertical red line, Figure 2) was approached. During the screening, 142 articles were
48 labelled relevant. After the inter-reviewer agreement check, 142 articles proceeded to the full
49 text reviewing phase, of which 65 were excluded because these were no articles with an original
50
51
52
53
54
55
56
57
58
59
60

1
2
3 research format, and three because the full text could not be retrieved. After full text reviewing
4
5 of the remaining 74 articles, 18 articles from 13 individual studies were included in our review.
6

7
8 After snowballing, one additional article from a study already included was added.
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

DISCUSSION

In our systematic review, the AI tool considerably reduced the number of articles in the screening process. Since the AI tool is offered open source, many researchers may benefit from its time-saving potential in selecting articles. Choices in several scenarios regarding the use of AI, however, are still left open to the researcher, and need consideration to prevent pitfalls. These include the choice whether or not to use AI, the importance of deduplication, double screening to check inter-reviewer agreement, a data-driven stopping criterion to optimally utilise the algorithm's predictive performance, and quality of reporting of the AI-related methodology chosen. This communication paper is, to our knowledge, the first elaborately explaining and discussing these choices regarding the application of this AI tool in a systematic review.

The main advantage of using the AI tool is the amount of time saved. Indeed, in our study, only 23% of the total number of articles were screened before the predefined stopping criterion was met. Assuming that all relevant articles were found, the AI tool saved 77% of the time for title and abstract screening. An additional advantage is that research questions previously unanswerable due to the insurmountable number of articles to screen in a 'classic' (i.e., manual) review, now actually are possible to answer. An example of the latter is a review screening over 60,000 articles (25), which would probably never have been performed without AI supporting the article selection.

Since the introduction of the ASReview tool in 2021 it was applied in seven published reviews (25–31). An important note to make is that only one (25) clearly reported AI-related choices in the methods and a complete and transparent flowchart reflecting the study selection process in the results section. Two reviews reported a relatively small number (< 400) of articles to screen (26,27), of which more than 75% of the articles were screened before the stopping criterion was met, so the amount of time saved was limited. Also, three reviews

1
2
3 reported many initial articles (> 6,000) (25,28,29) and one reported 892 articles (31), of which
4
5 only 5 to 10% needed to be screened. So in these reviews, the AI tool saved an impressive
6
7 amount of screening time. In our systematic review, 3% of the articles were labelled relevant
8
9 during the title and abstract screening and eventually, less than 1% of all initial articles were
10
11 included. These percentages are low, and are in line with the three above-mentioned reviews
12
13 (1-2% and 0-1%, respectively) (25,28,29). Still, relevancy and inclusion rates are much lower
14
15 when compared with 'classic' systematic reviews. A study evaluating the screening process in
16
17 25 'classic' systematic reviews showed that approximately 18% was labelled relevant and 5%
18
19 was actually included in the reviews (32). This difference is probably due to more narrow
20
21 literature searches in 'classic' reviews for feasibility purposes compared with AI-supported
22
23 reviews, resulting in a higher proportion of included articles.
24
25
26
27

28
29 In this paper we show how we applied the AI tool, but we did not evaluate it in terms
30
31 of accuracy. This means that we have to deal with a certain degree of uncertainty. Despite the
32
33 data-driven stopping criterion there is a chance that relevant articles were missed, as 77% was
34
35 automatically excluded. Considering this might have been the case, firstly, this could be due to
36
37 wrong decisions of the reviewer that would have undesirably influenced the training of the
38
39 algorithm by which the articles were labelled as (ir)relevant and the order in which they were
40
41 presented to the reviewer. Relevant articles could have therefore remained unseen if the
42
43 stopping criterion was reached before they were presented to the reviewer. As a remedy, of the
44
45 20% of the articles that screened by the first reviewer, relevancy was also assessed by another
46
47 reviewer to assess inter-reviewer reliability, which was high. It should be noted, though, that
48
49 'classic' title and abstract screening is not necessarily better than using AI, as medical-scientific
50
51 researchers tend to assess one out of nine abstracts wrongly (32). Secondly, the AI tool may
52
53 not have properly ranked highly relevant to irrelevant articles. However, given that simulations
54
55 proved this AI tool's accuracy before (9–11) this was not considered plausible. Since our study
56
57
58
59
60

1
2
3 applied, but did not evaluate, the AI tool, we encourage future studies evaluating the
4 performance of the tool. This could not only enrich the knowledge about the AI tool, but also
5 increase certainty about using it.
6
7
8
9

10 Although various researcher-in-the-loop AI tools for title and abstract screening have
11 been developed over the years (9,23,33), they often do not develop into usable mature software
12 (33), which impedes AI to be permanently implemented in research practice. For medical-
13 scientific research practice, it would therefore be helpful if large systematic review institutions,
14 like Cochrane and PRISMA, would consider to ‘officially’ make AI part of systematic
15 reviewing practice. When guidelines on the use of AI in systematic reviews are made available
16 and widely recognised, AI-supported systematic reviews can be uniformly and transparently
17 reported. Only then we can really benefit from AI’s time-saving potential and reduce our
18 research time waste.
19
20
21
22
23
24
25
26
27
28
29
30
31
32

33 **CONCLUSION**

34 Our experience with the AI tool during the title and abstract screening was positive as it has
35 highly accelerated the literature selection process. However, users should consider applying
36 appropriate remedies to scenarios that may form a threat to the methodological quality of the
37 review. We provided an overview of these scenarios, their pitfalls and remedies. These
38 encourage reliable use and transparent reporting of AI in systematic reviewing. To ensure the
39 continuation of conducting systematic reviews in the future, and given their importance for
40 medical guidelines and practice, we consider this tool as an important addition in the review
41 process.
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

DECLARATIONS

Contributors

SvD proposed the methodology and conducted the study selection. MBK, CD and AL critically reflected on the methodology. MBK and AL contributed substantially to the study selection. CB, JvdP and CD contributed to the study selection. The manuscript was primarily prepared by SvD and critically revised by all authors. All authors read and approved the final manuscript.

Funding

The systematic review is conducted as part of the RE-SAMPLE project. RE-SAMPLE has received funding from the European Union's Horizon 2020 research and innovation programme (grant agreement no. 965315).

Competing interests

None declared.

Protocol registration

PROSPERO CRD42022283952

REFERENCES

1. Bornmann L, Mutz R. Growth Rates of Modern Science: A Bibliometric Analysis Based on the Number of Publications and Cited References. *J Am Soc Inf Sci Technol*. 2015;66(11):2215–22.
2. Bornmann L, Haunschild R, Mutz R. Growth rates of modern science: a latent piecewise growth curve approach to model publication numbers from established and new literature databases. *Humanit Soc Sci Commun*. 2021;8(1).
3. Michels C, Schmoch U. The growth of science and database coverage. *Scientometrics*. 2012;93(3):831–46.
4. Haghani M, Abbasi A, Zwack CC, Shahhoseini Z, Haslam N. Trends of research productivity across author gender and research fields: A multidisciplinary and multi-country observational study [Internet]. Vol. 17, *PLoS ONE*. 2022. Available from: <http://dx.doi.org/10.1371/journal.pone.0271998>
5. Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. The PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *J Clin Epidemiol* [Internet]. 2021;134:178–89. Available from: <https://doi.org/10.1016/j.jclinepi.2021.03.001>
6. Gurevitch J, Koricheva J, Nakagawa S, Stewart G. Meta-analysis and the science of research synthesis. *Nature*. 2018;555(7695):175–82.
7. Burns P, Rohrich R, Chung K. The Levels of Evidence and their role in Evidence-Based Medicine. *Plast Reconstr Surg* [Internet]. 2011;128(1):305–10. Available from: <http://www.alpbc.eu/>
8. Bastian H, Glasziou P, Chalmers I. Seventy-five trials and eleven systematic reviews a day: How will we ever keep up? *PLoS Med*. 2010;7(9).
9. van de Schoot R, de Bruin J, Schram R, Zahedi P, de Boer J, Weijdemans F, et al. An

- 1
2
3 open source machine learning framework for efficient and transparent systematic
4 reviews. *Nat Mach Intell* [Internet]. 2021;3(February):125–33. Available from:
5
6 <http://dx.doi.org/10.1038/s42256-020-00287-7>
7
8
9
10. Ferdinands G, Schram R, de Bruin J, Bagheri A, Oberski D, Tummers L, et al. Active
11 learning for screening prioritization in systematic reviews: A simulation study. 2020.
12
13
14
11. Ferdinands G. AI-Assisted Systematic Reviewing: Selecting Studies to Compare
15 Bayesian Versus Frequentist SEM for Small Sample Sizes. *Multivariate Behav Res*
16 [Internet]. 2020;56(1):153–4. Available from:
17
18 <https://doi.org/10.1080/00273171.2020.1853501>
19
20
21
22
23
12. van Dijk SHB, Brusse-Keizer MGJ, Bucsan CC, Doggen CJM, van der Palen J,
24 Lenferink A. Distinguishing acute heart failure from exacerbations of COPD: An AI-
25 supported systematic literature review. *PROSPERO* [Internet].
26
27 2022;CRD42022283952. Available from:
28
29 https://www.crd.york.ac.uk/prospero/display_record.php?ID=CRD42022283952
30
31
32
33
34
13. Havrlant L, Kreinovich V. A simple probabilistic explanation of term frequency-
35 inverse document frequency (tf-idf) heuristic (and variations motivated by this
36 explanation). *Int J Gen Syst* [Internet]. 2017;46(1):27–36. Available from:
37
38 <http://dx.doi.org/10.1080/03081079.2017.1291635>
39
40
41
42
43
44
14. Wang R, Cao Z, Li Y, Yu K. Utility of N-terminal pro B-type natriuretic peptide and
45 mean platelet volume in differentiating congestive heart failure from chronic
46 obstructive pulmonary disease. *Int J Cardiol* [Internet]. 2013;170(2):e28–9. Available
47
48 from: <http://dx.doi.org/10.1016/j.ijcard.2013.10.048>
49
50
51
52
53
15. Ouanes I, Jalloul F, Ayed S, Dachraoui F, Ouanes-Besbes L, Fekih Hassen M, et al. N-
54 terminal proB-type natriuretic peptide levels aid the diagnosis of left ventricular
55 dysfunction in patients with severe acute exacerbations of chronic obstructive
56
57
58
59
60

- 1
2
3 pulmonary disease and renal dysfunction. *Respirology*. 2012;17(4):660–6.
4
5
6 16. Andrijevic I, Milutinov S, Lozanov Crvenkovic Z, Matijasevic J, Andrijevic A,
7
8 Kovacevic T, et al. N-Terminal Prohormone of Brain Natriuretic Peptide (NT-
9
10 proBNP) as a Diagnostic Biomarker of Left Ventricular Systolic Dysfunction in
11
12 Patients with Acute Exacerbation of Chronic Obstructive Pulmonary Disease
13
14 (AECOPD). *Lung* [Internet]. 2018;196(5):583–90. Available from:
15
16 <http://dx.doi.org/10.1007/s00408-018-0137-3>
17
18
19 17. Hawkins NM, Khosla A, Virani SA, McMurray JJV, FitzGerald JM. B-type natriuretic
20
21 peptides in chronic obstructive pulmonary disease: A systematic review. *BMC Pulm*
22
23 *Med* [Internet]. 2017;17(1). Available from: [http://dx.doi.org/10.1186/s12890-016-](http://dx.doi.org/10.1186/s12890-016-0345-7)
24
25 [0345-7](http://dx.doi.org/10.1186/s12890-016-0345-7)
26
27
28 18. Clarivate Analytics. EndNote X9. 2018.
29
30 19. Covidence systematic review software [Internet]. Melbourne, Australia: Veritas Health
31
32 Innovation; Available from: www.covidence.org
33
34
35 20. McKeown S, Mir ZM. Considerations for conducting systematic reviews: evaluating
36
37 the performance of different methods for de-duplicating references. *Syst Rev*.
38
39 2021;10(1):4–11.
40
41
42 21. Ishibashi H, Hino H. Stopping Criterion for Active Learning Based on Error Stability.
43
44 2021;1:1–32. Available from: <http://arxiv.org/abs/2104.01836>
45
46
47 22. Wang W, Cai W, Zhang Y. Stability-Based Stopping Criterion for Active Learning.
48
49 *Proc - IEEE Int Conf Data Mining, ICDM*. 2014;2015-Janua(January):1019–24.
50
51
52 23. Blaizot A, Veettil SK, Saidoung P, Moreno-Garcia CF, Wiratunga N, Aceves-Martins
53
54 M, et al. Using artificial intelligence methods for systematic review in health sciences:
55
56 A systematic review. *Res Synth Methods*. 2022;13(3):353–62.
57
58
59 24. McHugh ML. Interrater reliability: the kappa statistic. *Biochem Med* [Internet].
60

- 2012;22(3):276–82. Available from:
<https://www.ncbi.nlm.nih.gov/pmc/articles/PMC3900052/>
25. Bernardes RC, Botina LL, Araújo R dos S, Guedes RNC, Martins GF, Lima MAP. Artificial Intelligence-Aided Meta-Analysis of Toxicological Assessment of Agrochemicals in Bees. *Front Ecol Evol.* 2022;10(May).
26. Silva GFS, Fagundes TP, Teixeira BC, Chiavegatto Filho ADP. Machine Learning for Hypertension Prediction: a Systematic Review. *Curr Hypertens Rep [Internet].* 2022;(3). Available from: <https://doi.org/10.1007/s11906-022-01212-6>
27. Miranda L, Paul R, Pütz B, Koutsouleris N, Müller-Myhsok B. Systematic Review of Functional MRI Applications for Psychiatric Disease Subtyping. *Front Psychiatry.* 2021;12(October).
28. Schouw HM, Huisman LA, Janssen YF, Slart RHJA, Borra RJH, Willemsen ATM, et al. Targeted optical fluorescence imaging: a meta-narrative review and future perspectives. *Eur J Nucl Med Mol Imaging [Internet].* 2021;48(13):4272–92. Available from: <https://doi.org/10.1007/s00259-021-05504-y>
29. Bakkum L, Schuengel C, Sterkenburg PS, Frielink N, Embregts PJCM, de Schipper JC, et al. People with intellectual disabilities living in care facilities engaging in virtual social contact: A systematic review of the feasibility and effects on well-being. *J Appl Res Intellect Disabil.* 2022;35(1):60–74.
30. Huang Y, Procházková M, Lu J, Riad A, Macek P. Family Related Variables' Influences on Adolescents' Health Based on Health Behaviour in School-Aged Children Database, an AI-Assisted Scoping Review, and Narrative Synthesis. *Front Psychol.* 2022;13(August).
31. Zhang W, Huang S, Lam L, Evans R, Zhu C. Cyberbullying definitions and measurements in children and adolescents: Summarizing 20 years of global efforts.

- 1
2
3 Front Public Heal. 2022;10.
4
5
6 32. Wang Z, Nayfeh T, Tetzlaff J, O'Blenis P, Murad MH. Error rates of human reviewers
7 during abstract screening in systematic reviews. PLoS One [Internet]. 2020;15(1):1–8.
8 Available from: <http://dx.doi.org/10.1371/journal.pone.0227742>
9
10
11
12 33. Marshall IJ, Wallace BC. Toward systematic review automation: A practical guide to
13 using machine learning tools in research synthesis. Syst Rev. 2019;8(1):1–10.
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

For peer review only

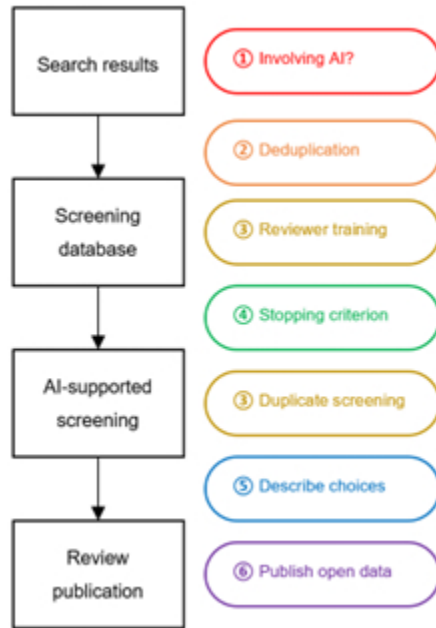
1
2
3 **FIGURE CAPTIONS**
4

5 **Figure 1** Flowchart showing when and where to act upon when using ASReview in systematic
6 reviewing
7
8
9

10
11 **Figure 2** Proportion of relevant articles identified after a certain number of titles and abstracts
12 were screened using the AI tool
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

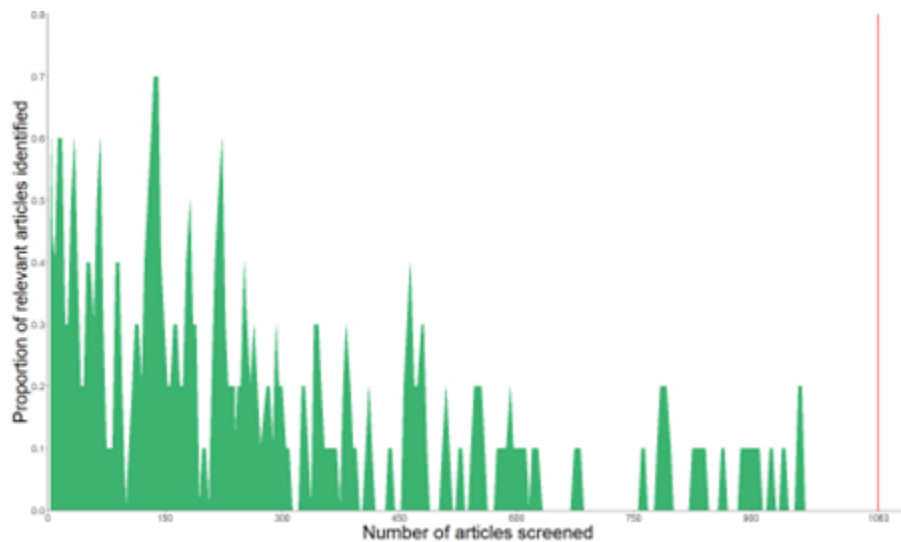
For peer review only

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60



Flowchart showing when and where to act upon when using ASReview in systematic reviewing

147x211mm (38 x 38 DPI)



Proportion of relevant articles identified after a certain number of titles and abstracts were screened using the AI tool

409x242mm (28 x 28 DPI)

BMJ Open

Artificial intelligence in systematic reviews: promising when appropriately used

Journal:	<i>BMJ Open</i>
Manuscript ID	bmjopen-2023-072254.R2
Article Type:	Communication
Date Submitted by the Author:	16-Jun-2023
Complete List of Authors:	van Dijk, Sanne H B; University of Twente Technical Medical Centre, Health Technology & Services Research; Medisch Spectrum Twente, Pulmonary Medicine Brusse-Keizer, Marjolein; Medisch Spectrum Twente, Medical School Twente; University of Twente Technical Medical Centre, Health Technology & Services Research Bucsán, Charlotte; Medisch Spectrum Twente, Pulmonary Medicine; University of Twente Faculty of Behavioural Sciences, Cognition, Data & Education van der Palen, Job; Medisch Spectrum Twente, Medical School Twente; University of Twente Faculty of Behavioural Sciences, Cognition, Data & Education Doggen, Carine J.M.; University of Twente Technical Medical Centre, Health Technology & Services Research; Rijnstate Hospital, Clinical Research Centre Lenferink, Anke; University of Twente Technical Medical Centre, Health Technology & Services Research; Medisch Spectrum Twente, Pulmonary Medicine
Primary Subject Heading:	Communication
Secondary Subject Heading:	Research methods
Keywords:	Systematic Review, STATISTICS & RESEARCH METHODS, Information technology < BIOTECHNOLOGY & BIOINFORMATICS

SCHOLARONE™
Manuscripts



I, the Submitting Author has the right to grant and does grant on behalf of all authors of the Work (as defined in the below author licence), an exclusive licence and/or a non-exclusive licence for contributions from authors who are: i) UK Crown employees; ii) where BMJ has agreed a CC-BY licence shall apply, and/or iii) in accordance with the terms applicable for US Federal Government officers or employees acting as part of their official duties; on a worldwide, perpetual, irrevocable, royalty-free basis to BMJ Publishing Group Ltd ("BMJ") its licensees and where the relevant Journal is co-owned by BMJ to the co-owners of the Journal, to publish the Work in BMJ Open and any other BMJ products and to exploit all rights, as set out in our [licence](#).

The Submitting Author accepts and understands that any supply made under these terms is made by BMJ to the Submitting Author unless you are acting as an employee on behalf of your employer or a postgraduate student of an affiliated institution which is paying any applicable article publishing charge ("APC") for Open Access articles. Where the Submitting Author wishes to make the Work available on an Open Access basis (and intends to pay the relevant APC), the terms of reuse of such Open Access shall be governed by a Creative Commons licence – details of these licences and which [Creative Commons](#) licence will apply to this Work are set out in our licence referred to above.

Other than as permitted in any relevant BMJ Author's Self Archiving Policies, I confirm this Work has not been accepted for publication elsewhere, is not being considered for publication elsewhere and does not duplicate material already published. I confirm all authors consent to publication of this Work and authorise the granting of this licence.

1
2
3
4 **1 TITLE PAGE**

5
6
7 **2 *Title:***

8
9 3 Artificial intelligence in systematic reviews: promising when appropriately used

10
11
12
13
14 **5 *Authors and affiliations:***

15
16 6 Sanne H B van Dijk^{1,2}, Marjolein G J Brusse-Keizer^{1,3}, Charlotte C Bucsán^{2,4}, Job van der
17
18 7 Palen^{3,4}, Carine J M Doggen^{1,5}, Anke Lenferink^{1,2,5}

19
20
21 8 ¹ Health Technology & Services Research, Technical Medical Centre, University of Twente,
22
23 9 Enschede, the Netherlands

24
25 10 ² Department of Pulmonary Medicine, Medisch Spectrum Twente, Enschede, the Netherlands

26
27 11 ³ Medical School Twente, Medisch Spectrum Twente, Enschede, the Netherlands

28
29
30 12 ⁴ Cognition, Data & Education, Faculty of Behavioural, Management & Social Sciences,
31
32 13 University of Twente, Enschede, the Netherlands

33
34 14 ⁵ Clinical Research Centre, Rijnstate Hospital, Arnhem, the Netherlands

35
36
37
38
39 **16 *Corresponding author:***

40
41 17 Anke Lenferink, a.lenferink@utwente.nl

42
43
44 18

1 2 3 4 5 6 7 8 9 10 11 12 13 14 15 16 17 18 19 20 21 22 23 24 25 26 27 28 29 30 31 32 33 34 35 36 37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60

1 ABSTRACT

2 **Background:** Systematic reviews provide a structured overview of the available evidence in
3 medical-scientific research. However, due to the increasing medical-scientific research output,
4 it is a time-consuming task to conduct systematic reviews. To accelerate this process, artificial
5 intelligence (AI) can be used in the review process. In this communication paper, we suggest
6 how to conduct a transparent and reliable systematic review using the AI tool ‘ASReview’ in
7 the title and abstract screening.

8 **Methods:** Use of the AI tool consisted of several steps. First, the tool required training of its
9 algorithm with several prelabelled articles prior to screening. Next, using a researcher-in-the-
10 loop algorithm, the AI tool proposed the article with the highest probability of being relevant.
11 The reviewer then decided on relevancy of each article proposed. This process was continued
12 until the stopping criterion was reached. All articles labelled relevant by the reviewer were
13 screened on full text.

14 **Results:** Considerations to ensure methodological quality when using AI in systematic reviews
15 included: the choice of whether to use AI, the need of both deduplication and inter-reviewer
16 agreement, how to choose a stopping criterion, and the quality of reporting. Using the tool in
17 our review resulted in much time saved: only 23% of the articles were proposed by the AI tool.

18 **Conclusion:** The AI tool is a promising innovation for the current systematic reviewing
19 practice, as long as it is appropriately used and methodological quality can be assured.

21 **STRENGTHS AND LIMITATIONS OF THIS STUDY**

22 - Potential pitfalls regarding the use of artificial intelligence in systematic reviewing were
23 identified.

- 1
- 2
- 3 1 - Remedies for each pitfall were provided to ensure methodological quality. - A time-efficient
- 4
- 5 2 approach is suggested on how to conduct a transparent and reliable
- 6
- 7
- 8 3 systematic review using an artificial intelligence tool.
- 9
- 10
- 11 4 - The artificial intelligence tool described in the paper was not evaluated for its accuracy.
- 12
- 13
- 14
- 15 5
- 16
- 17
- 18
- 19
- 20
- 21
- 22
- 23
- 24
- 25
- 26
- 27
- 28
- 29
- 30
- 31
- 32
- 33
- 34
- 35
- 36
- 37
- 38
- 39
- 40
- 41
- 42
- 43
- 44
- 45
- 46
- 47
- 48
- 49
- 50
- 51
- 52
- 53
- 54
- 55
- 56
- 57
- 58
- 59
- 60

For peer review only

1 **BACKGROUND**

2 Medical-scientific research output has grown exponentially since the very first medical papers
3 were published (1–3). The output in the field of clinical medicine increased and keeps doing
4 so (4). To illustrate, a quick PubMed search for “cardiology” shows a fivefold increase in
5 annual publications from 10,420 (2007) to 52,537 (2021). Although the medical-scientific
6 output growth rate is not higher when compared to other scientific fields (1–3), this field creates
7 the largest output (3). Staying updated by reading all published articles is therefore not feasible.
8 However, systematic reviews facilitate up-to-date and accessible summaries of evidence, as
9 they synthesise previously published results in a transparent and reproducible manner (5,6).
10 Hence, conclusions can be drawn that provide the highest considered level of evidence in
11 medical research (5,7). Therefore, systematic reviews are not only crucial in science, but they
12 have a large impact on clinical practice and policy-making as well (6). They are, however,
13 highly labour-intensive to conduct due to the necessity of screening a large amount of articles,
14 which results in a high consumption of research resources. Thus, efficient and innovative
15 reviewing methods are desired (8).

16 An open-source artificial intelligence (AI) tool ‘ASReview’ (9) was published in 2021
17 to facilitate the title and abstract screening process in systematic reviews. Applying this tool
18 facilitates researchers to conduct systematic reviews: simulations already showed its time-
19 saving potential (9–11). We used the tool in the study selection of our own systematic review
20 and came across scenarios that needed consideration to prevent loss of methodological quality.
21 In this communication paper, we provide a reliable and transparent AI-supported systematic
22 reviewing approach.

23

1 METHODS

2 We first describe how the AI tool was used in a systematic review conducted by our research
3 group. For more detailed information regarding searches and eligibility criteria of the review,
4 we refer to the protocol (PROSPERO registry: CRD42022283952) (12). Subsequently, when
5 deciding on the AI screening-related methodology, we applied appropriate remedies against
6 foreseen scenarios and their pitfalls to maintain a reliable and transparent approach. These
7 potential scenarios, pitfalls and remedies will be discussed in the result section.

8 In our systematic review, the AI tool ‘ASReview’ (version 0.17.1) (9) was used for the
9 screening of titles and abstracts by the first reviewer (SvD). The tool uses an active researcher-
10 in-the-loop machine learning algorithm to rank the articles from high to low probability of
11 eligibility for inclusion by text mining. The AI tool offers several classifier models by which
12 the relevancy of the included articles can be determined (9). In a simulation study using six
13 large systematic review datasets on various topics, a Naïve Bayes (NB) and a term frequency-
14 inverse document frequency (TF-IDF) outperformed other model settings (10). The NB
15 classifier estimates the probability of an article being relevant, based on TF-IDF measurements.
16 TF-IDF measures the originality of a certain word within the article relative to the total number
17 of articles the word appears in (13). This combination of NB and TF-IDF were chosen for our
18 systematic review.

19 Before the AI tool can be used for the screening of relevant articles, its algorithm needs
20 training with at least one relevant and one irrelevant article (i.e., prior knowledge). It is assumed
21 that the more prior knowledge, the better the algorithm is trained at the start of the screening
22 process, and the faster it will identify relevant articles (9). In our review, the prior knowledge
23 consisted of three relevant articles (14–16) selected from a systematic review on the topic (17)
24 and three randomly picked irrelevant articles .

1
2
3 1 After training with the prior knowledge, the AI tool made a first ranking of all
4
5 2 unlabelled articles (i.e., articles not yet decided on eligibility) from highest to lowest
6
7 3 probability of being relevant. The first reviewer read the title and abstract of the number one
8
9 4 ranked article and made a decision ('relevant' or 'irrelevant') following the eligibility criteria.
10
11 5 Next, the AI tool took into account this additional knowledge and made a new ranking. Again,
12
13 6 the next top ranked article was proposed to the reviewer, who made a decision regarding
14
15 7 eligibility. This process of AI making rankings and the reviewer making decisions, which is
16
17 8 also called 'researcher-in-the-loop', was repeated until the predefined data-driven stopping
18
19 9 criterion of – in our case - 100 subsequent irrelevant articles was reached. After the reviewer
20
21 10 rejected what the AI tool puts forward as 'most probably relevant' a hundred times, it was
22
23 11 assumed that there were no relevant articles left in the unseen part of the dataset.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

12 The articles that were labelled relevant during the title and abstract screening were each
13 screened on full text independently by two reviewers (SvD & MBK, AL, JvdP, CD, CB) to
14 minimise the influence of subjectivity on inclusion. Disagreements regarding inclusion were
15 solved by a third independent reviewer.
16

1 RESULTS

2 *How to maintain reliability and transparency when using AI in title and abstract screening*

3 A summary of the potential scenarios, and their pitfalls and remedies, when using the AI tool
4 in a systematic review is given in Table 1. These potential scenarios should not be ignored, but
5 acted upon to maintain reliability and transparency. Figure 1 shows when and where to act
6 upon during the screening process reflected by the PRISMA flowchart (18), from literature
7 search results to publishing the review.

8 In our systematic review, by means of broad literature searches in several scientific
9 databases, a first set of potentially relevant articles was identified, yielding 8,456 articles,
10 enough to expect the AI tool to be efficient in the title and abstract screening (scenario ①) was
11 avoided, see Table 1). Subsequently, this complete set of articles was uploaded in reference
12 manager EndNote X9 (19) and review manager Covidence (20), where 3,761 duplicate articles
13 were removed. Given that EndNote has quite low sensitivity in identifying duplicates,
14 additional deduplication in Covidence was considered beneficial (21). Deduplication is usually
15 applied in systematic reviewing (21), but is increasingly important prior to the use of AI. Since
16 multiple decisions regarding a duplicate article weigh more than one, this will
17 disproportionately influence classification and possibly the results (Table 1, scenario ②). In
18 our review, a deduplicated set of articles was uploaded in the AI tool. Prior to the actual AI-
19 supported title and abstract screening, the reviewers (SvD & AL, MBK) trained themselves
20 with a small selection of 74 articles. The first reviewer became familiar with the ASReview
21 software, and all three reviewers learnt how to apply the eligibility criteria, to minimise
22 personal influence on the article selection (Table 1, scenario ③).

23 Defining the stopping criterion used in the screening process is left to the reviewer (9).
24 An optimal stopping criterion in active learning is considered a perfectly balanced trade-off
25 between a certain cost (in terms of time spent) of screening one more article versus the

1
2
3 1 predictive performance (in terms of identifying a new relevant article) that could be increased
4
5 2 by adding one more decision (22). The optimal stopping criterion in systematic reviewing
6
7 3 would be the moment that screening additional articles will not result in more relevant articles
8
9
10 4 being identified (23). Therefore, in our review, we predetermined a data-driven stopping
11
12 5 criterion for the title and abstract screening as ‘100 consecutive irrelevant articles’ in order to
13
14 6 prevent the screening from being stopped before or a long time after all relevant articles were
15
16 7 identified (Table 1, scenario ④).

17
18
19 8 Due to the fact that the stopping criterion was reached after 1,063 of the 4,695 articles,
20
21 9 only a part of the total number of articles was seen. Therefore, this approach might be sensitive
22
23 10 to possible mistakes when articles are screened by only one reviewer, influencing the
24
25 11 algorithm, possibly resulting in an incomplete selection of articles (Table 1, scenario ③) (24).
26
27
28 12 As a remedy, second reviewers (AL, MBK) checked 20% of the titles and abstracts seen by the
29
30 13 first reviewer. This 20% had a comparable ratio regarding relevant versus irrelevant articles
31
32 14 over all articles seen. The percentual agreement and Cohen’s Kappa (κ), a measure for the
33
34 15 inter-reviewer agreement above chance, were calculated to express the reliability of the
35
36 16 decisions taken (25). The decisions were agreed in 96% and κ was 0.83. A κ equal of at least
37
38 17 0.6 is generally considered high (25), and thus it was assumed that the algorithm was reliably
39
40 18 trained by the first reviewer.

41
42
43
44 19 The reporting of the use of the AI tool should be transparent. If the choices made
45
46 20 regarding the use of the AI tool are not entirely reported (Table 1, scenario ⑤), the reader will
47
48 21 not be able to properly assess the methodology of the review, and review results may even be
49
50 22 graded as low-quality due to the lack of transparent reporting. The ASReview tool offers the
51
52 23 possibility to extract a data file providing insight into all decisions made during the screening
53
54 24 process, in contrast to various other “black box” AI-reviewing tools (9). This file will be
55
56
57
58
59
60

1 published alongside our systematic review to provide full transparency of our AI-supported
 2 screening. This way, the screening with AI is reproducible (remedy to scenario ⑥, Table 1).

3
 4 **Table 1** Per-scenario overview of potential pitfalls and how to prevent these when using ASReview in a
 5 systematic review

Potential scenario	Pitfall	Remedy
① Only a small (i.e., manually feasible*) number of articles (with possibly a high proportion relevant) available for screening	Time wasted by considering AI-related choices, software training, and no time saved by using AI	Do not use AI: conduct manual screening
② Presence of duplicate articles in ASReview	Unequal weighing of labelled articles in AI-supported screening	Apply deduplication methods before using AI
③ Reviewer's own opinion, expertise or mistakes influence(s) AI algorithm on article selection	Not all relevant articles are included, potentially introducing selection bias	Reviewer training in title and abstract screening Perform (partial) double screening and check inter-reviewer agreement
④ AI-supported screening is stopped before or a long time after all relevant articles are found	Not all relevant articles are included, potentially introducing selection bias, or time is wasted	Formulate a data-driven stopping criterion (i.e., number of consecutive irrelevant articles)
⑤ AI-related choices not (completely) described	Irreproducible results, leading to a low-quality systematic review	Describe and substantiate the choices that are made
⑥ Study selection is not transparent	Irreproducible results (black box algorithm), leading to a low-quality systematic review	Publish open data (i.e., extracted file with all decisions)

6 * What is considered manually feasible is highly context-dependent (i.e., the intended workload and/or number
 7 reviewers available)
 8

9 ***Results of AI-supported study selection in a systematic review***

10 We experienced an efficient process of title and abstract screening in our systematic review.
 11 Whereas the screening was performed with a database of 4,695 articles, the stopping criterion
 12 was reached after 1,063 articles, so 23% were seen. Figure 2A shows the proportion of articles
 13 identified as being relevant at any point during the AI-supported screening process. It can be

1
2
3 1 observed that the articles are indeed prioritised by the active learning algorithm: in the
4
5 2 beginning, relatively many relevant articles were found, but this decreased as the stopping
6
7 3 criterion (vertical red line) was approached. Figure 2B compares the screening progress when
8
9 4 using the AI tool versus manual screening. The moment the stopping criterion was reached,
10
11 5 approximately 32 records would have been found when the titles and abstract would have been
12
13 6 screened manually, compared to 142 articles labelled relevant using the AI tool. After the inter-
14
15 7 reviewer agreement check, 142 articles proceeded to the full text reviewing phase, of which 65
16
17 8 were excluded because these were no articles with an original research format, and three
18
19 9 because the full text could not be retrieved. After full text reviewing of the remaining 74
20
21 10 articles, 18 articles from 13 individual studies were included in our review. After snowballing,
22
23 11 one additional article from a study already included was added.
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 DISCUSSION

2 In our systematic review, the AI tool considerably reduced the number of articles in the
3 screening process. Since the AI tool is offered open source, many researchers may benefit from
4 its time-saving potential in selecting articles. Choices in several scenarios regarding the use of
5 AI, however, are still left open to the researcher, and need consideration to prevent pitfalls.
6 These include the choice whether or not to use AI by weighing the costs versus the benefits,
7 the importance of deduplication, double screening to check inter-reviewer agreement, a data-
8 driven stopping criterion to optimally utilise the algorithm's predictive performance, and
9 quality of reporting of the AI-related methodology chosen. This communication paper is, to
10 our knowledge, the first elaborately explaining and discussing these choices regarding the
11 application of this AI tool in an example systematic review.

12 The main advantage of using the AI tool is the amount of time saved. Indeed, in our
13 study, only 23% of the total number of articles were screened before the predefined stopping
14 criterion was met. Assuming that all relevant articles were found, the AI tool saved 77% of the
15 time for title and abstract screening. However, time should be invested to become acquainted
16 with the tool. Whether the expected screening time saved outweighs this time investment is
17 context-dependent (e.g., researcher's digital skills, systematic reviewing skills, topic
18 knowledge). An additional advantage is that research questions previously unanswerable due
19 to the insurmountable number of articles to screen in a 'classic' (i.e., manual) review, now
20 actually are possible to answer. An example of the latter is a review screening over 60,000
21 articles (26), which would probably never have been performed without AI supporting the
22 article selection.

23 Since the introduction of the ASReview tool in 2021, it was applied in seven published
24 reviews (26–32). An important note to make is that only one (26) clearly reported AI-related
25 choices in the methods and a complete and transparent flowchart reflecting the study selection

1
2
3 1 process in the results section. Two reviews reported a relatively small number (< 400) of
4
5 2 articles to screen (27,28), of which more than 75% of the articles were screened before the
6
7 3 stopping criterion was met, so the amount of time saved was limited. Also, three reviews
8
9 4 reported many initial articles (> 6,000) (26,29,30) and one reported 892 articles (32), of which
10
11 5 only 5 to 10% needed to be screened. So in these reviews, the AI tool saved an impressive
12
13 6 amount of screening time. In our systematic review, 3% of the articles were labelled relevant
14
15 7 during the title and abstract screening and eventually, less than 1% of all initial articles were
16
17 8 included. These percentages are low, and are in line with the three above-mentioned reviews
18
19 9 (1-2% and 0-1%, respectively) (26,29,30). Still, relevancy and inclusion rates are much lower
20
21 10 when compared with 'classic' systematic reviews. A study evaluating the screening process in
22
23 11 25 'classic' systematic reviews showed that approximately 18% was labelled relevant and 5%
24
25 12 was actually included in the reviews (33). This difference is probably due to more narrow
26
27 13 literature searches in 'classic' reviews for feasibility purposes compared with AI-supported
28
29 14 reviews, resulting in a higher proportion of included articles.
30
31
32
33
34

35 15 In this paper we show how we applied the AI tool, but we did not evaluate it in terms
36
37 16 of accuracy. This means that we have to deal with a certain degree of uncertainty. Despite the
38
39 17 data-driven stopping criterion there is a chance that relevant articles were missed, as 77% was
40
41 18 automatically excluded. Considering this might have been the case, firstly, this could be due to
42
43 19 wrong decisions of the reviewer that would have undesirably influenced the training of the
44
45 20 algorithm by which the articles were labelled as (ir)relevant and the order in which they were
46
47 21 presented to the reviewer. Relevant articles could have therefore remained unseen if the
48
49 22 stopping criterion was reached before they were presented to the reviewer. As a remedy, in our
50
51 23 own systematic review, of the 20% of the articles screened by the first reviewer, relevancy was
52
53 24 also assessed by another reviewer to assess inter-reviewer reliability, which was high. It should
54
55 25 be noted, though, that 'classic' title and abstract screening is not necessarily better than using
56
57
58
59
60

1 AI, as medical-scientific researchers tend to assess one out of nine abstracts wrongly (33).
2 Secondly, the AI tool may not have properly ranked highly relevant to irrelevant articles.
3 However, given that simulations proved this AI tool's accuracy before (9–11) this was not
4 considered plausible. Since our study applied, but did not evaluate, the AI tool, we encourage
5 future studies evaluating the performance of the tool across different scientific disciplines and
6 contexts, since research suggests that the tool's performance depends on the context, for
7 example, the complexity of the research question (34). This could not only enrich the
8 knowledge about the AI tool, but also increase certainty about using it. Also, future studies
9 should investigate the effects of choices made regarding the amount of prior knowledge that is
10 provided to the tool, the number of articles defining the stopping criterion, and how duplicate
11 screening is best performed, to guide future users of the tool.

12 Although various researcher-in-the-loop AI tools for title and abstract screening have
13 been developed over the years (9,24,35), they often do not develop into usable mature software
14 (35), which impedes AI to be permanently implemented in research practice. For medical-
15 scientific research practice, it would therefore be helpful if large systematic review institutions,
16 like Cochrane and PRISMA, would consider to 'officially' make AI part of systematic
17 reviewing practice. When guidelines on the use of AI in systematic reviews are made available
18 and widely recognised, AI-supported systematic reviews can be uniformly conducted and
19 transparently reported. Only then we can really benefit from AI's time-saving potential and
20 reduce our research time waste.

21

22 **CONCLUSION**

23 Our experience with the AI tool during the title and abstract screening was positive as it has
24 highly accelerated the literature selection process. However, users should consider applying
25 appropriate remedies to scenarios that may form a threat to the methodological quality of the

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 review. We provided an overview of these scenarios, their pitfalls and remedies. These
2 encourage reliable use and transparent reporting of AI in systematic reviewing. To ensure the
3 continuation of conducting systematic reviews in the future, and given their importance for
4 medical guidelines and practice, we consider this tool as an important addition in the review
5 process.

6

For peer review only

1
2
3 **1 DECLARATIONS**
4

5
6 **2 *Contributors***
7

8 3 SvD proposed the methodology and conducted the study selection. MBK, CD and AL critically
9
10 4 reflected on the methodology. MBK and AL contributed substantially to the study selection.
11
12 5 CB, JvdP and CD contributed to the study selection. The manuscript was primarily prepared
13
14 6 by SvD and critically revised by all authors. All authors read and approved the final manuscript.
15
16

17 **7 *Funding***
18

19 8 The systematic review is conducted as part of the RE-SAMPLE project. RE-SAMPLE has
20
21 9 received funding from the European Union's Horizon 2020 research and innovation
22
23 10 programme (grant agreement no. 965315).
24
25

26 **11 *Competing interests***
27

28
29 12 None declared.
30

31 **13 *Protocol registration***
32

33 14 PROSPERO CRD42022283952
34
35
36 15
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1
2
3 **1** **REFERENCES**
4

- 5
6 2 1. Bornmann L, Mutz R. Growth Rates of Modern Science: A Bibliometric Analysis
7
8 3 Based on the Number of Publications and Cited References. *J Am Soc Inf Sci Technol*.
9
10 4 2015;66(11):2215–22.
11
12 5 2. Bornmann L, Haunschild R, Mutz R. Growth rates of modern science: a latent
13
14 6 piecewise growth curve approach to model publication numbers from established and
15
16 7 new literature databases. *Humanit Soc Sci Commun*. 2021;8:224.
17
18 8 3. Michels C, Schmoch U. The growth of science and database coverage. *Scientometrics*.
19
20 9 2012;93(3):831–46.
21
22 10 4. Haghani M, Abbasi A, Zwack CC, Shahhoseini Z, Haslam N. Trends of research
23
24 11 productivity across author gender and research fields: A multidisciplinary and multi-
25
26 12 country observational study. *Vol. 17, PLoS ONE*. 2022. e0271998 p.
27
28 13 5. Page MJ, McKenzie JE, Bossuyt PM, Boutron I, Hoffmann TC, Mulrow CD, et al. The
29
30 14 PRISMA 2020 statement: An updated guideline for reporting systematic reviews. *J*
31
32 15 *Clin Epidemiol*. 2021;134:178–89.
33
34 16 6. Gurevitch J, Koricheva J, Nakagawa S, Stewart G. Meta-analysis and the science of
35
36 17 research synthesis. *Nature*. 2018;555(7695):175–82.
37
38 18 7. Burns P, Rohrich R, Chung K. The Levels of Evidence and their role in Evidence-
39
40 19 Based Medicine. *Plast Reconstr Surg*. 2011;128(1):305–10.
41
42 20 8. Bastian H, Glasziou P, Chalmers I. Seventy-five trials and eleven systematic reviews a
43
44 21 day: How will we ever keep up? *PLoS Med*. 2010;7(9):e1000326.
45
46 22 9. van de Schoot R, de Bruin J, Schram R, Zahedi P, de Boer J, Weijdemans F, et al. An
47
48 23 open source machine learning framework for efficient and transparent systematic
49
50 24 reviews. *Nat Mach Intell*. 2021;3(February):125–33.
51
52 25 10. Ferdinands G, Schram R, de Bruin J, Bagheri A, Oberski D, Tummers L, et al. Active
53
54
55
56
57
58
59
60

- 1
2
3 1 learning for screening prioritization in systematic reviews: A simulation study. 2020.
4
5 2 11. Ferdinands G. AI-Assisted Systematic Reviewing: Selecting Studies to Compare
6
7 3 Bayesian Versus Frequentist SEM for Small Sample Sizes. *Multivariate Behav Res.*
8
9 4 2020;56(1):153–4.
10
11 5 12. van Dijk SHB, Brusse-Keizer MGJ, Bucsán CC, Doggen CJM, van der Palen J,
12
13 6 Lenferink A. Distinguishing acute heart failure from exacerbations of COPD: An AI-
14
15 7 supported systematic literature review. *PROSPERO.* 2022;CRD42022283952.
16
17 8 13. Havrlant L, Kreinovich V. A simple probabilistic explanation of term frequency-
18
19 9 inverse document frequency (tf-idf) heuristic (and variations motivated by this
20
21 10 explanation). *Int J Gen Syst.* 2017;46(1):27–36.
22
23 11 14. Wang R, Cao Z, Li Y, Yu K. Utility of N-terminal pro B-type natriuretic peptide and
24
25 12 mean platelet volume in differentiating congestive heart failure from chronic
26
27 13 obstructive pulmonary disease. *Int J Cardiol.* 2013;170(2):e28–9.
28
29 14 15. Ouanes I, Jalloul F, Ayed S, Dachraoui F, Ouanes-Besbes L, Fekih Hassen M, et al. N-
30
31 15 terminal proB-type natriuretic peptide levels aid the diagnosis of left ventricular
32
33 16 dysfunction in patients with severe acute exacerbations of chronic obstructive
34
35 17 pulmonary disease and renal dysfunction. *Respirology.* 2012;17(4):660–6.
36
37 18 16. Andrijevic I, Milutinov S, Lozanov Crvenkovic Z, Matijasevic J, Andrijevic A,
38
39 19 Kovacevic T, et al. N-Terminal Prohormone of Brain Natriuretic Peptide (NT-
40
41 20 proBNP) as a Diagnostic Biomarker of Left Ventricular Systolic Dysfunction in
42
43 21 Patients with Acute Exacerbation of Chronic Obstructive Pulmonary Disease
44
45 22 (AECOPD). *Lung.* 2018;196(5):583–90.
46
47 23 17. Hawkins NM, Khosla A, Virani SA, McMurray JJV, FitzGerald JM. B-type natriuretic
48
49 24 peptides in chronic obstructive pulmonary disease: A systematic review. *BMC Pulm*
50
51 25 *Med.* 2017;17(1).
52
53
54
55
56
57
58
59
60

- 1
2
3 1 18. Haddaway NR, Page MJ, Pritchard CC, McGuinness LA. PRISMA2020: An R
4
5 2 package and Shiny app for producing PRISMA 2020-compliant flow diagrams, with
6
7 3 interactivity for optimised digital transparency and Open Synthesis. *Campbell Syst*
8
9 4 *Rev.* 2022;18(2):1–12.
10
11
12 5 19. Clarivate Analytics. EndNote X9. 2018.
13
14 6 20. Covidence systematic review software. Melbourne, Australia: Veritas Health
15
16 7 Innovation;
17
18 8 21. McKeown S, Mir ZM. Considerations for conducting systematic reviews: evaluating
19
20 9 the performance of different methods for de-duplicating references. *Syst Rev.*
21
22 10 2021;10(1):4–11.
23
24 11 22. Ishibashi H, Hino H. Stopping Criterion for Active Learning Based on Error Stability.
25
26 12 2021;1:1–32. Available from: <http://arxiv.org/abs/2104.01836>
27
28 13 23. Wang W, Cai W, Zhang Y. Stability-Based Stopping Criterion for Active Learning. In:
29
30 14 *Proceedings - IEEE International Conference on Data Mining, ICDM.* IEEE; 2014. p.
31
32 15 1019–24.
33
34 16 24. Blaizot A, Veetil SK, Saidoung P, Moreno-Garcia CF, Wiratunga N, Aceves-Martins
35
36 17 M, et al. Using artificial intelligence methods for systematic review in health sciences:
37
38 18 A systematic review. *Res Synth Methods.* 2022;13(3):353–62.
39
40 19 25. McHugh ML. Interrater reliability: the kappa statistic. *Biochem Med.* 2012;22(3):276–
41
42 20 82.
43
44 21 26. Bernardes RC, Botina LL, Araújo R dos S, Guedes RNC, Martins GF, Lima MAP.
45
46 22 Artificial Intelligence-Aided Meta-Analysis of Toxicological Assessment of
47
48 23 Agrochemicals in Bees. *Front Ecol Evol.* 2022;10:845608.
49
50 24 27. Silva GFS, Fagundes TP, Teixeira BC, Chiavegatto Filho ADP. Machine Learning for
51
52 25 Hypertension Prediction: a Systematic Review. *Curr Hypertens Rep.* 2022;(3).
53
54
55
56
57
58
59
60

- 1
2
3 1 28. Miranda L, Paul R, Pütz B, Koutsouleris N, Müller-Myhsok B. Systematic Review of
4
5 2 Functional MRI Applications for Psychiatric Disease Subtyping. *Front Psychiatry*.
6
7 3 2021;12:665536.
8
9
10 4 29. Schouw HM, Huisman LA, Janssen YF, Slart RHJA, Borra RJH, Willemsen ATM, et
11
12 5 al. Targeted optical fluorescence imaging: a meta-narrative review and future
13
14 6 perspectives. *Eur J Nucl Med Mol Imaging*. 2021;48:4272–92.
15
16
17 7 30. Bakkum L, Schuengel C, Sterkenburg PS, Frielink N, Embregts PJCM, de Schipper
18
19 8 JC, et al. People with intellectual disabilities living in care facilities engaging in virtual
20
21 9 social contact: A systematic review of the feasibility and effects on well-being. *J Appl*
22
23 10 *Res Intellect Disabil*. 2022;35(1):60–74.
24
25
26 11 31. Huang Y, Procházková M, Lu J, Riad A, Macek P. Family Related Variables’
27
28 12 Influences on Adolescents’ Health Based on Health Behaviour in School-Aged
29
30 13 Children Database, an AI-Assisted Scoping Review, and Narrative Synthesis. *Front*
31
32 14 *Psychol*. 2022;13:871795.
33
34
35 15 32. Zhang W, Huang S, Lam L, Evans R, Zhu C. Cyberbullying definitions and
36
37 16 measurements in children and adolescents: Summarizing 20 years of global efforts.
38
39 17 *Front Public Heal*. 2022;10:1000504.
40
41
42 18 33. Wang Z, Nayfeh T, Tetzlaff J, O’Blenis P, Murad MH. Error rates of human reviewers
43
44 19 during abstract screening in systematic reviews. *PLoS One*. 2020;15(1):e0227742.
45
46
47 20 34. Muthu S. The efficiency of machine learning-assisted platform for article screening in
48
49 21 systematic reviews in orthopaedics. *Int Orthop*. 2023;47:551–6.
50
51
52 22 35. Marshall IJ, Wallace BC. Toward systematic review automation: A practical guide to
53
54 23 using machine learning tools in research synthesis. *Syst Rev*. 2019;8(1):1–10.
55
56 24
57
58 25
59
60

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60

1 **FIGURE CAPTIONS**

2 **Figure 1** Flowchart showing when and where to act upon when using ASReview in systematic
3 reviewing

4 Footnote: Adapted the PRISMA flowchart from Haddaway et al. (18).

5
6 **Figure 2** Relevant articles identified after a certain number of titles and abstracts were screened
7 using the AI tool compared with manual screening

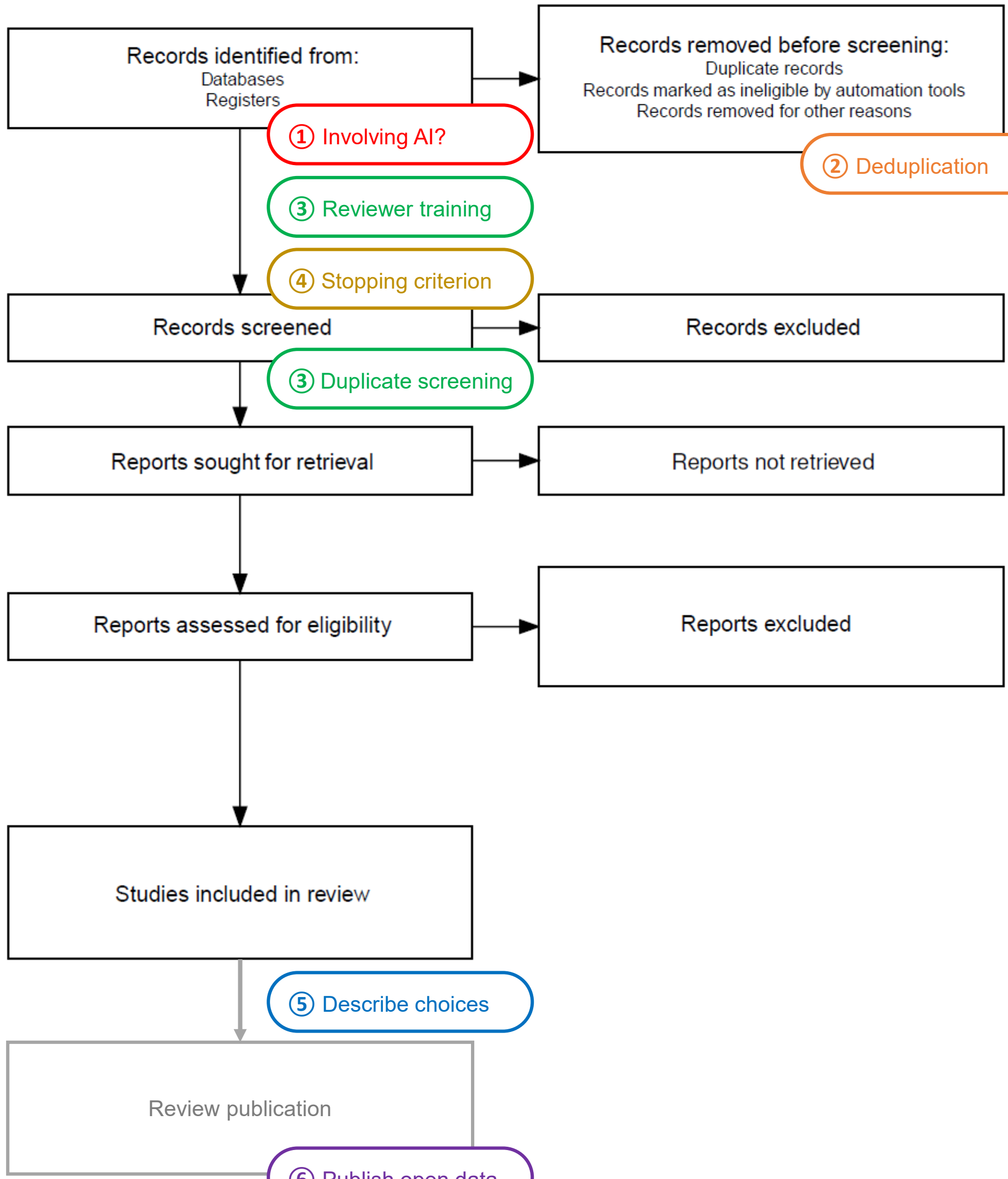
For peer review only

Identification of new studies via databases and registers

Identification

Screening

Included



1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46

