



Exploring responsible AI use in evidence synthesis: A framework for selecting AI tools

Ella Flemyng, Head of Editorial Policy & Research Integrity, Cochrane
Claudia Lenkewitz, Information Specialist, IQWiG

Trusted evidence.
Informed decisions.
Better health.





Declarations of interest: EF & CL are authors of Responsible AI use in evidence SynthEsis; EF convenor of joint AI Methods Group; no other interests to declare.

Source of support: EF & CL employers (Cochrane and IQWiG, respectively) have supported their involvement in this work; EF has been supported by funding via Destiny (Wellcome Trust [313586/Z/24/Z]).

Acknowledgement: Developed by Ella Flemyng, Claudia Lenkewitz, Sean Gardner, Jo-Ana Chase, James Thomas, Max Callaghan.



Session objectives

1. Highlight the four main expectations for evidence synthesists when using AI tools
2. Understand how the Responsible AI use in evidence SynthEsis (RAISE) guidance helps you navigate these four expectations
3. Introduce the responsible handover framework of an AI tool to evidence synthesists
4. Understand what details about an AI tool evidence synthesists should demand from tool developers



Introducing RAISE (Responsible AI use in evidence SynthEsis)





Recommendations and guidance

Three-paper RAISE collection

- 1** Responsible AI in Evidence synthesis 1: **Recommendations for practice**
- 2** Responsible AI in Evidence synthesis 2: **Building and evaluating evidence synthesis tools**
- 3** Responsible AI in Evidence synthesis 3: **Selecting and using evidence synthesis tools**



We need to support the wider adoption of AI

We need cross-field standards *and an evidence base*

We are part of an ecosystem made up of individuals, collaborations, and organisations

Each has a role to play in developing and using AI in a responsible way

One person / organisation may play multiple roles



AI use expectations for evidence synthesists



1

- Evidence synthesists are ultimately responsible for their research, including the decision to use AI and how it is used

2

- Evidence synthesists can use AI provided they can demonstrate it will not compromise the methodological rigor or integrity of their synthesis

3

- AI use should be fully and transparently reported

4

- AI should be used with human oversight

**Expectation 1: Evidence synthesists
are ultimately responsible for their
research, including the decision to use
AI and how it is used**

*Introducing the responsible handover framework for
an AI tool for evidence synthesists*

Responsible handover framework

Domain 1: What is the purpose of the AI tool?

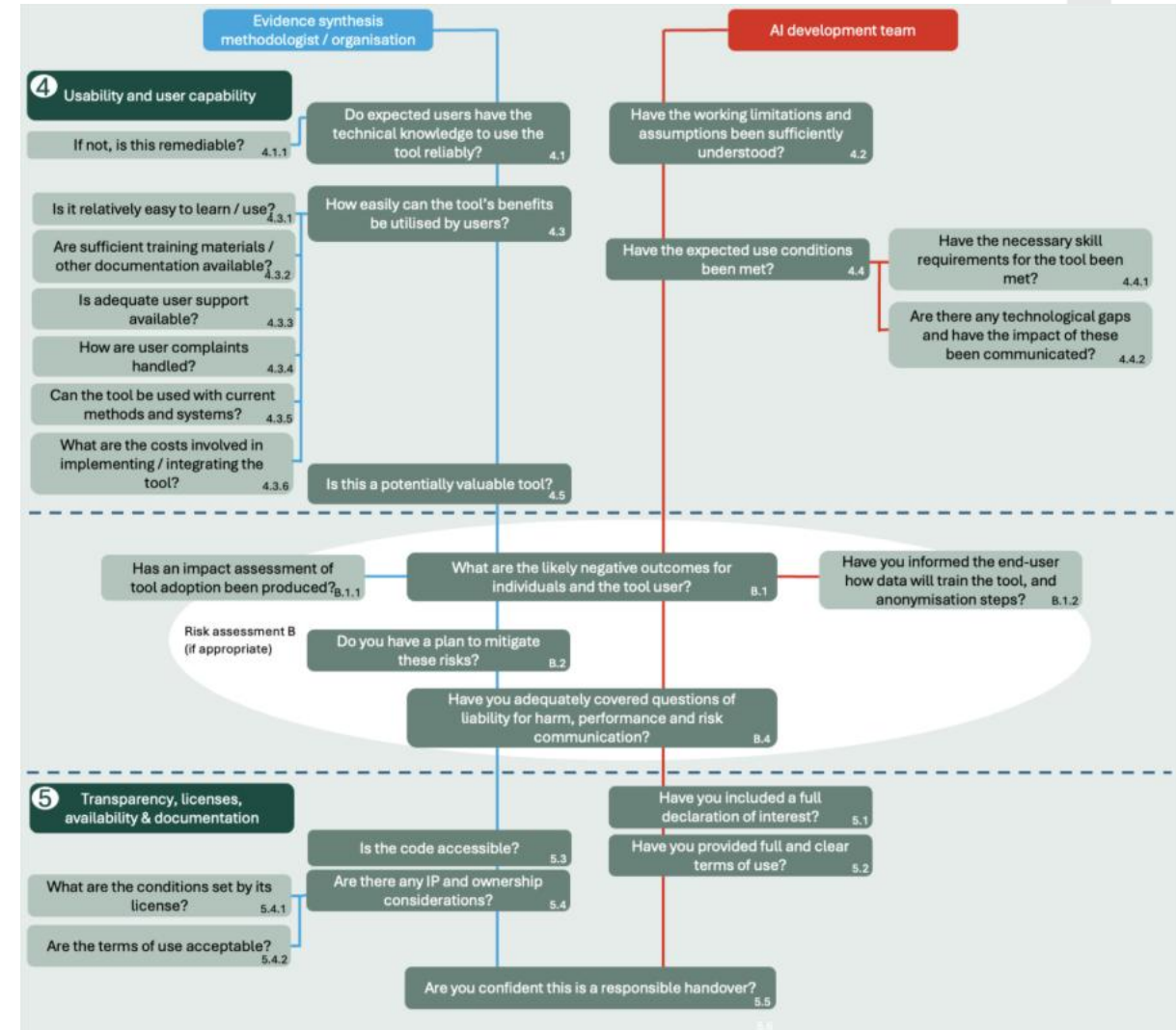
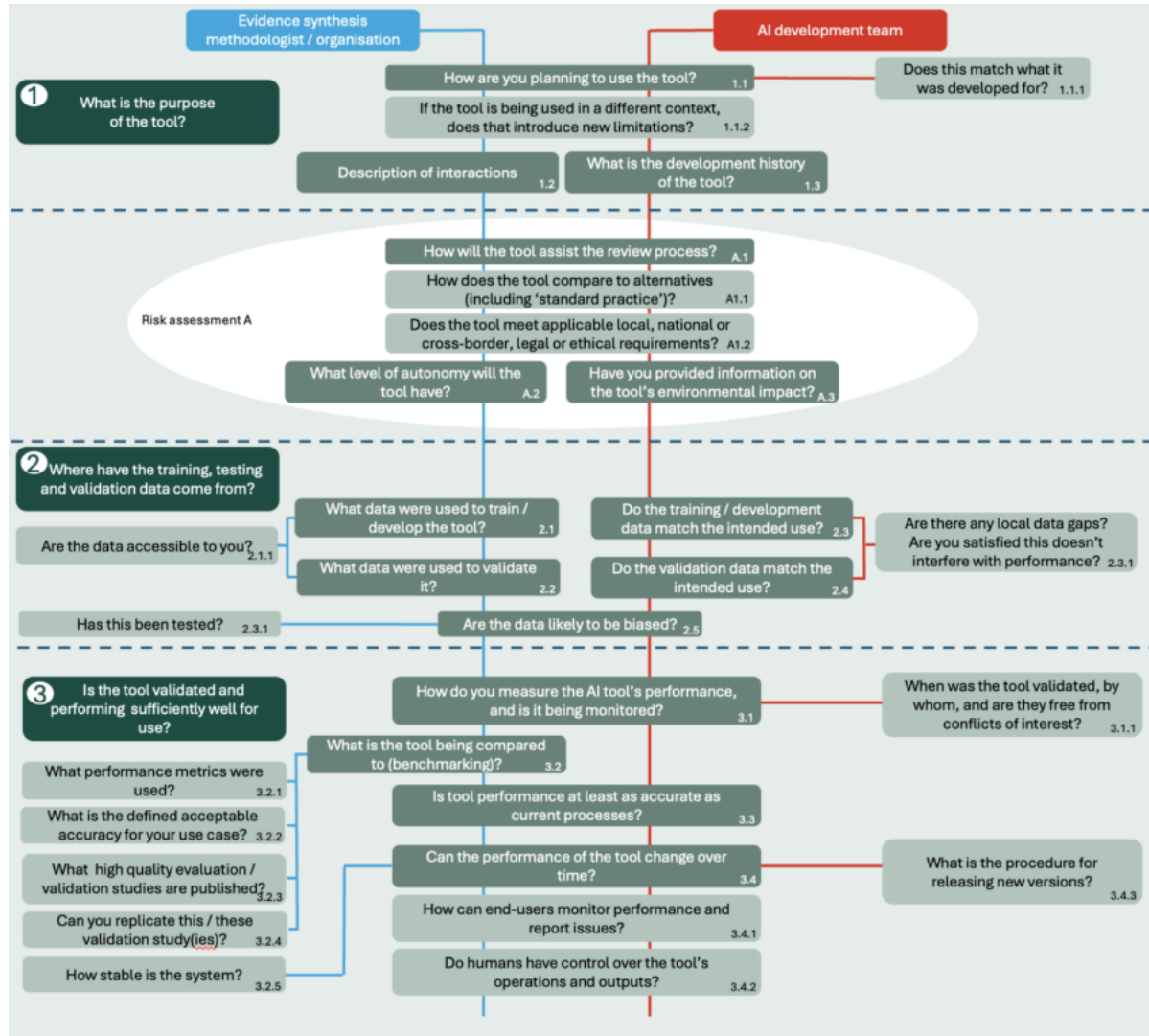
Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

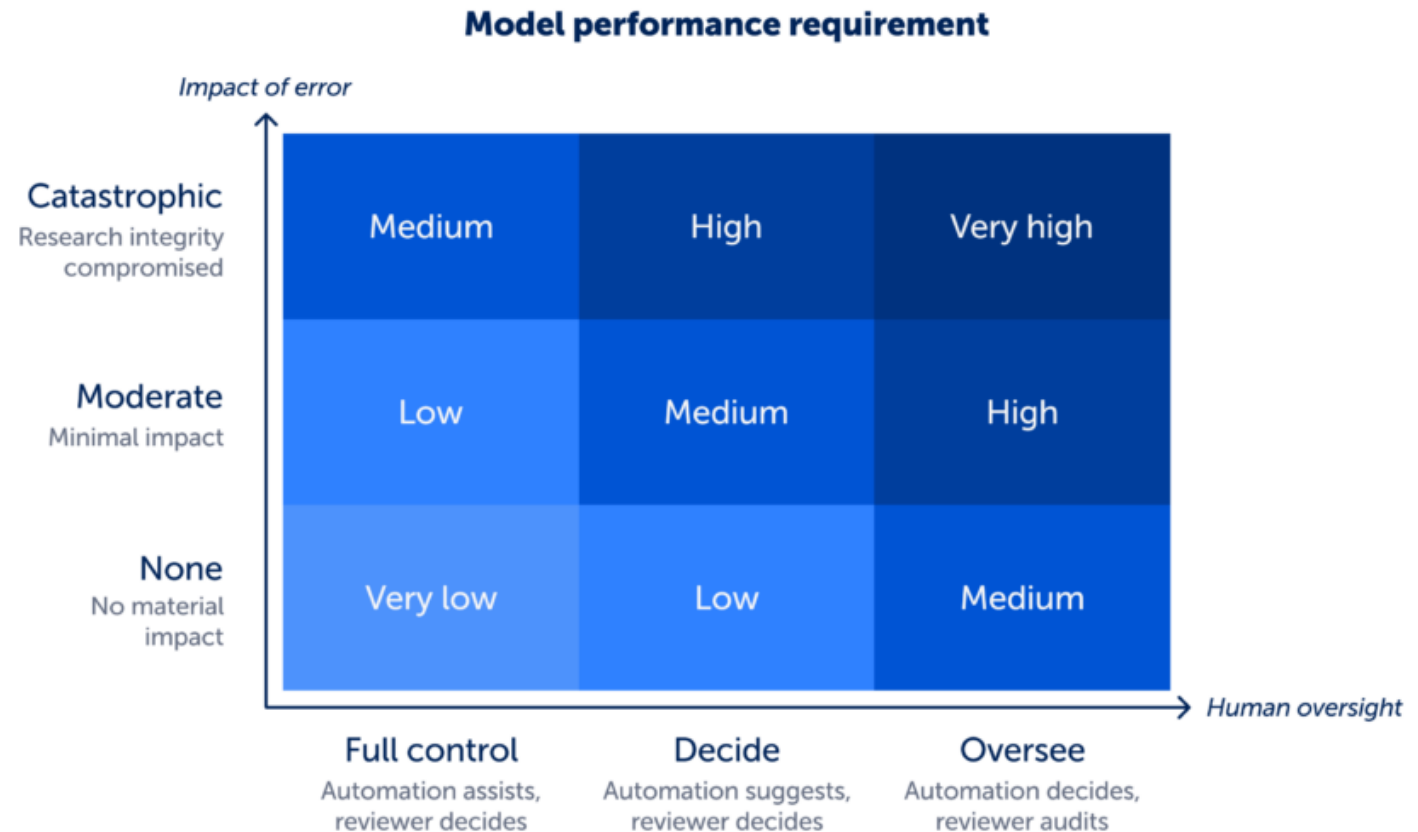
Domain 5: Transparency, licenses, availability and documentation





Considerations before using the framework

- Look critically; **all tools have limitations**
- Consider the acceptable model performance given the potential impact of errors and degree of human oversight
- Collaborate with those who have complementary backgrounds



Example evidence synthesis and task

- **Living systematic review**
- Focused on public health interventions addressing the **health impacts of climate change**
- Commissioned to **inform policy** at a national level
- Considerable time pressures

You will use AI for abstract and title screening

If the performance and use of the tool is acceptable then the AI could make the decisions (with disagreements flagged and ability to audit all decisions)

Or if not, you want the tool to assist or suggest with authors making the decisions

Responsible handover framework

Domain 1: What is the purpose of the AI tool?

Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

Domain 5: Transparency, licenses, availability and documentation



Domain 1: What is the purpose of the tool?

Main considerations for evidence synthesists:

- Is the intended use case and expected benefit clearly defined?
- Does it match your specific task and context?
- Is the level of human oversight appropriate?
- For LLMs, are the model version and prompts recorded?





Menti

Is it clear what the AI tool does and does it match our example review task?

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)



Public health interventions addressing the health impacts of climate change

Abstract and title screening

Depending on performance, the tool could make the decisions (disagreements flagged, auditability) or if not, could suggest

SuperScreener

Trained on evidence across diverse health fields, including climate-related health topics, and tested on 25 systematic reviews across health

Assist results

Set tool to make decisions

Audit and Trace

Decision Log

1015 PM

Cemihon Bont e ehri Vine deild bry
Táunðerhni

0:25 m

Decision berrtramps deplrenty
Táunðerhni

0:15 m

Audit and Trace

Decision Log

1013 PM

Cemihon bomos trodd Von deisd bry
Táunðerhni

0:25 m

Decision bomos vorh Vine deild bry
Táunðerhni

0:15 m

Decisiod prompt det deplrenty
Táunðerhni

0:25 m

User-defined prompts development

Iterative prompt development

Check results

Responsible handover framework

Domain 1: What is the purpose of the AI tool?

Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

Domain 5: Transparency, licenses, availability and documentation



Domain 2: Where have the training and testing data come from?

Main considerations for evidence synthesists:

- Are training and testing data sources disclosed and accessible?
- Do training and testing data match your evidence synthesis domain/context?
- Are there known limitations or biases (language, geography, date range, methods, etc.)?
- Is there a separation between the training, testing and validation datasets?



Menti

The importance of separation between training data and test data

Whose general medical knowledge are you more likely to trust?

Consider two medical students who have just sat their final exams:

1. Kai scored 95%, but you know that they had access to a leaked copy of the exam paper a week beforehand
2. Ren scored only 68%, but you know they did not cheat in this way

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)





Propliags

Arocept

Aadoraps

Cograles

Siguics



Performance Evaluation

SuperScreener Performance: Systematic Review Screening EVALUATION Scope & Methods

Diversified Review Selection

25 topics from <IMAGE 0>:

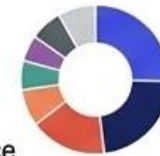
- Public Health
 - Population Health
 - Planetary Health
 - Infectious Disease Epi
 - Neurological
 - Autoimmune, NCDs
-

Standardized Criteria

- ✓ Open Access
- ✓ English
- ✓ 2020-2025
- ✓ Global Author Teams
- ✓ Global Study Geography

Comprehensive Study Designs

- RCTs
- Non-Randomized Studies (NRSs)
- Modelling Studies
- Qualitative Evidence



Data Integrity & Validation



Statement of Non-Contamination: A strict separation protocol was enforced. There was **NO data contamination** between the testing and validation datasets.

Public health interventions addressing the health impacts of climate change

Living review



Menti

Given there are gaps in the testing data and this is an LLM-based AI tool (so we do not know anything about the training data), with what types of mitigations would you consider if you proceeded with this tool?

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)





Responsible handover framework

Domain 1: What is the purpose of the AI tool?

Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

Domain 5: Transparency, licenses, availability and documentation





Domain 3: Is the AI tool validated and performs sufficiently for use?

Main considerations for evidence synthesists:

- Has performance been validated?
- Is the validation published and / or peer-reviewed?
- Is it reported in sufficient enough detail to be transparent and replicable?
- Are metrics appropriate and does performance meet your acceptable threshold?
- Can performance change over time?





The question of acceptable thresholds

- What’s good enough?
- Lack of community consensus on the **appropriate level of confidence** and **appropriate level of performance**

Table 2. Stopping Boundaries and Decision Rules for Interim Analyses

Performance Metrics	Futility Boundaries (Point Estimate)*	Non-inferiority Margins (Upper Limit of 95% CI)*	Decision Rules
Screening			
Sensitivity	<80%**	<95%**	Stop if either boundary is crossed
Specificity (for full-text screening only)	<50%***	<60%***	Stop if either boundary is crossed
Data Extraction			
Sensitivity	<92%**	<97%**	Stop if either boundary is crossed
Major Error Proportion	>3%†	>2%†	Stop if either boundary is crossed
Usability			
System Usability Scale (score)	<57‡	<75‡	Stop if threshold is not met



Responsible handover framework

Domain 1: What is the purpose of the AI tool?

Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

Domain 5: Transparency, licenses, availability and documentation



Domain 4: Usability and user capability

Main considerations for evidence synthesists:

- Does your team have the skills to use it reliably?
- Is training/documentation available?
- What are the costs (financial, time, infrastructure)?
- Is user support available?

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)



Menti

If you got to this stage of the framework and there was adequate training, guidance and support from the AI tool developer, and no reason not to proceed from the previous domains, **how confident do you think you would be in using LLM-based or genAI tools for review tasks?**

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)





Responsible handover framework

Domain 1: What is the purpose of the AI tool?

Domain 2: Where have the training and testing data come from?

Domain 3: Is the AI tool validated and perform sufficiently for use?

Domain 4: Usability and user capability

Domain 5: Transparency, licenses, availability and documentation



Domain 5: Transparency, licenses, availability and documentation

Main considerations for evidence synthesists:

- Are the terms of use clear and acceptable?
 - Plagiarism and provenance; copyright and intellectual property; jurisdiction and licensing; data security and usage rights; and data protection and confidentiality
- Will uploaded content be used for training? Can you opt out?
- Does the tool meet your legal and ethical requirements?
- Is there a conflict-of-interest disclosure?

Menti – quick fire questions

1. You hereby grant (...) to Company a (...) license to reproduce (...) and otherwise use and exploit any and all information and content that you submit to, or use with the Site.
2. When you use our services for individuals such as ChatGPT, Codex, and Sora, we may use your content to train our models.
3. We do not and will not permit any third party to use your content or User Data to improve or train AI models.

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)



Example using Cochrane's RCT classifier



From RAISE 3, Supplementary file 2

Expectation 1: Evidence synthesists are ultimately responsible for their research

Supplementary file 2: example of completed responsible handover tool

What type of handover is this?
<p>Who is completing this form?</p> <p>For example, is it a tool development team handing a tool over to the evidence synthesis community for their use? Is it an evidence synthesis methodologist reporting the results of an evaluation and recommending a tool for use in a specific way? Or is it the result of collaboration between a tool development team and an evidence synthesis methodologist, who are jointly reporting their understanding about the suitability of a tool for use? Record the tool name and version (including model version, where relevant) and the date this form was completed. If any question cannot be answered, record this explicitly (for example, 'Unknown' or 'Not disclosed') and treat missing information as a potential risk signal. Schedule a reassessment date (e.g. annually or when major version updates occur) and record this in the form.</p>
<p>This form is the joint work of information retrieval methodologists at an evidence synthesis organisation (Cochrane), and an AI development team.</p>
(1) What is the purpose of the tool?
<p>Methodologists and AI development teams: agree a statement that defines how the tool should be used (1.1) in an evidence synthesis project. Record the tool name, version, and deployment context (platform, local installation, API, etc.). If the AI tool uses generative AI, state the model provider and model version, and (where relevant) the prompt set or template version. Then, jointly consider whether there are contextual limitations (1.1.2), informed by input from the AI development team, and whether the intended context of use matches that originally envisaged (1.1.1). State whether the use case is to replace an existing task or method, or if it is to be used in parallel with current methods. If the latter, state which existing task or method this new approach should complement or be used alongside, and what human is required.</p>
<p>1.1: The system estimates the likelihood of a given study report describing a Randomized Controlled Trial (RCT), using only the title and abstract. This allows studies that are very unlikely to be RCTs to be automatically screened out during the development of systematic reviews for which this is desirable.</p>
<p>The RCT classifier is designed for use as a binary classifier. The cut-point is set to achieve $\geq 99\%$ recall (to identify at least 99% of reports of RCTs). Within Cochrane, the classifier is implemented as a step in the Screen4Me workflow, with all studies not excluded at this step being screened by a human before an inclusion decision. See figure from Thomas et al 2021. Deployed in this way, the RCT classifier aims to significantly reduce manual</p>

Expectation 2:

**Demonstrate AI will not compromise
the methodological rigor or integrity of
the synthesis**



Decisions after using the framework

Proceed

- AI was validated for its proposed use; supporting evidence is strong
- Risk of well-understood with moderate errors or inconsequential differences
- Limitations known and reasonably manageable

Proceed with mitigations

- Tool shows promise but has gaps in evidence,
- Requires additional monitoring, or presents moderate risks that can be actively managed.

Do not proceed

- AI not validated for its proposed used; missing or weak supporting evidence
- Transparency is insufficient for meaningful assessment
- Risks cannot be adequately mitigated; AI unreliable

Menti – quick fire questions

Would you proceed if...?

The validation study only partially matches your review context

The validation study was conducted by tool developers

The tool requires your consent to train an AI model on your data

The tool's workflow is not transparent nor reproducible

The tool produces transparent outputs and limitations are well understood

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)



Consider not proceeding if...

- No published validation in a relevant context
- Concerns about the replicability of the validation
- Concerns about performance claims based only on the developer's own validation and /or methodological limitations
- Lack of legal or policy compliance at the evidence synthesist's or tool user's organisation, or at the national or international level
- Terms allow use of your content for model training without opt-out
- Inappropriate level of human oversight available, e.g. no way to monitor or audit outputs
- The AI tool developer is unresponsive to questions

Current state of AI tools



Table 2: current (February 2026) state of AI tools

Task	Tool Class	Detail and considerations	Example tools	Recommendation
Writing a protocol				
Question formulation	Generative LLMs	Asking LLMs to provide novel questions for synthesis may support early question development. However, suggestions may be incomplete, irrelevant, subject to bias (based on its sources), or overlap with past reviews.	ChatGPT, CoPilot, Claude, Gemini, DeepSeek	Human verification required
Drafting	Generative LLMs	Pre-trained LLMs can provide an outline using well-established protocol formats. Users may also provide a format / direct the LLM to resources to support this.	ChatGPT, CoPilot, Claude, Gemini, DeepSeek	Human verification required
The Search				
Exploring the literature	Unsupervised	Topic modelling tools aid in identifying clusters of evidence quickly to get a sense of key themes/areas of interest.	Carrot2 (https://search.carrot2.org)	Acceptable for use
	Agentic AI	AI agents develop, refine, and perform searches based on natural language queries. Highly dependent on data sources the tool has access to and requires human input at each stage to guide agent. May be helpful to gain a sense of the literature at an early stage but should not be used as part of any formal evidence retrieval.	Undermind (https://www.undermind.ai/), Elicit (https://elicit.com/), Asta Find Papers (https://asta.allen.ai/)	Acceptable for use
Search strategy development	Rule-based	Tools analyse frequency of keywords and/or controlled vocabularies in search results. Specialised tools are required to cover indexing from different bibliographic databases. May provide additional keywords to inform search strategy but should be used in combination with other search development methods.	Yale MeSH Analyzer (https://mesh.med.yale.edu/), TERA WordFreq (https://tera-tools.com/word-freq), PubReMiner (https://hgserver2.a)	Acceptable for use

Current state of AI tools

Task	Tool Class	Detail and considerations
Writing a protocol		
Question formulation	Generative LLMs	Asking LLMs to provide novel question development. How subject to bias (based on its
Drafting	Generative LLMs	Pre-trained LLMs can provide formats. Users may also provide support this.
The Search		
Exploring the literature	Unsupervised	Topic modelling tools aid in sense of key themes/areas of
	Agentic AI	AI agents develop, refine, and queries. Highly dependent on human input at each stage to the literature at an early stage evidence retrieval.
Search strategy development	Rule-based	Tools analyse frequency of keywords results. Specialised tools are bibliographic databases. May strategy but should be used in combination with other search development methods.


Recommendation	
Acceptable for use	AI outputs may be used directly within the review workflow, if any limitations or potential biases are acknowledged and accounted for.
Human verification required	AI outputs may be used to support review tasks but must be carefully checked by humans before use. The degree of checking required may vary, but typically this will require a human to read and possibly make amendments to the entirety of the output.
Requires validation within the review	AI outputs may be used if their performance is explicitly evaluated within the context of the review itself and deemed adequate (e.g. comparable to human performance).
Exploratory and supplementary use	AI outputs may be used for developing ideas or as a starting point to support understanding. All outputs should be extensively refined by human reviewers prior to use for a review task. Alternatively, outputs may be appropriate for use as an additional, supplementary approach, but without replacing established processes.
Not acceptable for use	The current state of technology means that these AI outputs have such serious limitations, that they should not be relied upon.

	WordFreq (https://tera-tools.com/word-freq); PubReMiner (https://hgserver2.am)
--	---

A word about validation vs. verification in your review

Requires validation
within the review

Human verification
required



medRxiv
THE PREPRINT SERVER FOR HEALTH SCIENCES



[Follow this preprint](#) [Previous](#)

Cochrane Evaluation of (Semi-) Automated Review (CESAR) Methods: Protocol for an adaptive platform study within reviews

Posted April 15, 2026.

[Download PDF](#)
[Print/Save Options](#)
[Author Declarations](#)
[Data/Code](#)

Gerald Gartlehner, Susan Banda, Max Callaghan, Jo-Ana Chase, Andreea Dobrescu, Angelika Eisele-Metzger, Ella Flemyng, Sean Gardner, Ursula Griebler, Bartosz Helfer, Pawel Jemiolo, Biljana Macura, Jan C Minx, Anna Noel-Storr, Noosheen Rajabzadeh Tahmasebi, Amin Sharifan, Joerg Meerpohl, James Thomas

doi: <https://doi.org/10.64898/2026.04.13.26350802>

This article is a preprint and has not been certified by peer review [what does this mean?]. It reports new medical research that has yet to be evaluated and so should not be used to guide clinical practice.

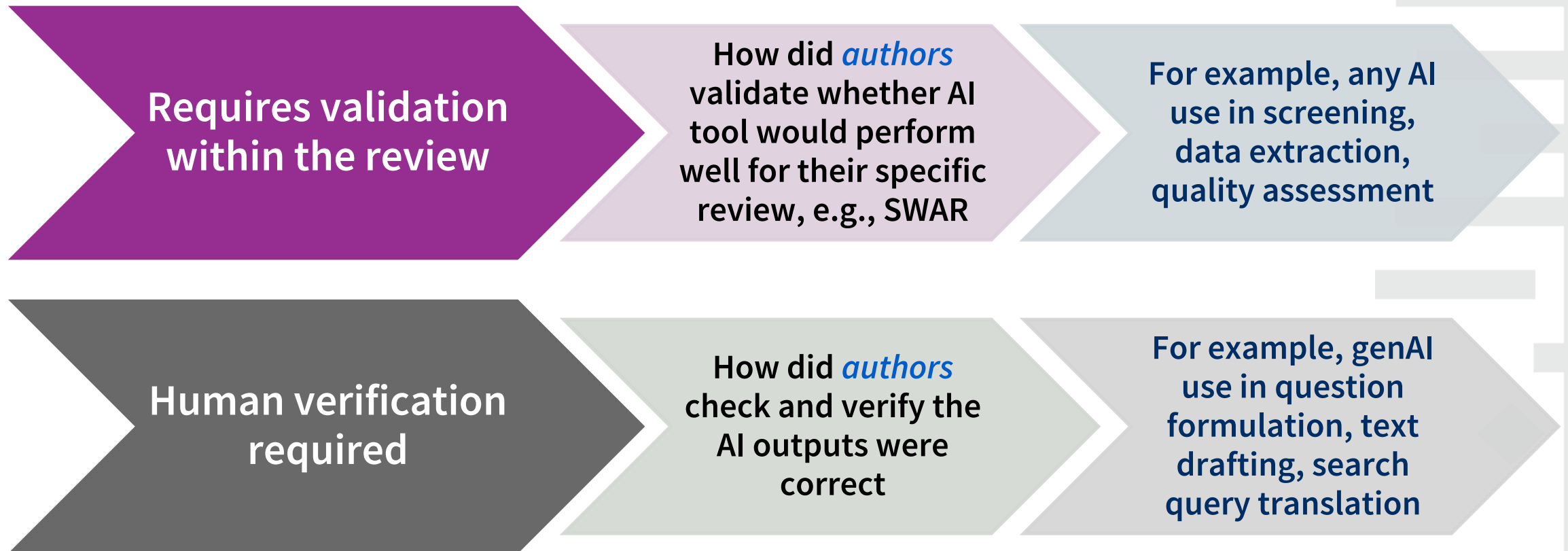
[Abstract](#) [Info/History](#) [Metrics](#) [Preview PDF](#)

[Download PDF](#)

Reviews and Context

- 0 Comment
- 0 TRIP Peer Review

A word about validation vs. verification in your review



**Expectation 3:
Full and transparent reporting of AI
use**



Disclosure of AI use

Name and purpose of AI tool

We will use [AI system/tool/approach name, version, date of use] developed by [organization/developer] for [specific purpose(s)] in [the evidence synthesis process]. The [AI system/tool/approach] will [state it will be used according to the user guide, and include reference, and/or briefly describe any customization, training, or parameters to be applied].

Disclosure of AI use

Name and purpose of AI tool

We will use [*AI system/tool/approach name, version, date of use*] developed by [*organization/developer*] for [*specific purpose(s)*] in [*the evidence synthesis process*]. The [*AI system/tool/approach*] will [*state it will be used according to the user guide, and include reference, and/or briefly describe any customization, training, or parameters to be applied*].

Degree of human oversight

Outputs from the [*AI system/tool/approach*] are justified for use in our synthesis because:

- [*state the degree of human oversight such as any steps taken to review, verify, or override AI-generated outputs.*]

Disclosure of AI use

Name and purpose of AI tool

We will use [AI system/tool/approach name, version, date of use] developed by [organization/developer] for [specific purpose(s)] in [the evidence synthesis process]. The [AI system/tool/approach] will [state it will be used according to the user guide, and include reference, and/or briefly describe any customization, training, or parameters to be applied].

Degree of human oversight

Outputs from the [AI system/tool/approach] are justified for use in our synthesis because:

- [state the degree of human oversight such as any steps taken to review, verify, or override AI-generated outputs.]
- [describe how you have determined it is methodologically sound and will not undermine the trustworthiness or reliability of the synthesis or its conclusions (e.g., model validation, feature validation)]
- [describe how it has been validated or calibrated to ensure that it is appropriate for use in the context of the specific evidence synthesis, to include degree of author involvement, if not covered in the user guide, evaluations or elsewhere (e.g., real-world effectiveness)].

Justify use of AI system or tool (conduct of review)

Limitations [of the AI system/tool/approach] include [describe known limitations, potential biases, and ethical concerns]/ [are included as a supplementary material]. [If applicable] A detailed description of the methodology, including parameters and validation procedures, is available in [supplementary materials].

Other considerations for reporting AI use

- **Methods:**
 - PRISMA items include a note about reporting any automation tools
- **Discussion > Limitations of the review process:**
 - detail any limitations or potential biases, consider the potential impact of errors, limitations or generalizability
- **Declaration of interest**
 - declare any financial and non-financial interests related to the tool, including any relevant interests in the organization(s) that own or fund them

**Expectation 4:
AI should be used with human
oversight**

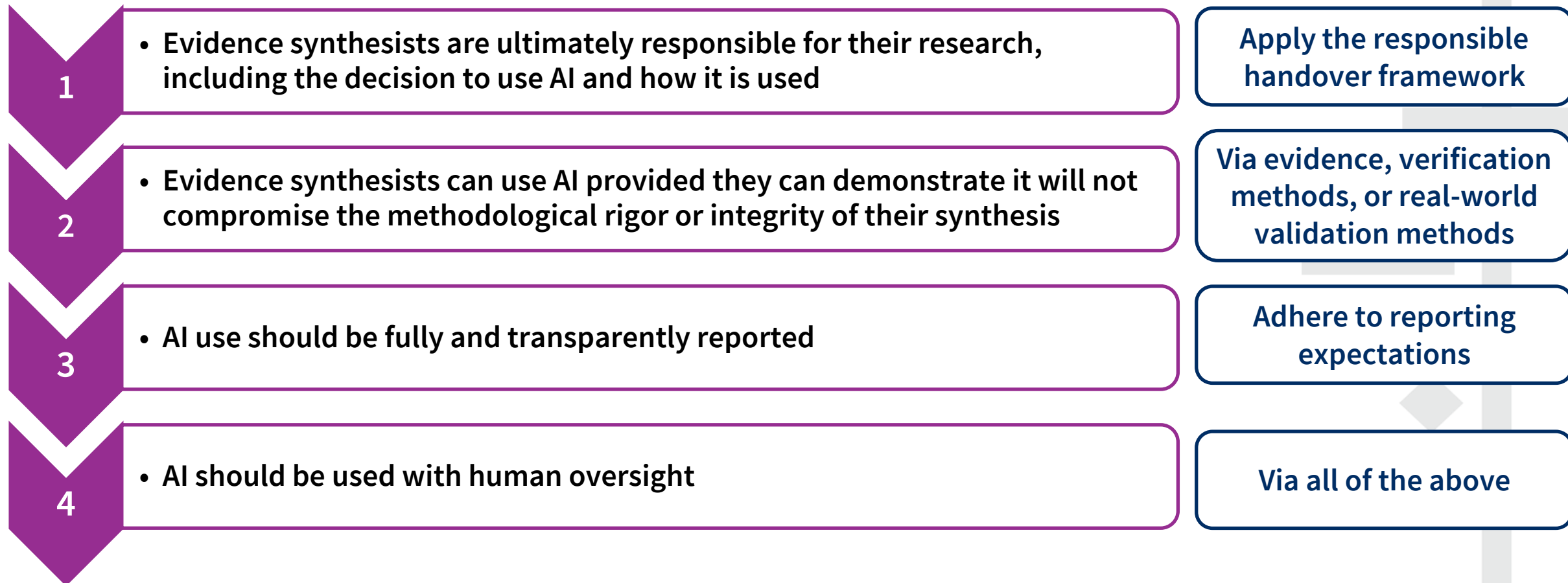


Human oversight and accountability



AI should be a companion, not a replacement

AI use expectations for evidence synthesists



Menti

What was your biggest learning from this framework / this session?

[menti.com / Code: 6181 2178](https://menti.com/Code:61812178)





Thank you!

Questions?

