



THE UNIVERSITY of EDINBURGH  
Institute for Neuroscience and  
Cardiovascular Research



# AI-enabled Systematic Review Platforms

## How do they perform?

Dr Charis Wong

Honorary Senior Clinical Lecturer, University of Edinburgh

Consultant Neurologist, NHS Fife

[charis.wong@ed.ac.uk](mailto:charis.wong@ed.ac.uk)

UHMLG Spring Forum | 23 April 2026

 [edin.ac/incr](https://edin.ac/incr)

## Disclosure

---

Nested Knowledge has provided us with a free pilot of their 'Enterprise' tier for the purpose of this evaluation.



# About me (and how I got interested in this...)

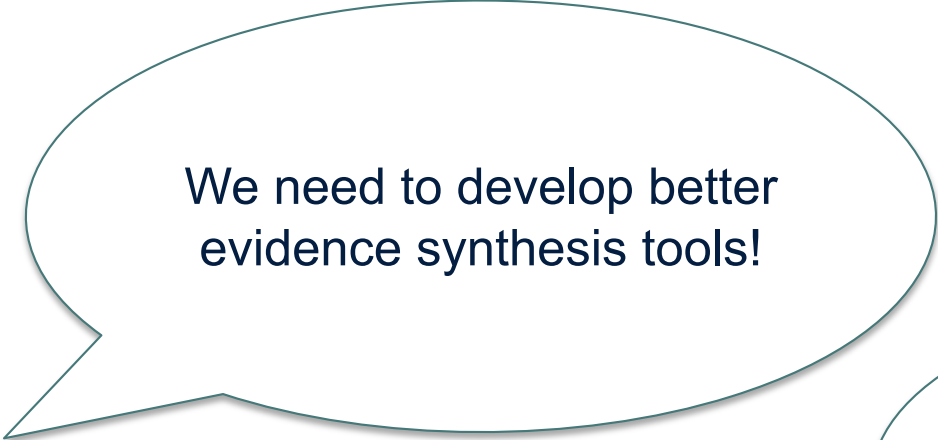


RESULTS BY YEAR

>20k  
publication/year  
for dementia on  
PubMed alone



THE UNIVERSITY of EDINBURGH  
Institute for Neuroscience and  
Cardiovascular Research



We need to develop better  
evidence synthesis tools!



Why don't we use  
[insert AI-SRP name]?

## Aim

---

- To evaluate the performance of three AI-enabled systematic review platforms in replicating three recently published systematic reviews
- Protocol pre-registered on OSF 1st July 2025 (<https://osf.io/8nf6h/files/qadcw>)



## Methods: SR selection

---

- Identify SR published in past year by
  - CAMARADES (non-human animal SR)
  - JBI
  - Cochrane
- Dual screen SRs against inclusion criteria – SR containing quantitative analysis
- Random selection of 1 SR per category from included SR



## Methods: SRP evaluation

---

We selected three Systematic Review Platforms (SRP) with embedded AI techniques for evaluation.

Elicit.AI  
Nested Knowledge  
Scispace

We considered discussions from other researchers/collaborators/evidence synthesis communities, available evaluation data, and publications citing platforms in selection of SRP

We perform evaluation of these platforms on each SR in a 3x3 design.

	SR1	SR2	SR3
SRP1	Evaluation 1	Evaluation 4	Evaluation 7
SRP2	Evaluation 2	Evaluation 5	Evaluation 8
SRP3	Evaluation 3	Evaluation 6	Evaluation 9



# What do they claim?

## Elicit.AI

### 📈 Scale

Elicit can find up to 1,000 relevant papers and analyze up to 20,000 data points at once.

### 🎯 Accuracy

Elicit is the most accurate AI product for scientific research. Learn about how we validated Elicit's accuracy.

### 🔍 Transparency

Elicit supports all AI-generated claims with sentence-level citations from the underlying sources.

### + More than chat

Elicit goes beyond chat to provide rich, interactive tables and multi-step workflows.

## Nested Knowledge

Powerful evidence synthesis tools for medical researchers. Accelerate, collaborate, automate and share.

Decisions at the speed of AI: Get answers in days, not months, with Nested Knowledge®, your all-in-one solution for efficient literature reviews.

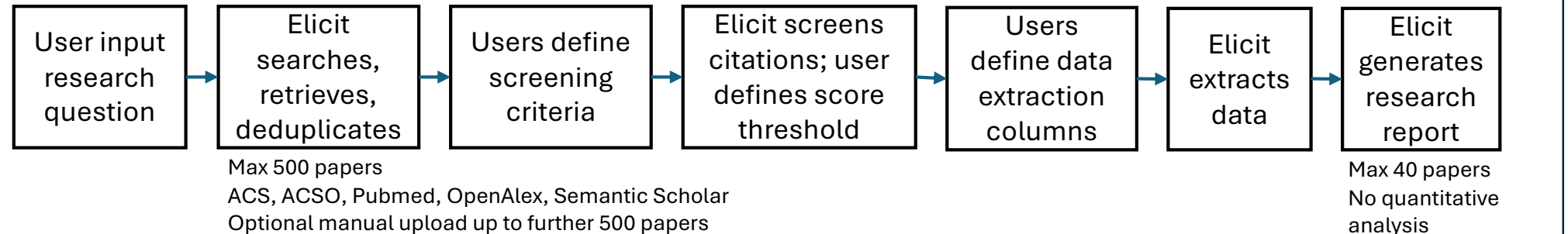
Start with AI, Dive Deeper with Experts

## Scispace

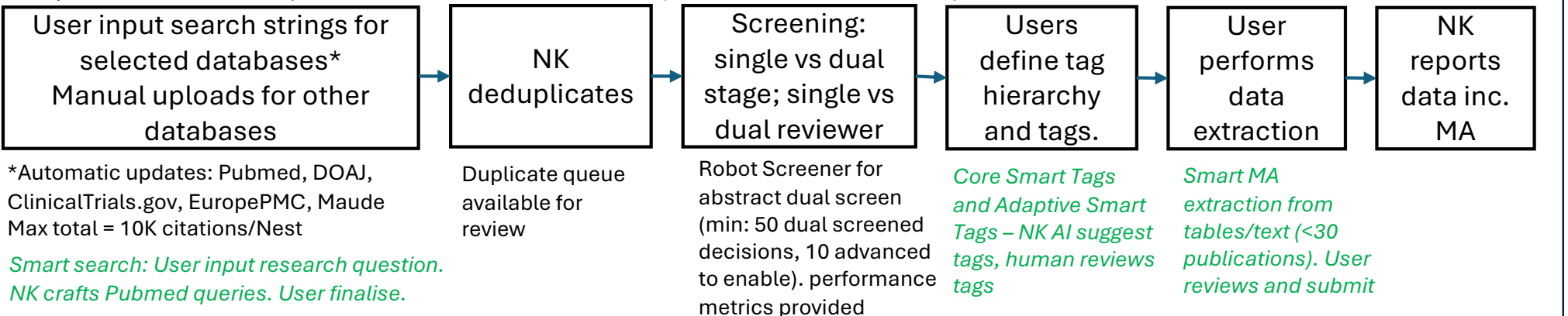
Feature	SciSpace Agent
Primary focus	Academic research workflows (writing, reviewing, data analysis, etc.)
Designed for	Researchers, PhD students, R&D professionals, and other academicians
Knowledge base	280M+ papers, 50M+ full-text PDFs
Use case coverage	Manuscript writing, peer review, grant proposal drafting, and more!
Accuracy wrt academic use cases	High (Deep integration with scholarly data)
Integration with academic databases	Integrated with 150+ academic tools & databases

## Platform features on versions available in July 2025

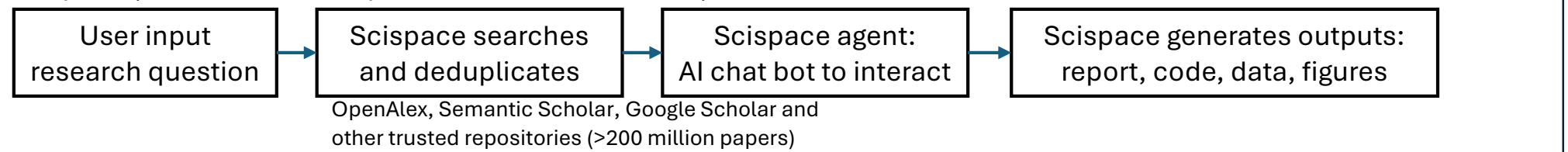
### Elicit (Team subscription; USD 65/month/user)



### NK (Academic subscription ~£18/month/user; *Enterprise \$695/month/user*)



### Scispace (Advanced subscription; USD 90/month/user)



# Data export issues with Scispace

https://scispace.com/chat/54f7d970-e0b5-431d-ad87-90851d4b3b1b

The screenshot shows the Scispace web interface. The browser address bar displays the URL: `scispace.com/chat/54f7d970-e0b5-431d-ad87-90851d4b3b1b`. The page title is "Effects Of TAAR1 Agonists".

The interface is divided into three main sections:

- Left Sidebar:** Contains navigation options such as "Home", "My Library", "Agent Gallery", "AI Writer", "Chat with PDF", "Literature Review", "Find Topics", "Paraphraser", "Citation Generator", "Extract Data", and "AI Detector". It also lists "Recent Chats" with "Effects Of TAAR1 Agonists" selected.
- Chat Area:** Shows the chat history with the text: "supporting TAAR1 agonist development as next-generation antipsychotic therapy." Below this, four files are listed: "main.pdf", "enhanced\_systematic\_review.", "final\_included\_studies.csv", and "inclusion\_exclusion\_criteria.m". A "Follow up suggestions" section offers tasks like "Conduct network meta-analysis comparing different TAAR1 agonist compounds", "Search for TAAR1 agonist clinical trial data", "Create interactive dashboard displaying meta-analysis results", "Generate publication-ready figures for journal submission", and "Develop PROSPERO protocol for living systematic review".
- Output Panel:** Displays an error message: "Something went wrong" with a red exclamation mark icon. Below the message, it states: "We encountered an error while loading the generated file. Please try refreshing the page." There are two buttons: "Try Again" and "Refresh Page".

At the bottom of the chat area, there is a text input field with the placeholder "Ask anything or give follow up task..." and a "Deep Search" button. The user's profile "aidesteval" is visible in the bottom left corner.

## Data export issues with Scispace

PMID	Title	Abstract	Journal	Authors	Publication	DOI
22641180	A growing understanding of the role of muscarinic receptors in the r	Pre-clinical models, postmortem and neuroimaging studies				
22641180	A new perspective for schizophrenia: TAAR1 agonists reveal antipsyc	Schizophre	Molecular	Revel, Mori	2013	10.1038/mg
22641180	A new perspective for schizophrenia: TAAR1 agonists reveal antipsyc	hotic-and antidepressant-like activity, improve cognition an				
30366004	Activation of trace amine-associated receptor 1 attenuates schedule	Trace Amin	Neurophar	Sukhanov,	2019	10.1016/j.n
26441502	An animal model of differential genetic risk for methamphetamine in	The questic	Frontiers in	Phillips, Sh	2015	10.3389/fr
26441502	Avenues for the development of therapeutics that target trace amine associated receptor 1 (TAAR1) miniperspective					
29520239	Behavioral Phenotyping of Dopamine Transporter Knockout Rats: Cr	Alterations	Frontiers in	Cinque, Zoi	2018	10.3389/fr
29520239	Beyond dopamine receptor antagonism: new targets for schizophre	Treatment of schizophrenia (SCZ) historically relies on the us				
29520239	Binding of SEP-363856 within TAAR1 and the 5HT1A receptor: impli	Current medications for schizophrenia typically modulate d				
22763617	Brain-specific overexpression of trace amine-associated receptor 1 a	Trace amin	Neuropsych	Revel, Mew	2012	10.1038/n

## Data export issues with Scispace

PMID	Title	Abstract	Journal	Authors
22641180	A new perspective for schizophrenia: TAAR1 agonists reveal antipsychotic- and antidepressant-like activity, improve cognition and control	Schizophrenia	Molecular psychiatry	Revel, Morea
	Binding of SEP-363856 within TAAR1 and the 5HT1A receptor: implications for the design of novel antipsychotics	Current medications for schizophrenia typically act on dopamine and serotonin receptors		
	Effect of co-treatment of olanzapine with SEP-363856 in mice models of schizophrenia	Olanzapine is a commonly used drug in the treatment of schizophrenia		
	In-vivo pharmacology of trace-amine associated receptor 1			
32382134	Reproducing the dopamine pathophysiology of schizophrenia and approaches to ameliorate it with the novel TAAR1 agonist SEP-363856, a novel psychotropic agent with a unique, non-D2 receptor mechanism of action	Patients with schizophrenia have a dopamine dysregulation hypothesis	Molecular psychiatry	Kokkinou, Irwin
36100653	TAAR1 dependent and independent actions of the potential antipsychotic and dual TAAR1/5-HT1A receptor agonist SEP-363856	For the past 50 years, the clinical efficacy of antipsychotics has been limited by their side effects	Neuropsychopharmacology	Saarinen, Marjaana
	The potential of trace amines and their receptors for treating neurological and psychiatric disorders	Mining of the human genome has revealed a large number of novel receptors		
	The trace amine-associated receptor 1 modulates methamphetamine's neurochemical and behavioral effects	The newly discovered trace amine-associated receptor 1 (TAAR1) is a potential target for the development of new antipsychotics		
	Trace amine associated receptor 1 modulates behavioral effects of ethanol	Background: Few treatment options for alcohol use disorder (AUD) are available		
	Trace amine-associated receptor 1 (TAAR1) agonism as a new treatment strategy for schizophrenia and related disorders			
	Trace amine-associated receptor 1 as a target for the development of new antipsychotics: characterization of the novel TAAR1 agonist SEP-363856	Schizophrenia is a mental illness associated with dopamine dysregulation		
	Trace amine-associated receptor 1 contributes to diverse functional actions of O-Phenyl-lodotyramine in mice but not to the effects of amphetamine			
	Trace amine-associated receptor 1 modulation of dopamine system	Trace amine-associated receptor 1 (TAAR1) is a potential target for the development of new antipsychotics		
	Trace amine-associated receptor 1 partial agonism reveals novel paradigm for neuropsychiatric therapeutics			
	Trace amine-associated receptor 1: a multimodal therapeutic target for neuropsychiatric disorders	Introduction: The trace amines, endogenous neurotransmitters, are involved in a wide range of physiological functions		
	Unlocking the therapeutic potential of ulotaront as a trace amine-associated receptor 1 agonist	All antipsychotics currently used in clinical practice act on dopamine and serotonin receptors		

# Search and screening results for SR1

non human SR in “Trace amine-associated receptor 1 (TAAR1) agonism for psychosis: a living systematic review and meta-analysis of human and non-human data” <https://wellcomeopenresearch.org/articles/9-182> 15 FT included in original review

Original review - FT included	SRP decision	Elicit	NK (full text dual-human screened; Robot Screener used for TiAb dual screen)	NK Robot screener (TiAb)	Scispace
In	In	6	15	5	6
In	Out	9	0	2	9
Out	In	21	6	14	11
Sensitivity		40%	100%	71%	40%
Precision		22%	94%	26%	35%

## Nested Knowledge Robot Screener Internal Cross Validation Metrics for SR1

History	Records (adv)	Recall	Precision	F1	Accuracy	AUC
2025-08-08	1889 (29)	0.60	0.33	0.39	0.97	0.92
2025-08-04	1889 (29)	0.59	0.28	0.37	0.97	0.93
2025-08-03	1859 (24)	0.89	0.38	0.52	0.98	0.99
2025-08-03	1769 (24)	0.75	0.35	0.48	0.98	0.98
2025-08-03	1017 (22)	0.72	0.43	0.53	0.97	0.96
2025-08-03	990 (21)	0.74	0.43	0.53	0.97	0.99
2025-08-01	477 (11)	0.72	0.38	0.45	0.96	0.92
2025-08-01	477 (11)	0.61	0.28	0.32	0.96	0.87

View Recommendations: [Advance](#) [Exclude](#) [Delete Model](#)

### How the Screening Model Works

At a high level, the model is a Decision Tree- a series of Yes/No questions about characteristics of records that lead to different probabilities of inclusion/advancement.

In more detail, the model is a gradient-boosted decision tree ensemble. Its hyperparameters, particularly around model complexity (number of trees, tree depth) are optimized using a cross validation grid search. The model produces posterior probabilities and is optimized on logistic loss. SMOTE oversampling is employed as a correction to highly imbalanced classes frequently seen in screening.

# SR1: Additional inclusions

Publications before original search performed	Appropriate inclusion	Elicit	NK	Scispace
TRUE	Yes	2 <sup>a,b</sup>	1 <sup>b</sup>	0
	No	15	0	11
FALSE	Yes	1	5	0
	No	3	0	0

(a) 1 paper was identified within the TAAR1 agonist review update 2024 via psychosis-SOLES. (b) 1 paper was identified in the original search but excluded following human reviewer dual screening – this paper was included following further review by senior reviewer prompted by this evaluation.

# Search and screening results for SR2

JBI review - “Comparison of diagnostic accuracy of rapid antigen tests for COVID-19 compared to the viral genetic test in adults: a systematic review and meta-analysis. DOI: [10.11124/JBIES-23-00291](https://doi.org/10.11124/JBIES-23-00291)  
 Original review screening stages – TiAb / full text / critical appraisal. 142 publications included in analysis

Original review - FT included	SRP decision	Elicit	NK (full text dual-human screened; Robot Screener used for TiAb dual screen)	NK Robot screener (TiAb)	NK Robot screener against original TiAb
In	In	58	99	111	323
In	Out	84	43	37	256
Out	In	367	85	549	337
Sensitivity		40%	70%	75%	56%
Precision		14%	54%	17%	49%

# SR2-Elicit

JBIR review - “Comparison of diagnostic accuracy of rapid antigen tests for COVID-19 compared to the viral genetic test in adults: a systematic review and meta-analysis. DOI: [10.11124/JBIES-23-00291](https://doi.org/10.11124/JBIES-23-00291)

Included publications by Elicit not within original review		
Publications before original search performed (Last update 12 July 2022)	Appropriate inclusion	N
TRUE	Yes	50
	No	116
FALSE	Yes	64
	No	23

Multiple duplicates within Elicit search (5 using DOI alone; 23 further if matching preprint to published DOI); Duplicates (including preprint with linked publications) have been removed for this analysis.

# Annotations

## *Percentage of tags in original review accurately tagged by SRP*

<b>Evaluation</b>	<b>Population</b>	<b>Intervention</b>	<b>Comparison</b>	<b>Outcomes</b>	<b>ROB</b>	<b>Overall accuracy</b>
<b>SR1-Elicit</b>	60.0%	55.6%	66.7%	NA	51.7%	54.5%
<b>SR1-NK</b>	53.3%	33.3%	66.7%	20.0%	NA	43.3%
<b>SR2-Elicit</b>	88.9%	85.8%	61.7%	72.0%	55.1%	70.3%

## *Percentage of publications given additional tags by SRP (percentage appropriate in parenthesis)*

<b>SR1-Elicit</b>	80% (0%)	55.6% (20%)	66.7% (0%)	NA (NA)	0% (0%)
<b>SR1-NK</b>	26.7% (75%)	60% (11%)	25% (33%)	73.3% (45%)	NA
<b>SR2-Elicit</b>	3.9% (0%)	5.8% (0%)	0% (0%)	12.1% (61%)	44.8% (6%)

Challenges in tagging pharmacologically induced models (assigning intervention rather than population tag); tagging unrelated terms in rest of text.

Many of additional outcomes tagged by SRPs were not primary outcomes of original SR and hence were appropriately not tagged there.

Multiple tags per category and hierarchical tagging likely to affect planned calculations of kappa agreement.

## Discussion: AI-enabled systematic review platforms

---

- Variation in
  - Performance in SR tasks – sensitivity an issue currently across platforms relying on AI methods; annotation/extraction variable; analysis limited
  - Transparency
  - Explainability
  - Human supervision / Customisability
  - Hallucinations
  - Interoperability / export of code/data/metadata
  - Scalability
  - Costs
- Generalisability: need more studies



## Discussion: Evaluating systematic review platforms

---

- Challenging!
- Challenges with replicating SRs
  - Data and metadata FAIRness can affect replication (e.g. lack of bibliographic data of excluded studies; incomplete/inaccurate metadata; systematic search strategies / results / bibliographic data in PDF format)
  - SR selection matters: Consider size, scope, how well questions are defined vs feasibility / resources
- Challenges with evaluating SRPs
  - Variation in extent of human effort needed to perform SR in SRPs
  - Lack of gold standards on data, metadata / code, format leads to significant variation in SRP outputs; cleaning/extracting data for evaluation of SRP can be time and resource intensive
  - Results limited to version of SRP evaluated; quickly superseded



## Discussion: Looking ahead (1)

---

- Responsible use of AI in evidence synthesis crucial. See RAISE guidance.
- Each SR provides data and opportunity to validate AI digital evidence synthesis tools both retrospectively and prospectively – we need to maximise this.
  - increasing FAIRness of data/metadata at each step of the SR to improve replicability:
  - Roles of individuals, groups, community, institutions
  - Systems/infrastructure to maximise this – using ‘routinely collected’ data in SRs for evaluation of AI evidence synthesis tools
- Methods development in evaluating AI-enabled tools and platforms



## Further questions

---

- Should SRPs be regulated? if so, who by? (*possible trade off with innovation?*)
- Who should be responsible for evaluation of tools and platforms? Cf. onus on drug developers/pharma to demonstrate evidence of efficacy/safety for regulatory submissions
- Ethical, legal and environmental implications of use of these platforms?



## Disclaimers

THE SERVICE AND ALL RELATED MATERIALS AND CONTENT ARE PROVIDED “AS-IS” AND “AS AVAILABLE,” AND WE DISCLAIM ANY AND ALL WARRANTIES, WHETHER EXPRESS OR IMPLIED, INCLUDING WITHOUT LIMITATION IMPLIED WARRANTIES OF TITLE, MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE, AND NON-INFRINGEMENT. WE CANNOT GUARANTEE AND DO NOT PROMISE AVAILABILITY OF THE SERVICE OR ANY SPECIFIC RESULTS FROM USE OF THE SERVICE.

WE ARE NOT RESPONSIBLE FOR ANY OF YOUR CONTENT OR ANY OTHER USER’S CONTENT, INCLUDING CONTENT THAT VIOLATES THESE TERMS OF SERVICE OR THAT IS INCORRECT, INCOMPLETE, INACCURATE, OR OFFENSIVE, WHETHER SUCH CONDITION IS CAUSED BY USERS OF THE SERVICE, MEMBERS, OR BY ANY OF THE SOFTWARE OR INFRASTRUCTURE ASSOCIATED WITH OR USED TO PROVIDE THE SERVICE. WE ARE NOT RESPONSIBLE FOR ANY THIRD PARTY’S SERVICES OR PERFORMANCE, INCLUDING ANY THIRD PARTY PROVIDING SERVICES, DATA, OR INFRASTRUCTURE FOR THE SERVICE.

UNDER NO CIRCUMSTANCES WILL WE BE RESPONSIBLE FOR ANY LOSS OR DAMAGE, INCLUDING PERSONAL INJURY OR DEATH, RESULTING FROM ANYONE’S USE OF THE SERVICE, ANY CONTENT POSTED ON THE SERVICE OR TRANSMITTED TO OTHER USERS, OR ANY INTERACTIONS BETWEEN USERS OF THE SERVICE, WHETHER ONLINE OR OFFLINE.



THE UNIVERSITY of EDINBURGH  
Institute for Neuroscience and  
Cardiovascular Research



## Acknowledgements

Harris Khalid, Kat Stefanska, Florence Do, Mohammed Abomostafa, Gemma Dalgety, Malcolm Macleod

Screeners for Nested Knowledge review:

Leonor Rodriguez, Matheus Gallas-Lopez, Stephen Malden, Kaitlyn Hair, Gillian Currie, Yvonne Khor, Nadia Soliman, Khrystyna Stoska

✉ [charis.wong@ed.ac.uk](mailto:charis.wong@ed.ac.uk)