

# Pharmacoevidence

Improving patient lives by supporting evidence base for healthcare interventions



## Gen AI in Systematic Literature Reviews: First Case Study on the use of Gen AI in a NICE Submission

ISPOR EU 2025 | Glasgow

### Moderator

- Barinder Singh, Pharmacoevidence

### Speakers:

- Sven Klijn, Bristol Myers Squibb
- Dilip Makhija, Gilead Sciences Inc.
- Raphael Sonabend-Friend, National Institute for Health and Care Excellence, UK



Pharmaco<sup>®</sup>  
**Evidence**

Mohali India | London UK | Canada





## MODERATOR



**Barinder Singh**  
Pharmacoevidence

Introduction

## PANELISTS



**Sven Klijn**  
Bristol Myers Squibb

AI-augmented evidence  
generation



**Dilip Makhija**  
Gilead Sciences Inc.

Real-world experience with  
automated SLRs



**Dr. Raphael Sonabend-Friend**  
NICE, UK

HTA perspective on AI-  
driven evidence synthesis

# SLRs lie at the top of the evidence hierarchy



## What

- A structured and transparent process to identify, assess, and synthesize existing evidence on a specific research question
- Follows a transparent and replicable process
- Key features:
  - Predefined research question
  - Inclusion/exclusion criteria
  - Quality assessment of studies
  - Clear synthesis and reporting



## Why

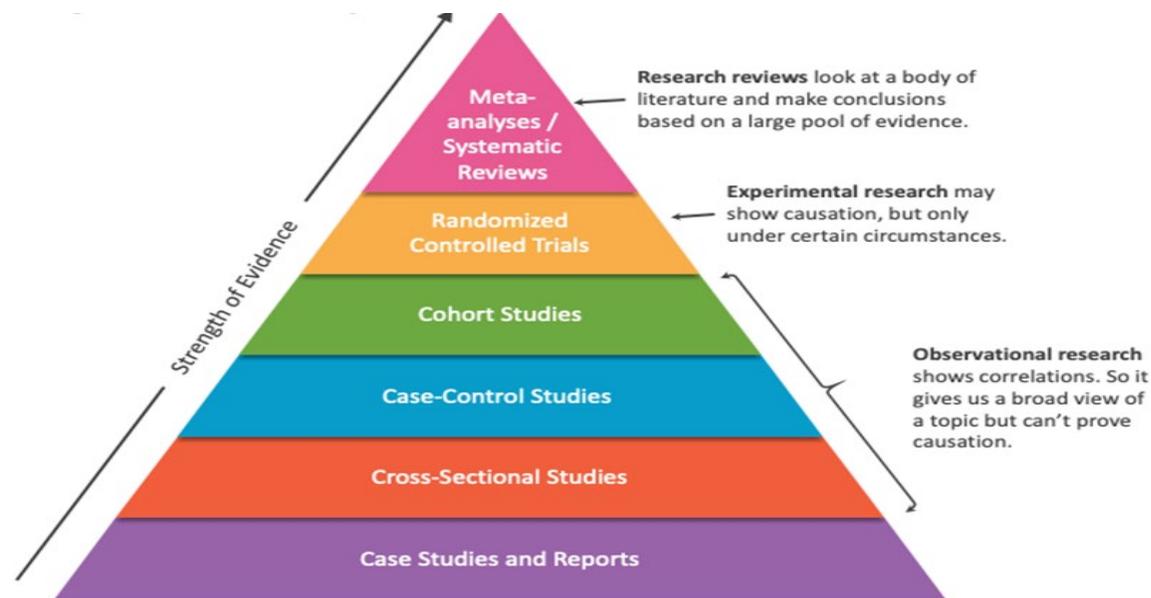
- To provide a **comprehensive** and **unbiased** summary of available data to guide **clinical, regulatory, or research decision-making**
- Compare treatment efficacy and safety across multiple studies
- Synthesize high-quality data to **strengthen scientific conclusions**
- Facilitate meta-analyses by providing **structured data for statistical pooling**



## How

- Define a focused research question using structured frameworks (e.g., PICO)
- Develop a review protocol with clear inclusion/exclusion criteria
- Search and screen literature across multiple databases
- Assess study quality using standardized tools
- Extract and synthesize data through narrative or statistical analysis

# SLRs lie at the top of the evidence hierarchy



## SLRs are The Evidence Foundation for HTA

Health Technology Assessment (HTA) relies on SLRs to ensure that all clinical and economic evidence is identified, appraised, and synthesized transparently and reproducibly

- HTA agencies (e.g., NICE, CADTH, EUnetHTA/JCA) require SLR-based evidence to support comparative effectiveness and cost-effectiveness analyses
- SLRs provide a transparent, auditable trail of how evidence was sourced and selected
- They minimize bias and duplication, ensuring that payer and regulatory decisions are grounded in complete and reliable evidence
- In the upcoming EU Joint Clinical Assessment (JCA), SLRs form the core evidence package for all new therapies

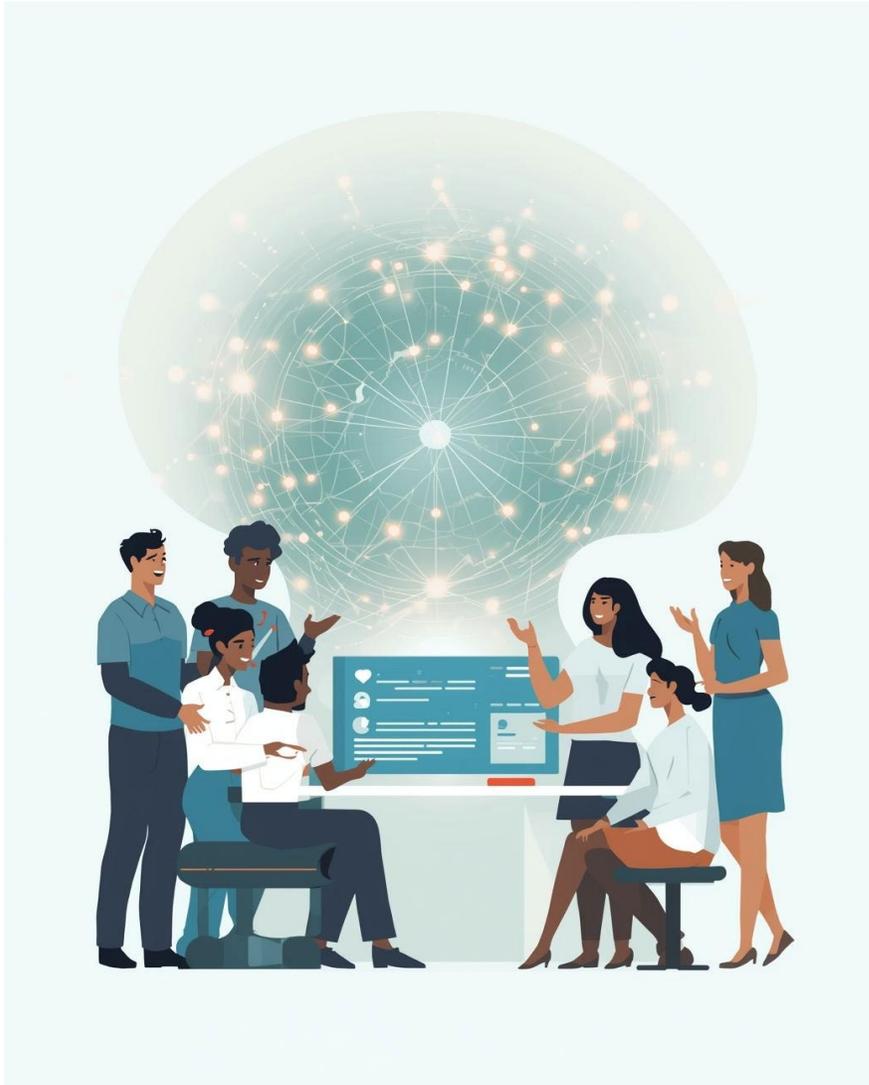
# More Data. More Complexity.



- Healthcare generates nearly 30% of the world's data<sup>1</sup>
- The volume of healthcare data is projected to grow at a 36% Compound Annual Growth Rate (CAGR) by 2025 — faster than any other industry (Arcadia-2023)<sup>1,2</sup>



- “Is it possible to deliver *faster, scalable* outputs without compromising on methodological rigor?”
- “Can generative AI truly be trusted to conduct systematic literature reviews for HTA submissions?”



- HTA bodies now recognize GenAI as an emerging contributor to evidence generation and evaluation
- However, concerns remain around the appropriateness, transparency, and reliability of AI

It is important to consider the use of AI methods carefully to ensure the anticipated benefits are balanced against the known concerns<sup>1</sup>



## NICE UK position statement

- ✓ Engage early with NICE and inform them of your plans to use AI
- ✓ Apply AI only when it adds demonstrable value
- ✓ Keep human reviewers in the loop; AI assists, never replaces
- ✓ Disclose algorithms, data sources, and validation tools.
- ✓ Perform and document validation, bias, and security tests.
- ✓ Comply with UK AI rules, licenses, and IP requirements

## CDA position statement

- ✓ Engage early with CDA-AMC to discuss plans for using AI
- ✓ All uses of AI should align with Canada's Voluntary Code of Conduct on the Responsible Development and Management of Advanced Generative AI Systems
- ✓ Principle of augmentation, not replacement of humans
- ✓ Ensure compliance with licensing agreements, maintain explainability of AI methods, and promote transparency in applications
- ✓ Ensure that all algorithms, models, datasets, and data processing pipelines used are in alignment with relevant Canadian regulations and guidance for AI use
- ✓ Any potential risks should be mitigated by adhering to established guidance and checklists
- ✓ Must adhere to the ethical principles for the use of AI in health



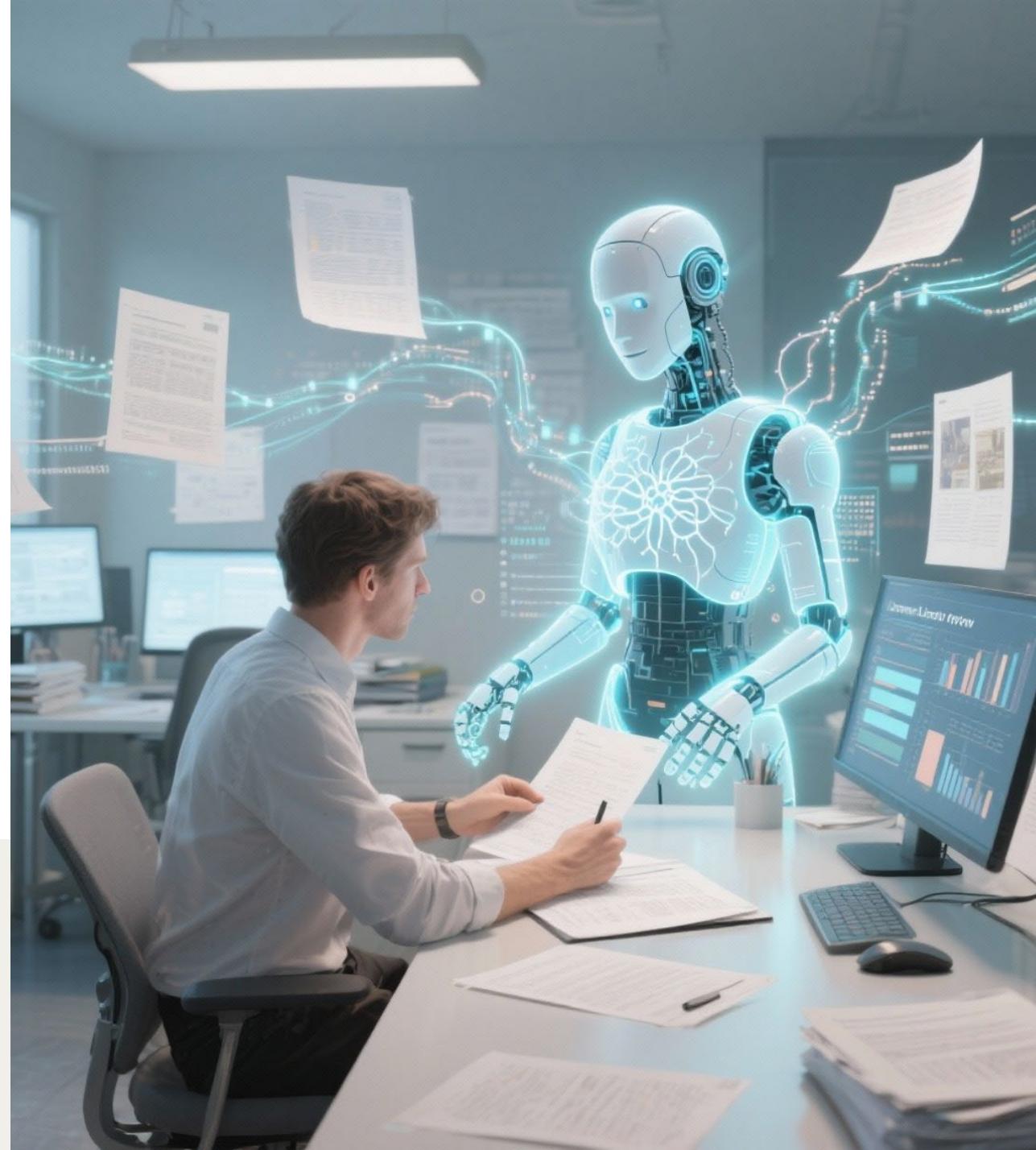
- "Have we reached a point where AI-assisted SLRs could stand up to regulatory scrutiny in real-world HTA submissions?"
- "Are there any real-world examples of Gen AI-assisted SLRs in HTA submissions?"

# GenAI in SLRs: From Dual Review to Human-AI Collaboration

November 11, 2025

ISPOR Glasgow

Sven Klijn





The opinions expressed in this presentation and on the following slides are solely those of the presenter and do not necessarily represent opinions of Bristol Myers Squibb.

# Dual human review

## Human 1

Formulating the research question



PICO definition



Search query



Title/abstract screening



Full text screening



Data extraction



Reporting

## Human 2

Title/abstract screening



Full text screening



Data extraction

## Human 3

Arbitration

Arbitration

Arbitration

# The gold standard is not perfect



Time-  
consuming



Limited  
scalability



Screening  
fatigue

**“The total error rate**  
(false inclusion and false exclusion)  
**was 10.76% (95% CI: 7.43% to 14.09%)”**

Wang et al. (2020). Error rates of human reviewers during abstract screening in systematic reviews. DOI: 10.1371/journal.pone.0227742

# Human-AI Collaboration

Human 1

Formulating the research question



PICO definition



Search query



Title/abstract screening



Full text screening



Data extraction



Reporting

AI

Title/abstract screening



Full text screening



Data extraction

Human 2

Arbitration

Arbitration

Arbitration



# AI Governance



## Model Documentation

- Architecture
- Training data
- Performance metrics
- Logging



## Life Cycle Maintenance

- Version control
- Bias & fairness assessments:  
Systematic evaluation of outputs
- Reproducibility

### Best practice principle 6:

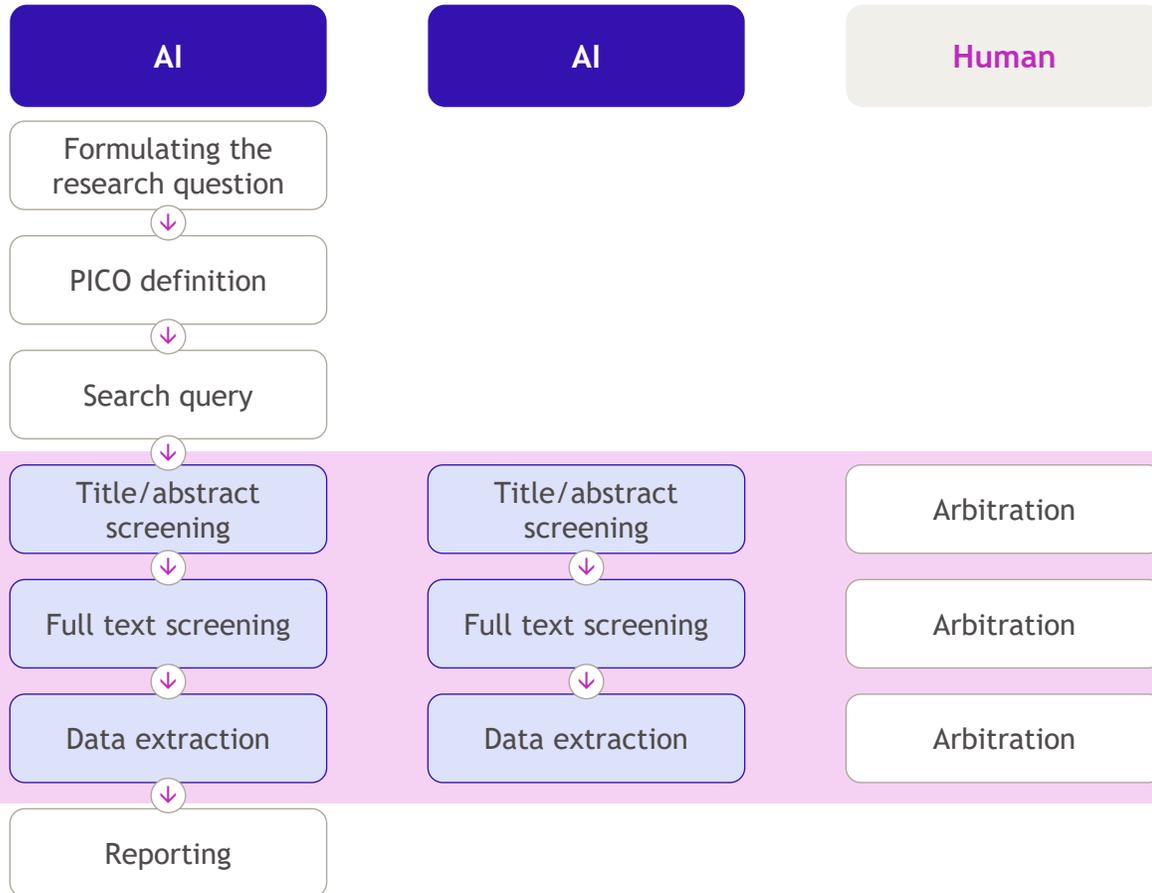
*“Ensure GenAI-assisted analyses are reproducible and auditable where possible.”*

NICE (2025). Generative Artificial Intelligence (GenAI) best practice principles

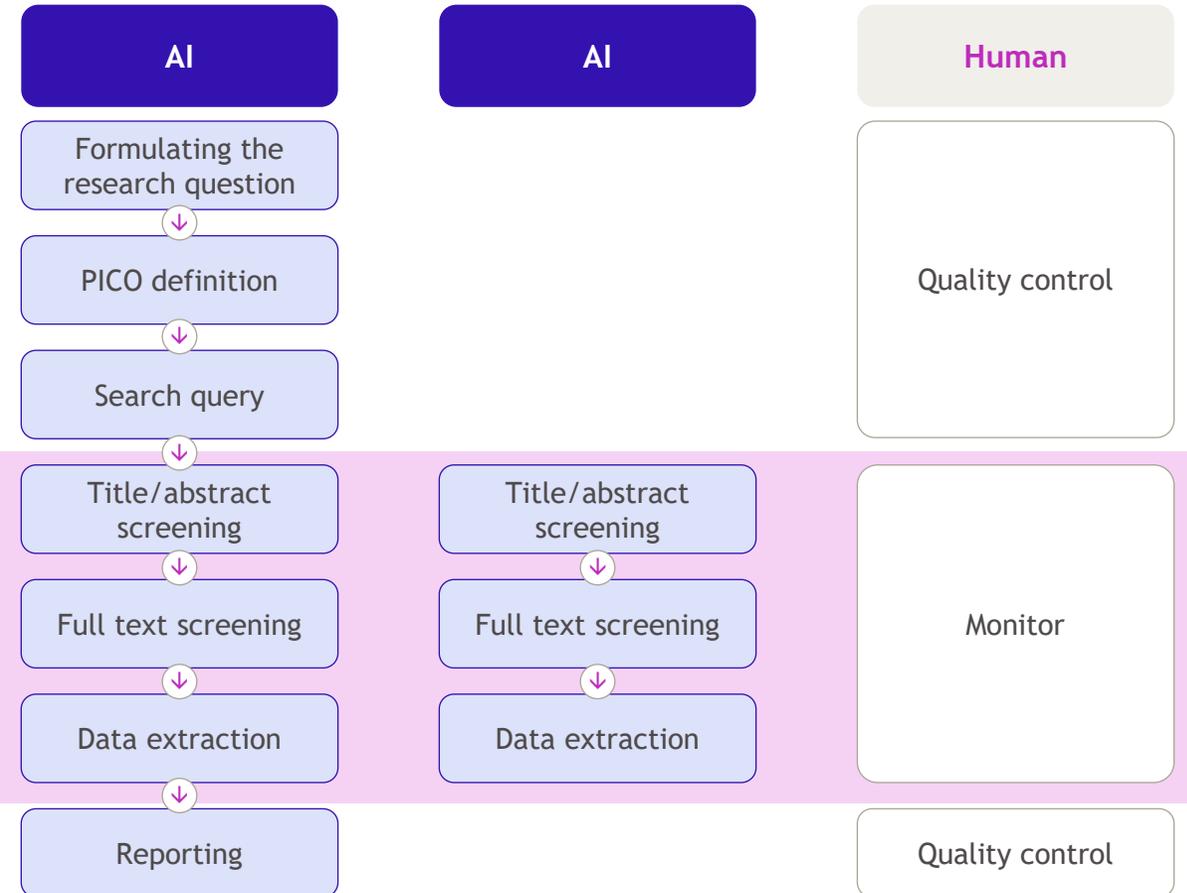


# Alternative models

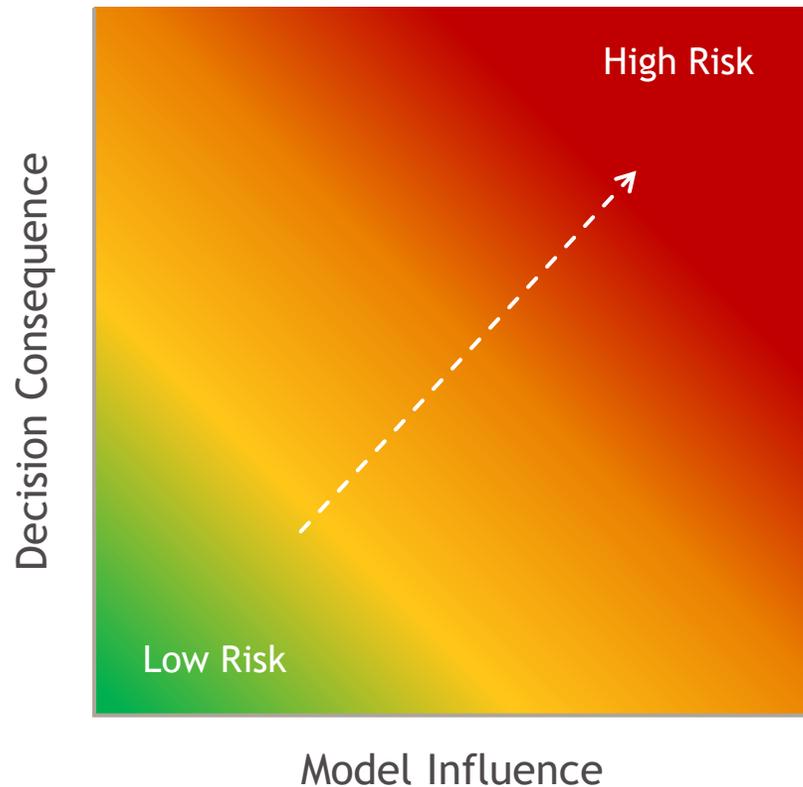
## Human Arbitration



## Human Monitoring



# Model risk matrix



## Decision consequence depends on:

- Type of publication
- Type of error made:  
false inclusion vs. false exclusion



## Model influence depends on:

- Model used:  
Human-AI / Human Arbitration / Human Monitoring
- Guardrails and mitigation strategies



# Reflection

---



Discussion required on  
where the bar should be set

---



Strong model governance is critical

---



Role of the human is essential,  
but differs by implementation model

# Implementation of AI tools: A real-world example for the use of an automated SLR in HTA submissions

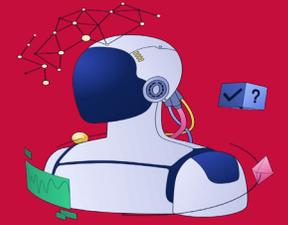
Dilip Makhija | ISPOR EU Glasgow Panel 2025

# Disclaimer

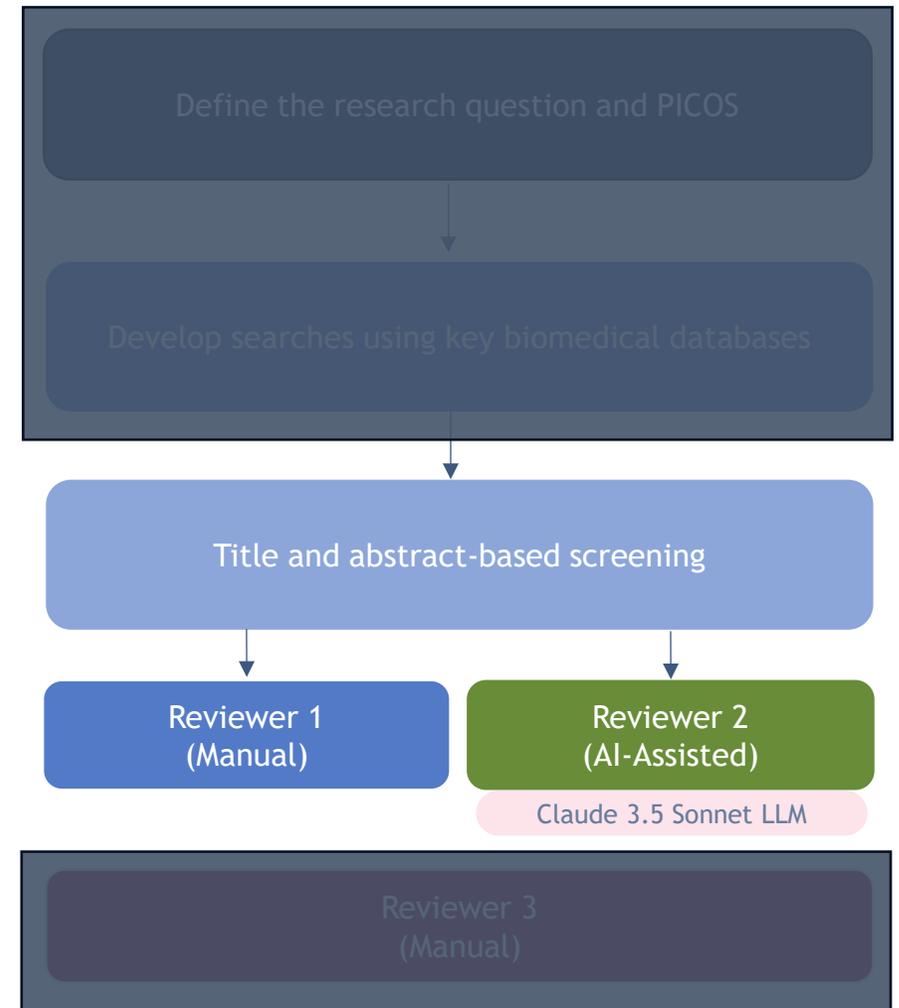
The opinions expressed on the following slides are solely those of the presenter and not those of Gilead Sciences, Inc.



# An AI tool was employed as a second reviewer to manage the large volume of citations retrieved from the literature search, helping to adhere to the tight timelines for the UK NICE submission



- AI was implemented as a complementary second reviewer, in alignment with NICE and CDA position statement, and the Cochrane guidance, i.e., human-in-the-loop process
- The approach adhered to NICE methodological standards as outlined in the technical support document (TSDs) on identifying, selecting, and synthesizing clinical evidence with a two-reviewer screening and robust quality control to reduce bias and errors



# Integration of the AI tool helped meet ambitious timelines and save ~50% time



## Rationale and value of using Gen AI for SLR

- ✓ The rationale for integrating the AI tool was driven by the **substantial volume of citations** retrieved during literature searches, **presenting significant timeline challenges for submission**
- ✓ AI-based screening **reduced the time** required for the Title/Abstract-based screening by **~2 weeks**

Published evidence on the effectiveness of the tool

Level of agreement between Human and AI



# Over 10 published studies that had previously employed the automated tool were provided to support its validation and justification for use

Rationale and value of using Gen AI for SLR



**Published evidence on the effectiveness of the tool and Early engagement with HTA**

Level of agreement between Human and AI

- Early engagement with NICE regarding the use of AI is ideal, but we couldn't do that this time because of tight submission timelines
- We submitted comprehensive evidence supporting the tool's effectiveness, including over 10 published studies demonstrating its accuracy, specificity, sensitivity, and reviewer agreement rates

# Human-AI agreement was high, and the AI tool did not exclude any relevant citation

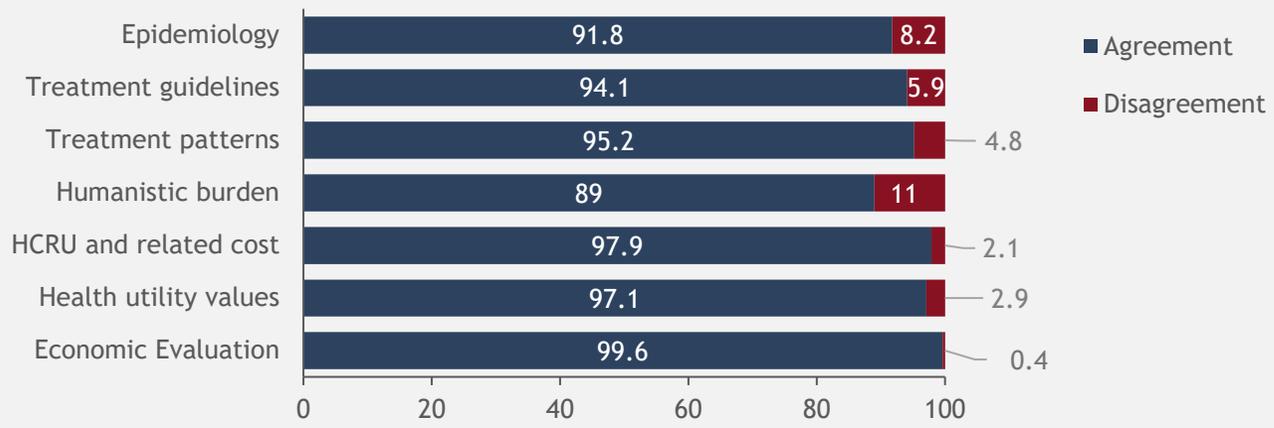
Rationale and value of using Gen AI for SLR

Published evidence on the effectiveness of the tool

- ✓ In standard practice for SLRs, an **ideal agreement** between two human reviewers is **approximately 95%**. However, **real-world experience** typically demonstrates around or **below 90% agreement**, mainly due to reviewer expertise and judgment variations
- ✓ In our AI-assisted review process, the overall average agreement between AI and human reviewers was **94.97% (median: 95.24%)**, ranging from 89.01% (Humanistic Burden Review) to 99.59% (Economic Evaluation Review)

  
**Level of agreement between Human and AI**

Decision alignment between human and AI



Did AI exclude any relevant citations? **No, it didn't**



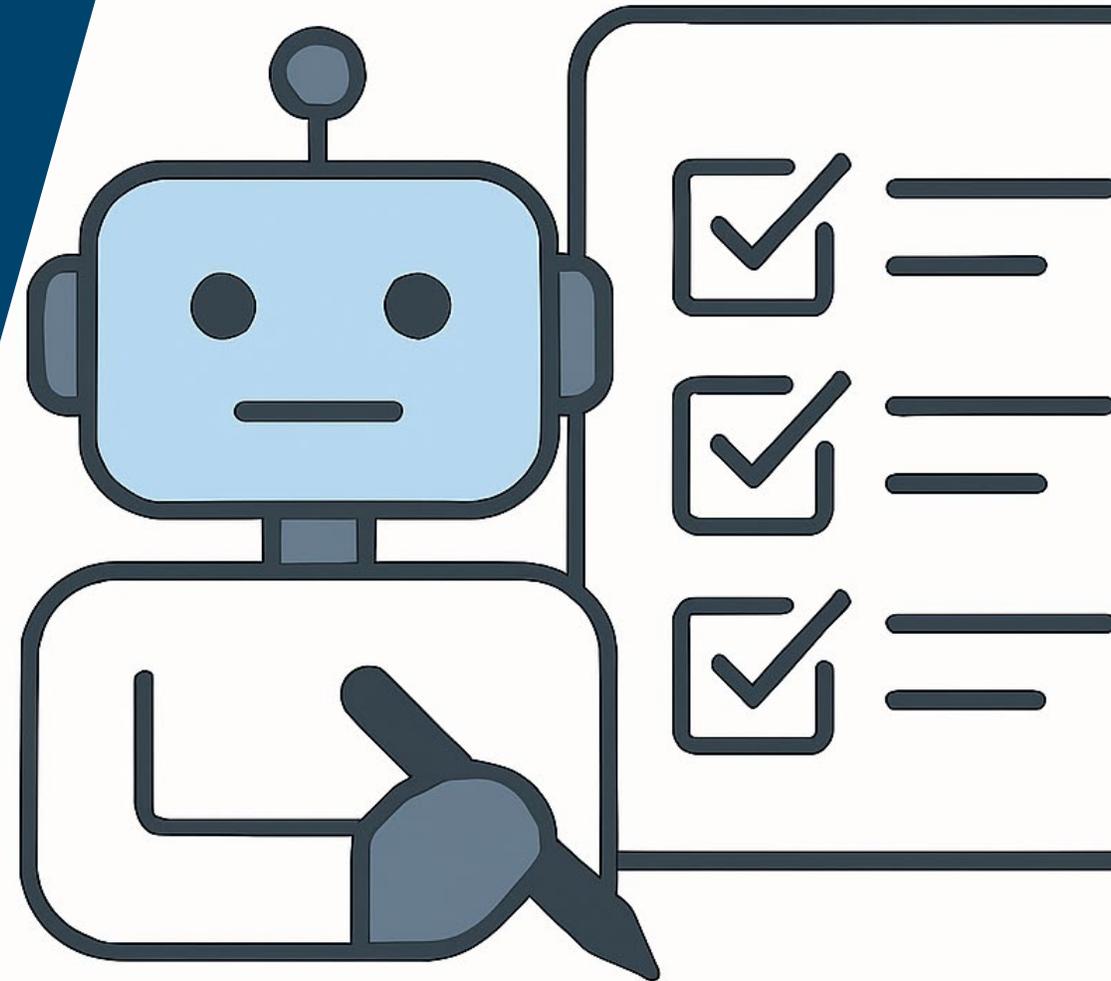
# NICE was in agreement and aligned on the use of the AI-assisted tool

- ✓ The AI-assisted reviewer functioned as a complementary second reviewer
- ✓ Any disagreements between the human and AI reviewers (e.g., one includes, one excludes) were resolved by an independent QC human reviewer
- ✓ To ensure robust quality assurance, all excluded citations, even those where both reviewers agreed, were rechecked by a separate QA reviewer to confirm no relevant studies were missed
- ✓ AI as second reviewer and layered QC/QA framework ensured accuracy, transparency, and methodological rigor, maintaining alignment with NICE's "human-in-the-loop" guidance and confirming the reliability of AI-assisted screening
- ✓ This submission established an important benchmark for the integration of Generative AI in evidence generation and represents the first AI-assisted SLR accepted by NICE (Singh 2025<sup>1</sup>)



# AI in systematic reviews

Raphael Sonabend-Friend  
Scientific Adviser



# AI-assisted citation screening accepted in a NICE submission [ID6429]

1

Pharmacoevidence AI tool used as a second reviewer for screening citations with an independent human expert checking discrepancies.

2

EAG asked for the rationale behind AI integration and for validation of the AI's effectiveness.

3

"Given the QA processes described by the company, the EAG believed that the use of the AI/ML tool was appropriate". [\[ID6429\]](#)

# AI Statement of Intent and Position Statement



Statement  
of Intent

Intention to develop our approach to: providing guidance on best practice development; evaluation; internal use.



Evaluation

In-house experimentation; living systematic reviews; external consultations (HTA Lab).



Position  
Statement

**Living** document outlining NICE's position on use of AI in generation and reporting of evidence.

# AI updates to NICE submission templates

1

Has AI been used in the generation of data or research presented in this submission?

2

Has AI been used in the reporting or synthesis of research presented in this evidence submission?

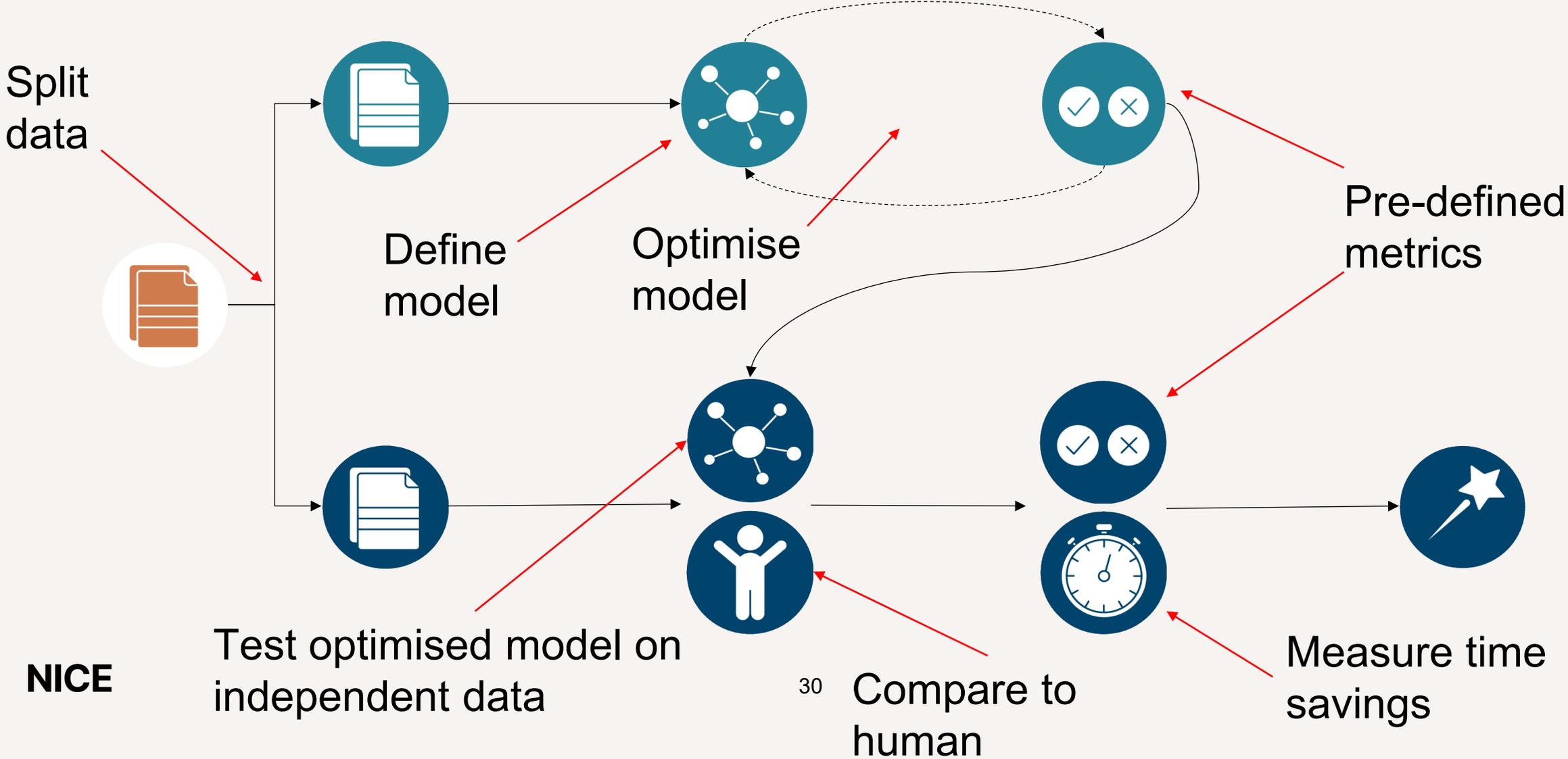
3

Please indicate where AI was used, any potential risks with its use and what steps were taken to address these risks.

4

Please indicate the level of human oversight for these uses.

# Best practice evaluation

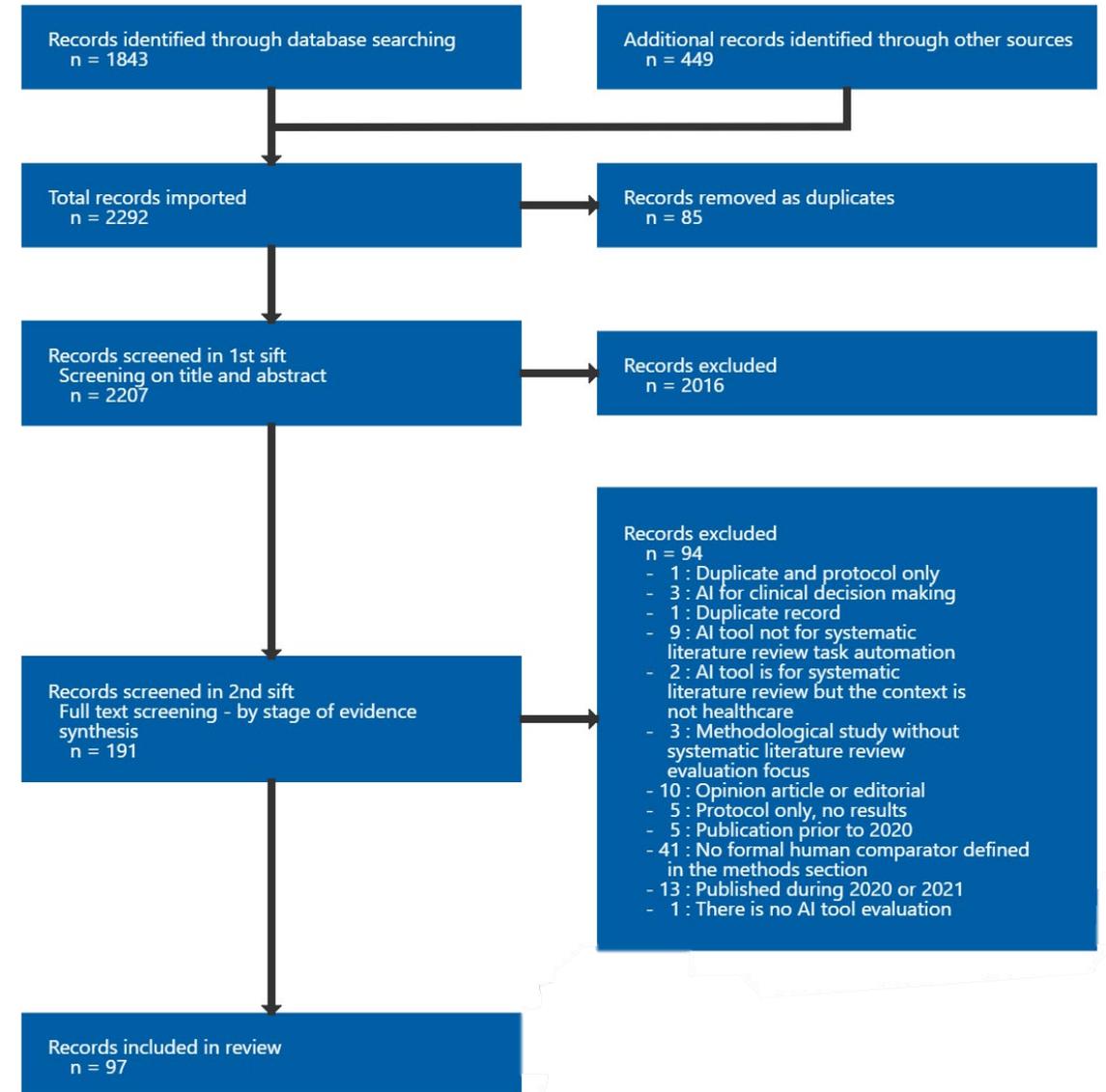


# AI in SLR

## Inclusion criteria

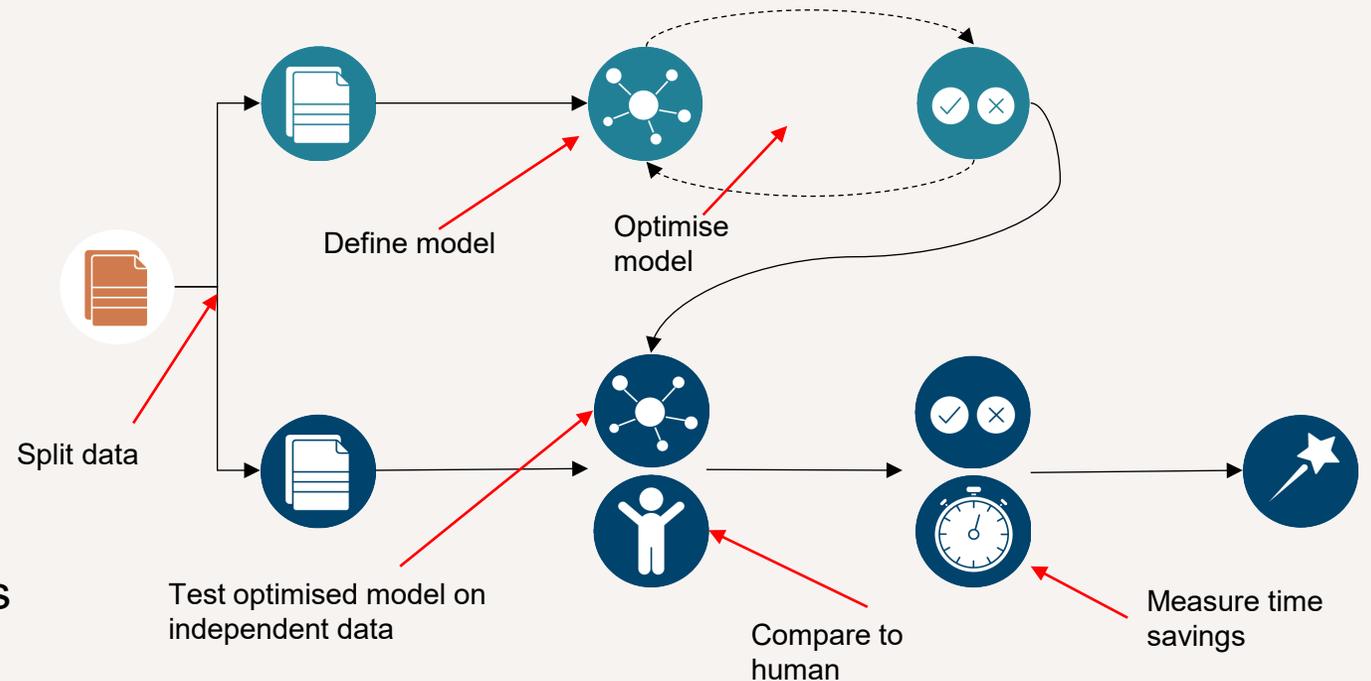
1. **English** language studies
2. **Healthcare** context
3. Published from **2022**
4. Evaluation of **AI** performance on **systematic literature review** tasks
5. Formal **comparison to human** reviewers.

NICE



# Common problems

1. No consistent approach to model testing.
2. No measure of time savings and/or human comparator.
3. No model description.
4. No model optimisation.
5. No data splitting.
6. No interoperability or standardisation across tasks.



# Risks and biases

| Risk/bias        | Description   | Potential Mitigations                        |
|------------------|---|--|
| Context bias     | Occurs when LLMs are biased towards the prompter, often occurring when the prompt template is too complex and the LLM learns the 'voice' of the user. | Independently verify prompts.                |
| Benchmark bias   | Occurs when the benchmark standard itself has errors. For example, making mistakes when using K-shot prompting.                                       | Independently verify prompts.                |
| Explainability   | LLMs are 'black boxes', the results should not only be plausible but also verifiable and justifiable.   | Output justifications and exact quotes.      |
| Hallucinations   | When LLMs 'make up' plausible but incorrect results.  | Prompt engineering: RAG, self-consistency    |
| Algorithmic bias | Occurs when a model favours one subset of data over another, for example medical domains.   | Present performance across multiple domains. |

# Thank you.

Stay up to date



- Sign up to our newsletters - scan the QR code or visit [nice.org.uk/newsletters](https://nice.org.uk/newsletters)
- Follow us on social media - visit [link.tr.ee/nicecomms](https://link.tr.ee/nicecomms)



- “Is it possible to deliver *faster, scalable* outputs without compromising on methodological rigor?”
  - ✓ **Yes, using Gen AI assisted evidence synthesis frameworks**
- “Can generative AI truly be trusted to conduct systematic literature reviews for HTA submissions?”
  - ✓ **Yes, by keeping human-in-loop and adhering to HTA guidance on the use of Gen AI**
- “Are there any real-world examples of Gen AI-assisted SLRs in HTA submissions?”
  - ✓ **Yes, Singh 2025 and Makhija 2025 posters provide the NICE landscape assessment on the use of AI-assisted SLRs and a detailed example**
- “Have we reached a point where AI-assisted SLRs can stand up to regulatory scrutiny and truly save time and effort in real-world HTA submissions?”
  - ✓ **Yes, early implementations show around 50% savings in time and cost**



**Thank you**

