Nuffield Trust: Rapid evaluation in health care 2025





# About me

- Co-director of EPPI Centre, UCL
- Do a wide variety of evidence synthesis – mostly for Department of Health & Social Care
  - Addressing questions beyond effectiveness
  - Methodological development
- Evidence synthesis methods
- Long-standing area of work in making the review process more efficient using new technologies



# Acknowledgements and declaration of interests

- I am employed by University College London; receive funding from the funders below for this and related work; lead EPPI-Reviewer software development
- Cochrane roles: Review author; Co-convenor Joint Artificial Intelligence Methods Group<sup>NEW!</sup>; Co-Senior Scientific Editor Cochrane Handbook; support Cochrane with information technologies (EPPI-Reviewer and machine learning)
- Guidance for responsible use of AI in systematic reviews (RAISE)
- Parts of this work funded by: Wellcome Trust, National Institute for Health Research (NIHR), Education Endowment Foundation, Youth Endowment Foundation

# In this session

- Introduction to AI / machine learning / automation tools for evidence synthesis (and how they work)
- Hands-on experimentation with LLM tools for evidence synthesis
- Please feel free to ask questions as we go



+

0



https://eppi.ioe.ac.uk/cms/Default.aspx?tabid=3677

## **Evidence synthesis priorities**

Evidence syntheses are often used to inform decisions that affect people's lives

Evidence synthesists favour accuracy over efficiency

Highly sensitive searches are required to avoid selection bias

Highly accurate quality assurance processes are required to avoid human error



# Impact of these priorities

- An inefficient, resource-intensive process has evolved that produces reliable, but expensive and time consuming, reviews
- We cannot keep pace with the deluge of new research being published
- E.g. in the Cochrane Reviews published March 2014, > 163k citations were screened; 6,599 full text reports were read; and 703 studies were included
- That's about 2 million records per year

# This means

- Only a fraction of available studies are included in evidence syntheses
- Evidence synthesis does not cover all questions/ domains comprehensively
- We don't even know when reviews \*need\* to be updated



### 

### Are there AI tools we can use?



How To Automate Your Literature Review ETHICALLY Using ChatGPT (Prof. David Stuckler) 144 views - 4 months ago David Stuckler

. Writing the literature review How To Read Research Papers Effectively: https://youtu.be/WVv

Intro | Finding your research question | Developing an outline... 6 chapters 🗸

### 

### Are there AI tools we can use?



How To Automate Your Literature Review ETHICALLY Using ChatGPT (Prof. David Stuckler) 144K views • 4 months ago Wovid Stuckler

Writing the literature review How To Read Research Papers Effectively: https://youtu.be/WVv2

Intro | Finding your research question | Developing an outline... 6 chapters



INSTITUTE FOR Evidence-Based Healthcare



### **Review Wizard**



**Scope of action:** Write the method section for your review. **Purpose:** To Design and write your review methods section.

**Step 1:** You will need to be logged in and click on Review Wizard URL. <u>https://tera-tools.com/methods-wizard</u>



Consensus

Find the best science, faster.

### **Consensus Product Update 3/24**

We are excited to announce the launch of one of our most requested features ever: upload and chat with a PDF within Consensus <sup>1</sup>

Less PDF scrolling, more time-saving analysis. This new feature allows you to apply Consensus's models to your research paper library. Upload and chat with the full-text of your papers to ask about key figures, methodological details, novel insights and more!

This launch marks the start of a string of major changes to Consensus in the coming weeks. Upcoming changes will unlock a whole new level of AI analysis including full-text access, multi-paper upload & analysis, and more!



tion for your review. ew methods section.

d click on Review Wizard URL.



Try out our newest feature!

🕏 Elicit Tutorial

NEW

R

Consensus onsensus Product Up are excited to announce the launce load and chat with a PDF within C as PDF scrolling, more time-saving - nsensus's models to your research ur papers to ask about key figures, r s launch marks the start of a string eks. Upcoming changes will unlock	Find the best science, faster.	NSTITUTE FOR Evidence-Ba	Automa System Review	atic
t access, multi-paper upload & anal Try out our newest feature! Literature Revie Find Your Research Question	ysis, and more!	Lep 1: You wil https://tera-tou	I need to be logged in and click o	n Review Wizard URL.



### Are there AI tools we can use? There are a lot of AI-based evidence synthesis tools!

- Can we use them?
- Should we use them?
- And are we already being outevolved if we're not using AI?

 Important to understand a bit about how automation tools work to make good decisions about using them





# Four machine learning / automation paradigms

- Rules-based approaches
  - (strictly speaking, not *machine learning*)
- Unsupervised approaches
- Supervised approaches
- Generative approaches ('Gen AI')
- Covering in terms of technology not purpose, so we can consider their strengths and weaknesses more easily

# Rules-based approaches

As you might guess... a set of rules is constructed by humans and given to the machine

### For example

Look up a	Use of	lf a given phrase is	Many citation duplicate-
words	synonyms	present, apply a given code	checking algorithms



### Rules can be accurate... but fragile



If you stick within the rules, you get the anticipated results



If you stray outside – even a little bit – the rule can fail altogether



No grey area – it works, or completely fails



- The machine is given no rules...
- And simply identifies patterns in the data

– E.g.

- Relationships between words
- Clustering documents



Unsupervised approaches can help you explore patterns in your data Attractive visualisations are possible

 $\leftarrow$ 

Emb

DATA





### Top 100 results of about 268158 for smoking

### 1 Prospective, multi-centric benchmark study assessing delirium; prevalence, incidence and its correlates in hospitalized elderl Lebanese patients. 🐾 🗗 🔍

total cholesterol

pharmacist intervention

With the increase in the proportion of elderly Lebanese patients, little is known about delirium's prevalence, incidence and correlated factors. ... To identify the prevalence, incidence and factors associated with overall and incident delirium in hospitalized elderly Lebanese patients.

http://www.ncbi.nlm.nih.gov/pubmed/3120352

### Familial cancer of unknown primary, 🐾 🗗 🔍

Cancer of unknown primary site (CUP) is a deadly disease diagnosed through metastases at various organs without primary tumor identification. Despite the major molecular and technological advances, the carcinogenesis of CUP remains enigmatic which hampers adequate study design of treatments leading to survival improvement. To date, the pathogenesis of CUP is still debatable with one hypothesis considering CUP is simply a group of metastatic tumors with unidentified primaries and another considering it a distinct entity with specific genetic and phenotypic aberrations. Familial CUP seems to favor the first hypothesis due to common genetic predisposition factors between known primaries and CUP. Two clinical implications may be withdrawn from the pathogenesis of familial clustering of CUP. The detailed family history and environmental risk factors may orient towards the primary tumor identification. In cases of familial, smoking avoidance and adherence to general population guidelines for cancer screening would be strongly encouraged.

http://www.ncbi.nlm.nih.gov/pubmed/31203526

Discovery of biomarkers for glycaemic deterioration before and after the onset of type 2 diabetes: descriptive characteristics of the epidemiological studies within the IMI DIRECT Consortium, 🐁 🥵 🥥

Here, we describe the characteristics of the Innovative Medicines Initiative (IMI) Diabetes Research on Patient Stratification (DIRECT) epidemiological cohorts at baseline and follow-up examinations (18, 36 and 48 months of follow-up). http://www.ncbi.nlm.nih.gov/pubmed/31203377

4 Health behaviours and mental and physical health status in older adults with a history of homelessness; a cross-sectiona population-based study in England, 🞭 🗗 🔍

This study compared (1) levels of engagement in lifestyle risk behaviours and (2) mental and physical health status in individuals who have previously been homeless to those of individuals who have not. http://www.ncbi.nlm.nih.gov/pubmed/31203244

### 5 Combined effects of lung function, blood gases and kidney function on the exacerbation risk in stable COPD: Results from the COSYCONET cohort, 🐾 🗗 🔍

Alterations of acid-base metabolism are an important outcome predictor in acute exacerbations of COPD, whereas sufficient metabolic compensation and adequate renal function are associated with decreased mortality. In stable COPD there is, however, only limited information on the combined role of acid-base balance, blood gases, renal and respiratory function on exacerbation risk grading.

http://www.ncbi.nlm.nih.gov/pubmed/31203096

6 "Don't smoke in public, you look like trash": An exploratory study about women's experiences of smoking-related

v3.16.0-SNAPSHOT | build 277 | 2018-05-17 11:55 © 2002-2019 Stanislaw Osinski, Dawid Weiss

### Unsupervised approaches lack control



Very powerful – can reveal relationships in the data which are not necessarily obvious



Very efficient – data often need no preparation



But... you don't get to tell the machine which classifications to make

# Supervised approaches



Humans prepare 'training' data – containing data + labels which describe the desired classification



For example

Image recognition Text classification



### Good supervision is required....



Very dependent on quality and coverage of training data



Performance very dependent on context



# Example of study classification: RCT Classifier

- A classifier was built using more than 280,000 records from Cochrane Crowd
- It is 'simply' applying single classification (RCT / not RCT)
- It has been calibrated to achieve a recall = 99% on the McMaster 'Hedges' dataset
  - Calibration = ranking the 'test' dataset by score
  - $\odot\,\text{BUT}$  precision is low
- It is very accurate!

 But not all supervised learning can be so accurate, as lots of high-quality training data are needed



## **Generative approaches**



![](_page_24_Picture_3.jpeg)

![](_page_24_Picture_4.jpeg)

ChatGPT (or other LLM chatbot)

LLM-based database querying and summarisation

LLM-based information extraction

# **Building a GenAl chatbot**

Pretraining (unsupervised ML))

![](_page_25_Picture_3.jpeg)

'Naïve' model Cannot 'chat'; nextword prediction only Fine-tuning (supervised ML)

RLHF (supervised ML)

![](_page_25_Picture_7.jpeg)

Model can now 'chat' and answer questions Model produces 'better' and less toxic answers

Output

## **Generative LLM operation**

Input

![](_page_26_Picture_3.jpeg)

![](_page_27_Figure_2.jpeg)

![](_page_28_Figure_2.jpeg)

The selected word is added to the input

![](_page_29_Figure_2.jpeg)

![](_page_30_Figure_2.jpeg)

![](_page_31_Figure_2.jpeg)

Important to bear in mind that the system does not plan ahead ...and at no point does it check the accuracy of what is 'said'

![](_page_32_Figure_2.jpeg)

Is it about topic y?

Instead of the prompt containing "There's no place like…" it could contain a question about a passage of text that is also in the prompt

≜UC L

![](_page_33_Picture_1.jpeg)

### A major contrast with supervised machine learning:

## 'zero shot' or 'in context' learning

Image generated with the help of Microsoft Copilot

### Why zero-shot learning is a gamechanger

Development and evaluation of the Cochrane RCT Classifier (Using conventional supervised machine learning)

![](_page_34_Figure_2.jpeg)

Conventional machine learning model trained on 280,000 records from Cochrane Crowd

_		

o achieve 99% recall n a second Hedges') dataset
n a second Hedges') dataset
Hedges') dataset
~50,000 records)

Ο

Model was validated
on 92,000 studies
included in Cochrane
intervention reviews

Model was deployed for live use in Cochrane review workflows

### Why zero-shot learning is a gamechanger

Development and evaluation of the Cochrane RCT Classifier

![](_page_35_Figure_2.jpeg)

![](_page_35_Picture_3.jpeg)

Model was validated on 92,000 studies included in Cochrane intervention reviews Model was deployed for live use in Cochrane review workflows

With the new AI tools there's no need to create (expensive / hard to find) training data

### Why zero-shot learning is a gamechanger

Development and evaluation of a classification task using a language model

![](_page_36_Figure_2.jpeg)

Instead, a human writes some prompts for a large language model in their normal language

They check they work on their data

C	)
0~	6
	L
O	0

The language model can then apply the prompts to the remaining data

### Data (information) extraction

- Earlier language models lacked precision & limited context 'window'
- Newer models have larger windows and offer impressive early results
- E.g. Claude2, published by Anthropic

ANTHROP\C

Product Research Company News Careers

# AI <u>research</u> and <u>products</u> that put safety at the frontier

NEW

### Claude in Beta, now available!

Your friendly assistant. Fast, capable, and truly conversational.

Talk to Claude

ENTERPRISE

**Build with Claude** 

Start using Claude and unlock business value with AI.

Submit business interest

Smoking Cessation Internet Intervention for Low and High SES Groups ~ 198.93 KB

Here is a report of a randomized trial. State: the population recruited; the intervention(s) evaluated; the main outcome(s) measured; and the statistical results of the primary outcome

According to the paper, this randomized controlled trial had the following key details:

### Population:

- 4613 participants aged 18+ years who smoked daily and were willing to make a serious quit attempt in the UK.
- Recruited online, mainly via the English Department of Health website.
- Stratified into low socioeconomic status (n=2142) and high socioeconomic status (n=2471) based on occupation.

### Interventions:

- Interactive internet-based smoking cessation website called StopAdvisor.
- Control group received a 1-page static website with brief standard smoking cessation advice.

### Main outcome:

 Primary outcome was biochemically verified 6-month sustained abstinence from smoking (self-report of smoking <5 cigarettes in the past 6 months plus saliva cotinine <15 ng/mL).</li>

### **Results:**

- No overall difference between groups in 6-month abstinence (StopAdvisor 10% vs control 10%, RR 1.06, 95% CI 0.89-1.27, p=0.49).
- In the low SES subgroup, 6-month abstinence was higher in StopAdvisor vs control (8% vs 6%, RR 1.36, 95% Cl 1.00-1.86, p=0.0499).
- No difference between groups in the high SES subgroup (StopAdvisor 12% vs control 13%, RR 0.95, 95% CI 0.77-1.17, p=0.61).

### Internet-based intervention for smoking cessation (StopAdvisor) in people with low and high socioeconomic status: a randomised controlled trial

Jamie Brown, Susan Michie, Adam WA Geraghty, Lucy Yardley, Benjamin Gardner, Lion Shahab, John A Stapleton, Robert West

### Summary

( JT

**Background** Internet-based interventions for smoking cessation could help millions of people stop smoking at very low unit costs; however, long-term biochemically verified evidence is scarce and such interventions might be less effective for smokers with low socioeconomic status than for those with high status because of lower online literacy to engage with websites. We aimed to assess a new interactive internet-based intervention (StopAdvisor) for smoking cessation that was designed with particular attention directed to people with low socioeconomic status.

Methods We did this online randomised controlled trial between Dec 6, 2011, and Oct 11, 2013, in the UK. Participants aged 18 years and older who smoked every day were randomly assigned (1:1) to receive treatment with StopAdvisor or an information-only website. Randomisation was automated with an unseen random number function embedded in the website to establish which treatment was revealed after the online baseline assessment. Recruitment continued until the required sample size had been achieved from both high and low socioeconomic status subpopulations. Participants, and researchers who obtained data and did laboratory analyses, were masked to treatment allocation. The primary outcome was 6 month sustained, biochemically verified abstinence. The main secondary outcome was 6 month, 7 day biochemically verified point prevalence. Analysis was by intention to treat. Homogeneity of intervention effect across the socioeconomic subsamples was first assessed to establish whether overall or separate subsample analyses were appropriate. The study is registered as an International Standard Randomised Controlled Trial, number ISRCTN99820519.

**Findings** We randomly assigned 4613 participants to the StopAdvisor group (n=2321) or the control group (n=2292); 2142 participants were of low socioeconomic status and 2471 participants were of high status. The overall rate of smoking cessation was similar between participants in the StopAdvisor and control groups for the primary (237 [10%] vs 220 [10%] participants; relative risk [RR] 1.06, 95% CI 0.89–1.27; p=0.49) and the secondary (358 [15%] vs 332 [15%] participants; 1.06, 0.93–1.22; p=0.37) outcomes; however, the intervention effect differed across socioeconomic status subsamples (1.44, 0.99–2.09; p=0.0562 and 1.37, 1.02–1.84; p=0.0360, respectively). StopAdvisor helped participants with low socioeconomic status stop smoking compared with the information-only website (primary outcome: 90 [8%] of 1088 vs 64 [6%] of 1054 participants; RR 1.36, 95% CI 1.00–1.86; p=0.0499; secondary outcome: 136 [13%] vs 100 [10%] participants; 1.32, 1.03–1.68, p=0.0267), but did not improve cessation rates in those with high socioeconomic status (147 [12%] of 1233 vs 156 [13%] of 1238 participants; 0.95, 0.77–1.17; p=0.61 and 222 [18%] vs 232 [19%] participants; 0.96, 0.81–1.13, p=0.64, respectively).

![](_page_38_Picture_22.jpeg)

### Lancet Respir Med 2014

€ 🖉 🖗 🗲

Published **Online** September 25, 2014 http://dx.doi.org/10.1016/ S2213-2600(14)70195-X

See Online/Comment http://dx.doi.org/10.1016/ S2213-2600(14)70214-0

**Cancer Research UK Health** Behaviour Research Centre, Department of Epidemiology and Public Health (J Brown PhD, B Gardner DPhil, L Shahab PhD, Prof R West PhD) and Department of Clinical, Educational, and Health Psychology (Prof S Michie DPhil), University College London, London, UK; National Centre for Smoking Cessation and Training, London, UK (Prof S Michie, Prof R West); Primary Care and Population Sciences (A W A Geraghty PhD) and School of Psychology (Prof L Yardley PhD), University of Southampton, Southampton, UK; Addictions Department, Institute of Psychiatry, Kings College London, London, UK (J A Stapleton MSc) Correspondence to: Dr Jamie Brown, Health

between participants in the StopAdvisor and control, subsample (n=1687), the results were consistent with the groups for both the primary (237 [10%] vs 220 [10%] participants; relative risk [RR] 1.06, 95% CI 0.89–1.27; p=0.49) and the secondary (358 [15%] vs 332 [15%] participants; 1.06, 0.93 - 1.22; n=0.37) outcomes. However, , 818 participants; R analysis of the interaction between intervention and socioeconomic status showed clear evidence of nonignorable heterogeneity of intervention effect by both primary (RR 1.44, 95% CI 0.99-2.09; p=0.0562) and secondary  $(1 \cdot 37, 1 \cdot 02 - 1 \cdot 84; p=0 \cdot 0360)$  cessation measures. This finding was evident before and after adjustment for all other baseline characteristics (adjusted data not shown). Consequently, the analysis of outcome was done separately within each of the two socioeconomic status subsamples. In the subsample of participants with low

socioeconomic status, a benefit of StopAdvisor was evident for both primary and secondary measures compared with the information-only website, whereas in those with high socioeconomic status, no evidence of a difference was shown (table 2). Adjustment for all baseline characteristics had a negligible effect on these comparisons (table 2). In a post-hoc sensitivity analysis, we re-examined the effect of StopAdvisor on biochemically verified smoking cessation in participants with low socioeconomic status, after exclusion of those in full-time education from the classification (n=237). The benefit of StopAdvisor compared with the informationonly website remained in both unadjusted (primary outcome 87 [9%] of 967 vs 60 [6%] of 938 participants; RR 1.41, 95% CI 1.03-1.93; p=0.0346; secondary

Secondary outcome (point prev

The overall rate of smoking cessation was similar who did not have post-16 education. In this smaller primary analyses for StopAdvisor versus information on y, but were nor

> secondary outcome  $1 \cdot 21, 0 \cdot 93 - 1 \cdot 58; p =$ 1.27, 0.92-1.75;0.92 - 1.58; = 0.17analysis, we re-exai self-reported rather of smoking cessati basis of similar verification criteria new analyses show reported in table 2 ( socioeconomic sta compared with info not significant 141 [13%] of 1088 unadjusted RR 1.20 RR 1.23, 0.97-1 prevalence 227 [21%

RR 1.13, 0.95–1.34

p=0.07). That statistical

in participants with

failed to reach sign

decreased power

between the two gro

A\

fimary outcome

State, in tabular form, the outcomes of the sutdy

![](_page_39_Picture_5.jpeg)

Here is a table summarizing the key outcomes of the study:

	StopAdvisor	Control	Relative Risk (95% CI)	p-value
Overall				
6-month abstinence	10%	10%	1.06 (0.89-1.27)	0.49
Low SES				
6-month abstinence	8%	6%	1.36 (1.00-1.86)	0.0499
High SES				
6-month abstinence	12%	13%	0.95 (0.77-1.17)	0.61
🖞 Сору	1	1	, 	Ţ

	StopAdvisor	Control	Relative nsw (25% CI)	Odds ratio (95% CI)*	Percentage-point difference (95% CI)	p value†
Primary outcom	ne (abstinence for 6 r	nonths)				
High SES	147/1233 (12%)	156/1238 (13%)	0·95 (0·77 to 1·17)	0 )4 (0·74 to 1·19)	-0.68 (-3.27 to 1.91)	0.61
Adjusted			0·97 (0·78 to 1·19)‡	·95 (0·75 to 1·22)‡		0.75
L. V. SES	90/1088 (8%)	64/1054 (6%)	1·36 (1·00 to 1·86)	1·39 (1·00 to 1·94)	2·20 (0·02 to 4·38)	0.0499
Adjusteu			1·43 (1·05 •	1·46 (1·04 to 2·05)‡		0.0238

rates.

The study authors did not report a risk of bias assessment using the Cochrane RoB tool. However, here are some observations about potential risks of bias based on the details provided in the paper:

- Random sequence generation: Participants were randomized using an automated unseen random number function, so this should have been adequate.
- Allocation concealment: The randomization was automated with no experimenter involvement, so allocation was likely adequately concealed.
- Blinding of participants/personnel: Not reported, but blinding of participants may not have been feasible given the difference in intervention. Personnel were likely not blinded.
- Blinding of outcome assessment: Outcome assessors were reported to be blinded to treatment allocation.
- Incomplete outcome data: Follow-up rates were 77% overall and similar between groups, with missing data treated as smoking, so the risk appears low.
- Selective reporting: The outcomes specified in the protocol appear to be reported, suggesting low risk of bias.
- Other biases: The study appears to have been adequately powered and analyzed as pre-specified, with low risks of other biases.

So in summary, the lack of blinding of participants and personnel, as well as the lack of formal risk of bias assessment by the authors, suggest some uncertainties but overall the study seems reasonably well conducted.

A\

 $\nabla$ 

(JT)

![](_page_41_Picture_0.jpeg)

### Sounds amazing, but...

- Concerns about bias
- Appropriation of content for training without permission
- Environmental impact
- Lack of transparency
- Currently a lack of an evidence base
- A great deal of hype (and money) means claims are currently running ahead of capabilities

![](_page_42_Picture_7.jpeg)

### In summary

Rule-based	Unsupervised	Supervised	Generative
<ul> <li>Not fashionable</li> <li>Potentially powerful</li> <li>Very demanding in time</li> <li>Rules can be fragile</li> </ul>	<ul> <li>Very little time effort required to create rules or training data</li> <li>No control over classifications</li> </ul>	<ul> <li>Makes use of data created for other purposes</li> <li>Does not break as easily as rule-based approaches</li> <li>Can predict specific classification terms (unlike unsupervised)</li> <li>Can require large quantities of training data which may be unavailable / uneconomic to create</li> </ul>	<ul> <li>Considered current 'state of the art'</li> <li>Huge research focus</li> <li>Sometimes beats simpler models (though sometimes only marginally)</li> <li>Concerns about bias and other negative outcomes</li> </ul>

![](_page_44_Picture_1.jpeg)

Understanding their architectures is critical for understanding the strengths and limitations of these new tools

- ChatBot
- 'Semantic' search
- Retrieval Augmented Generation (RAG)

Should I used generative AI tools to help with evidence synthesis?

Yes, you *can* and *should consider* using generative AI tools to support evidence synthesis — but **with clear boundaries** and **critical oversight**. Here's a breakdown of how they can help, and what you should be careful about:

### Vays Generative AI Can Help with Evidence Synthesis

- 1. Search and Screening Support
  - Summarize abstracts quickly.
  - Assist in relevance screening by generating concise overviews.
  - Classify studies by population, intervention, outcome, etc.

### 2. Data Extraction Aid

- Pull structured information (e.g., study design, sample size) from articles.
- Convert narrative findings into standardized formats.
- Help identify outcome measures and comparator groups.
- 3. Thematic Analysis and Coding
  - Suggest initial codes for qualitative synthesis (e.g., for a meta-ethnography).
  - Identify emerging themes across studies.

### Ask anything

Search Q Reason A Deep research & Create image ...

Ŷ

ChatGPT can make mistakes. Check important info. See Cookie Preferences.

# Strengths and limitations: chatbot

![](_page_46_Picture_2.jpeg)

- Can be asked questions in standard prose
- Can provide accurate answers quickly
- But
- Frequency biased
- 'Hallucinate'
- Sounds confident, but is often wrong

### 'Semantic search'

![](_page_47_Figure_2.jpeg)

located and returned to the user

# Strengths and limitations: vector indexes

![](_page_48_Picture_2.jpeg)

- Can provide more semantically powerful searches
- Less 'fragile' than a Boolean search (and not necessary to know all relevant terms in advance)
- BUT
- Dependent on the right documents being available for indexing
- Dependent on the query being sufficiently 'similar' to the documents being retrieved
- Little in the way of an evidence base to support their use in evidence synthesis

### **Retrieval Augmented Generation**

![](_page_49_Figure_2.jpeg)

User queries are translated into vectors; the 'closest' chunks of documents to that query are located; the LLM then generates an answer to the user's query, based on the chunks of text returned

# Strengths and limitations: retrieval augmented generation

![](_page_50_Picture_2.jpeg)

- Can provide a powerful interactive experience where users can 'chat' to their documents using standard prose
- BUT
- Has many of the limitations of BOTH chatbots and vector indexes:
  - Can hallucinate
  - Requires good translation from query to retrieval AND question to the LLM
  - What if all the relevant documents are not retrieved?
  - What if irrelevant documents are retrieved?

Important questions to ask of LLM-based evidence synthesis tools

- For chatbots:
  - Can I verify its accuracy?
  - (Does it matter if not?)
- For search tools:
  - Are the records I need indexed?
  - How can I check that its retrieved everything it should?
- For 'RAG'-based approaches:
  - Are the right documents indexed?
  - Are the right documents retrieved?
  - Are incorrect documents avoided?
  - If present, does the summariser check that the research is reliable / that combining them is a valid thing to do?

# Now it's your turn!

A ST AND AND A ST A

to the st

Try out some tools

### 

...

Export

![](_page_53_Figure_1.jpeg)

### **Undermind Research Assistant**

Welcome! What research topic are you looking for?

I'll ask one or two questions to clarify your goals, then I'll do a deep search to find precisely relevant research papers for you.

You can tell me exactly what you want, like a colleague, and I'll understand. The more you explain, the better I can help, so please be as detailed and specific as possible.

I want to find...

③ How Undermind Works  $\equiv$  Examples

Try Undermind.AI

https://app.undermind.ai/

(8)

•••

÷

С

### Save your work

Logging in allows you to revisit questions and answers in the future.

ĥ

https://paperfinder.allen.ai/chat

→ Login

**Recent Searches** 

### ✤Ai2 Paper Finder

A research tool for paper discovery, with broad and deep coverage via a corpus of 8M+ full text papers and 108M+ abstracts. A project from Ai2.

### Q Papers introducing a dataset of Q Generative document retrieval models that ... Q Shallow marine ecosystem classification using . Q Long term memory in LLM agents Q Papers by Dan Weld about planning Q The Brown clustering paper

### Try the Ai2 paper finder and synthesis tool

https://paperfinder.allen. ai/chat

### Scholarly tools

Discovery

Find papers in your field, from popular ones to more niche and hard-to-find works.

### Ê

### Synthesis

Ask a question and get a comprehensive answer that synthesizes and cites multiple papers.

### **\***+

### More Coming Soon

We're actively developing new and exciting AI tools for scientists. Stay tuned!

### SEMANTIC SCHOLAR

Contact • Privacy Policy • Terms of Service • Responsible Use

James Thomas

Share

Details 7

Details 7

Details 🛛

 $\rightarrow$ 

Save PDF

A»

\$

Ø

Ζ

🛨 Upgrade

Report

Status

 $\oslash$ 

 $\oslash$ 

 $\oslash$ 

Chat

data

Try Elicit

https://elicit.com/

3

Gather papers

50 papers found

Screen papers

Extract data

Generate report

10 papers included

50 data points extracted

Ask anything about the report or its underlying

∽≞

Help

### 🕏 Elicit 🕒 Recent 🗌 Library

Optimal Papers for Meta-Regression Covaria...

### ••••• Research report View only $\checkmark$

MAY 13, 2025

### How many papers do I need per covariate in a meta regression?

Meta-regression paper requirements range from 8-40 studies, with higher numbers needed for complex models or high heterogeneity.

### ABSTRACT

This synthesis of simulation, theoretical, and methodological studies reveals no single formula for the number of papers (or effect sizes) needed per covariate in meta-regression. \* Fang and Zhang (2020) report that at least 20 effect sizes are required for parameter estimation, while Hedges et al. (2010) suggest a range of 20–40 studies ensures robust variance estimation. \* Jenkins and Quintana-Ascencio (2020) illustrate that when heterogeneity is low, as few as 8 studies may suffice, whereas high variance may require up to 25 studies. \* Mathur and VanderWeele (2020) indicate that, depending on the metric, 10–20 studies can be adequate. \* In addition, Tipton (2015) and Higgins and Thompson (2004) note that degrees of freedom depend on both the number and type of covariates, implying that more complex models may demand higher overall sample sizes. \* Thus, while explicit guidance on papers per covariate is rarely provided, the available evidence shows that requirements range from 8 to 40 studies overall, with adjustments based on model complexity and heterogeneity. \*

### $\textbf{METHODS} \rightarrow$

We analyzed 10 papers from an initial pool of 50, using 5 screening criteria. Each paper was reviewed for 5 key aspects that mattered most to the research question. More on methods

### RESULTS

### **Characteristics of Included Studies**

Study	<u> </u>	Study Design	<u> </u>	Primary Focus	<u> </u>	Statistical Approach = E	Sample Size Recommendations =	Full text retrieve	
									•

# Try out one or more tools

- Unsupervised machine learning
  - Carrot2 Workbench
- Generative machine learning
  - Ai2 paper finder and synthesis tools
  - Undermind Al
  - Elicit
  - RobotReviewer
- One of the tools for search strategy development
- Ask yourself
  - Is it clear how the tool works?
  - Can I tell whether it can find all the material that is relevant to my query?
  - How much do I trust its output?

https://eppi.ioe.ac.uk/cms/Default.aspx?tabid=3677

![](_page_58_Figure_1.jpeg)

https://eppi.ioe.ac.uk/cms/Default.aspx?tabid=3677

# Now it's your turn!

Carlos S

Discussion

### Conclusions

– More research is needed?

That's great! There's an evidence base that can inform this, right?

Right.

We're going to write guidance on using AI in evidence synthesis

### We were asked to write some guidance...

- ... about which tool to use, and when
- But found we couldn't!
- The evidence base on which to base our advice was very limited
- AI tools were being developed that were not engineered to be fit-for-purpose

# Vision: RAISE guidance for the responsible use of AI in evidence synthesis

- A draft of the guidance and recommendations is now online for consultation
- Our vision is for it to be a 'living' set of guidelines, that is updated through community input and helps to define roles & responsibilities within the ecosystem
- Should the ecosystem develop in this wellorganized way, we hope to see the development of AI tools that adhere to the principles of research integrity, and so enable evidence accessibility in equitable and rigorous ways

![](_page_62_Picture_4.jpeg)

![](_page_62_Picture_5.jpeg)

![](_page_62_Picture_6.jpeg)

### Roles-based ecosystem

- We need to support the wider adoption of AI to overcome the increasing burden of doing timely and cost-effective evidence synthesis
- We need cross-field standards to support the development of appropriate and responsible AI
- We anticipate an ecosystem made up of individuals, collaborations, and organisations which each have a role to play in developing and using Al in a responsible way
- (one person / organisation may play multiple roles)

![](_page_63_Figure_6.jpeg)

### How you can get involved (1)

- The link : <u>https://osf.io/fwaud/</u>
- Timetable for development
  - A new version will be published in the next few weeks
- Three documents:
  - Roles-based recommendations for practice
  - Guidance on building and evaluating Al tools
  - Guidance on selecting and using AI tools
- Do take a look and let us know what you think!

![](_page_64_Picture_10.jpeg)

### How you can get involved (2): 'Studies Within A Review' (SWARs)

👚 🐌 Page: 1 of 11

- + Automatic Zoom ÷

### Section 2: SWAR Title

### Title:-

Generative artificial intelligence (AI) tools versus conventional screening by humans for selecting eligible study reports for evidence synthesis: a living study within a review (living SWAR) – retrospective version.

### Section 3: Objective of This SWAR

Objective:-

To retrospectively assess the performance of generative AI tools for selecting eligible study reports for inclusion in systematic reviews or maps of research

### Section 4: Additional SWAR Details

Study Area (1):-STUDY IDENTIFICATION

Sample Type (1):-OTHER – Records / reports of studies

Estimated Funding Level Needed:-LOW

PERSPECTIVE	WILEY
Study within a review (SW	AR)
Deelen Devene <sup>123</sup> Nilita N. Burk	212   Shaun Trausalu <sup>4</sup>   Miles Clarke <sup>5</sup>
James Thomas <sup>6</sup> Andrew Booth <sup>7</sup>	Andrea C. Tricco <sup>8,9,10</sup> K. M. Saif-Ur-Rahman <sup>1,2</sup>
<sup>1</sup> Evidence Synthesis Ireland and Cochrane Ireland, University of Galw	/ay, Galway, Ireland
<sup>3</sup> HRB-Trials Methodology Research Network, University of Galway, Galway, G	Salway, Ireland
<sup>4</sup> Health Services Research Unit, University of Aberdeen, Aberdeen, U <sup>5</sup> Northern Ireland Methodology Hub, Queen's University Belfast, Bel	K fast LIK
<sup>6</sup> EPPI-Centre, UCL Social Research Institute, University College Lond	lon, London, UK
Research (ScHARR), University of Sh te, St. Michael's Hospital, Unity Hea stitute of Health Policy, Managemen ealth Care Quality: A JBI Centre of E	<ul> <li>More consistency in methods, tasks and questions</li> </ul>
hesis Ireland and Cochrane Ireland, t ityofgalway.ie	<ul> <li>Enabling cumulation across studies (which may be small-N)</li> </ul>
	<ul> <li>Invitation to join a 'living' SWAR evaluating the use of LLMs for title &amp; abstract / full text screening</li> </ul>
	<u>https://osf.io/g7mkb/</u>

Devane D, Burke NN, Treweek S, Clarke M, Thomas J, Booth A, Tricco AC, Saif-Ur-Rahman KM (2022) Study within a review (SWAR). *J Evid Based Med*; 15: 328-332 <u>https://doi.org/10.1111/jebm.12505</u>

### ESIC Stage 2: What do we need?

**Open consultation 12-19 March** 2025

### How you can get involved (3)

https://evidencesynthesis.atlassian.net/ wiki/spaces/ESE/overview

SUPPORTED BY

Evidence Synthesis Infrastructure Collaborative Summing up

- There are some great tools that may soon be ready for use
- But promise of GenAI will remain a promise until we have a good evidence base
- We need lots of rigorous evaluation before we can see the promise realized
- We need to increase our 'AI literacy' across the field to understand when and how to use (and not use) this new generation of tools

![](_page_67_Picture_5.jpeg)

![](_page_68_Picture_0.jpeg)

# 

### Thank you

### **James Thomas**

EPPI-Centre website: <u>http://eppi.ioe.ac.uk</u> Email <u>james.thomas@ucl.ac.uk</u> Twitter James\_M\_Thomas

### **EPPI-Centre**

Social Science Research Unit Institute of Education University of London 18 Woburn Square London WC1H 0NR

Tel +44 (0)20 7612 6397 Fax +44 (0)20 7612 6400 Email eppi@ioe.ac.uk Web eppi.ioe.ac.uk/

![](_page_68_Picture_8.jpeg)