

# Assessing AI and ML Tool Performance in SLRs: A Targeted Literature Review and Performance Benchmark Framework

Raju Gautam, PhD<sup>1</sup>, Saeed Anwar, MSc<sup>2</sup>, Ratna Pandey, MSc<sup>2</sup>, Khushbu Baranwal, MSc<sup>2</sup>, Tushar Srivastava, MSc<sup>1</sup>  
<sup>1</sup>ConnectHEOR, London, UK; <sup>2</sup>ConnectHEOR, Delhi, India. | Email: [raju.gautam@connectheor.com](mailto:raju.gautam@connectheor.com)

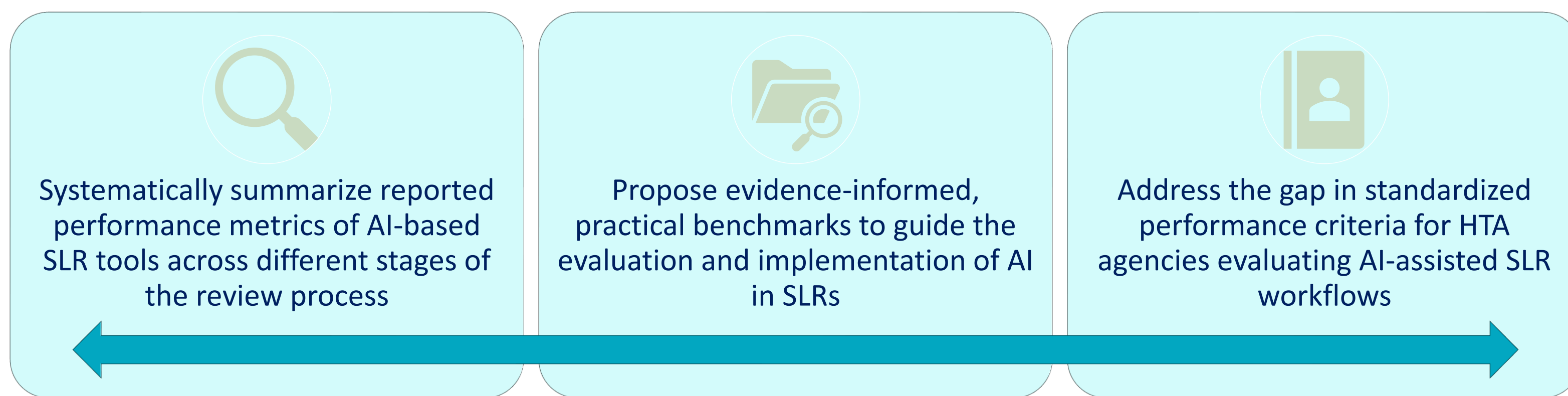
## KEY MESSAGE

Current evidence suggests that well-configured AI/ML tools can achieve performance comparable to human reviewers in SLRs, and this study presents the first evidence informed, HTA-aligned benchmark framework to guide evaluation of AI-assisted workflows. Based on 25 studies, title/abstract screening sensitivity ranged 14%–99% and specificity 19%–99%, while full-text screening showed consistently high sensitivity (76–99%). Establishing evidence-based performance benchmarks is essential for the adoption of AI-driven approaches with confidence by HTA agencies.

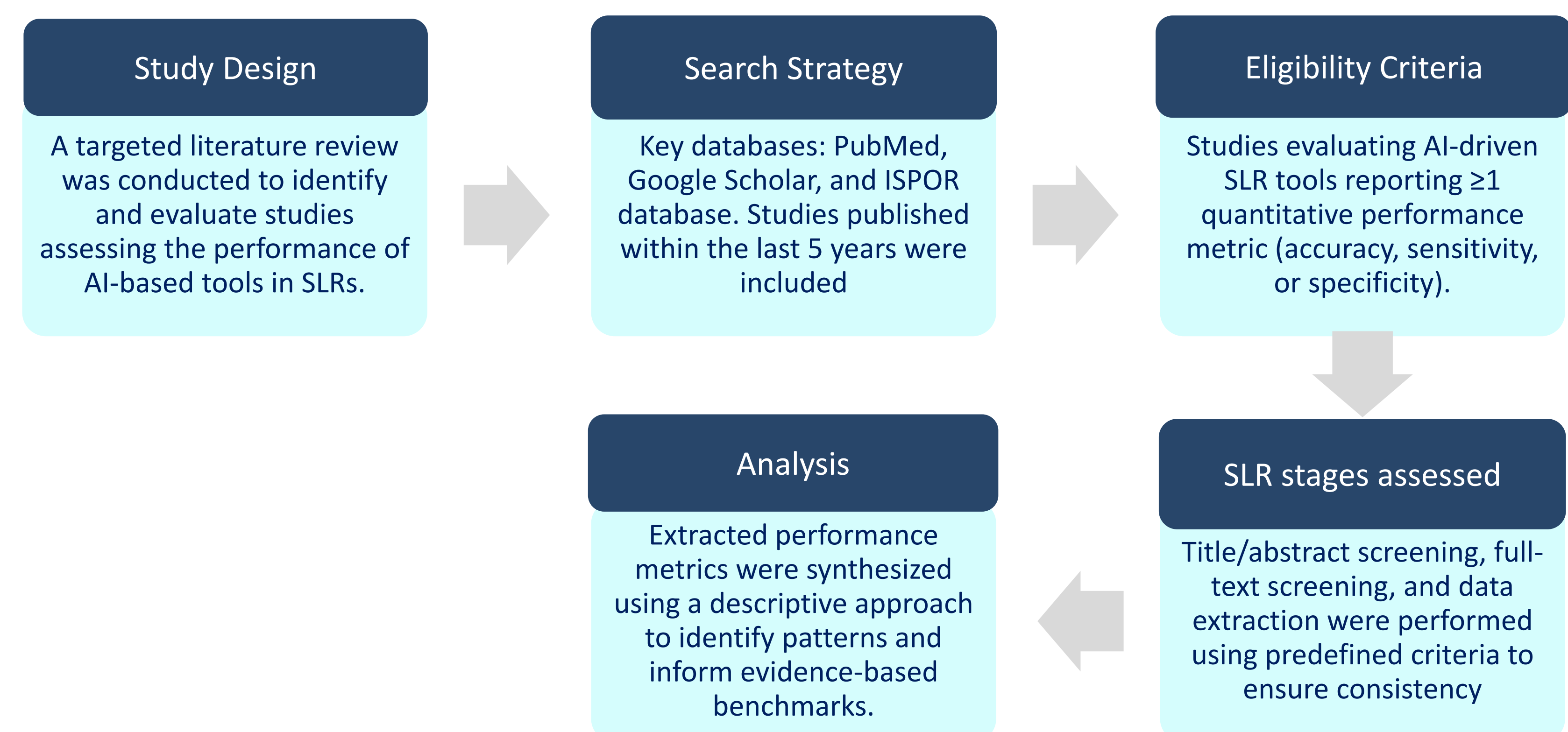
## BACKGROUND

- Systematic literature reviews (SLRs) are a cornerstone of evidence-based decision-making in healthcare, particularly for health technology assessment (HTA) agencies. However, traditional SLR processes are time-consuming and resource-intensive, often requiring extensive manual screening and data extraction.<sup>1</sup>
- With the growing volume of published research, AI/ML tools have emerged as promising solutions to improve efficiency and reduce workload. These tools can assist in key stages such as title/abstract screening, full-text review, and data extraction.<sup>2</sup>
- Despite increasing adoption, there remains uncertainty about their reliability and consistency compared to human reviewers. HTA agencies currently lack standardized benchmarks to evaluate AI performance in SLR workflows.<sup>2</sup>

## OBJECTIVES



## METHODS

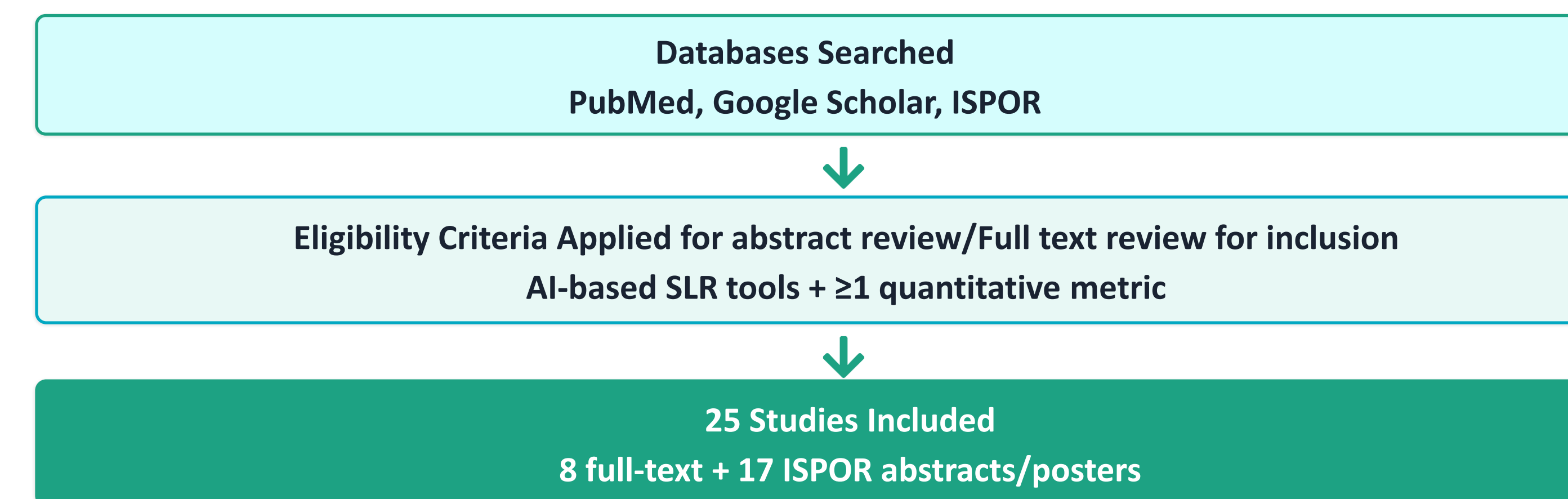


In the SLRs that are aimed to support HTA submissions, sensitivity is of utmost importance to minimize missed evidence, while maintaining acceptable specificity to reduce screening burden.

- Performance metrics:**
- Accuracy:** Agreement between AI and human reviewer decisions, reflecting overall classification performance
  - Sensitivity:** Ability of AI to correctly identify relevant studies, critical for minimizing missed evidence
  - Specificity:** Ability of AI to correctly exclude irrelevant studies, important for reducing screening workload

## RESULTS

Figure 1: Study Selection Flow



- Studies were identified through database searches and screened based on predefined eligibility criteria, resulting in 25 included studies for final analysis.
- Most of included evidence was derived from ISPOR abstracts/posters, reflecting the emerging and evolving nature of AI applications in SLRs.



Table 1: HTA aligned PROPOSED BENCHMARKS

This table presents the proposed evidence-informed, HTA-aligned performance benchmarks for AI-assisted SLR workflows across key stages, with emphasis on high sensitivity to minimize missed evidence while maintaining acceptable specificity and accuracy

Title/Abstract Screening		
≥95% Sensitivity	≥80% Specificity	≥85% Accuracy
Full-Text Screening		
≥90% Sensitivity	≥70% Specificity	≥80% Accuracy
Data Extraction		
— Sensitivity	— Specificity	≥85% Accuracy

Table 2: Performance Metrics by SLR Stage

- It summarizes the wide variability in AI performance across SLR stages, with consistently higher sensitivity in full-text screening and greater variability in specificity and accuracy across all stages.

SLR Stage	Accuracy	Sensitivity	Specificity
Title/Abstract Screening	10%–100%	14%–99%	19%–99%
Full-Text Screening	40%–98.5%	76%–99%	19%–77%
Data Extraction	40%–100%	—	—

Note: Only 5 studies reported data extraction accuracy

## RESULTS (CONTINUED)

Figure 2: Performance range in full-text screening

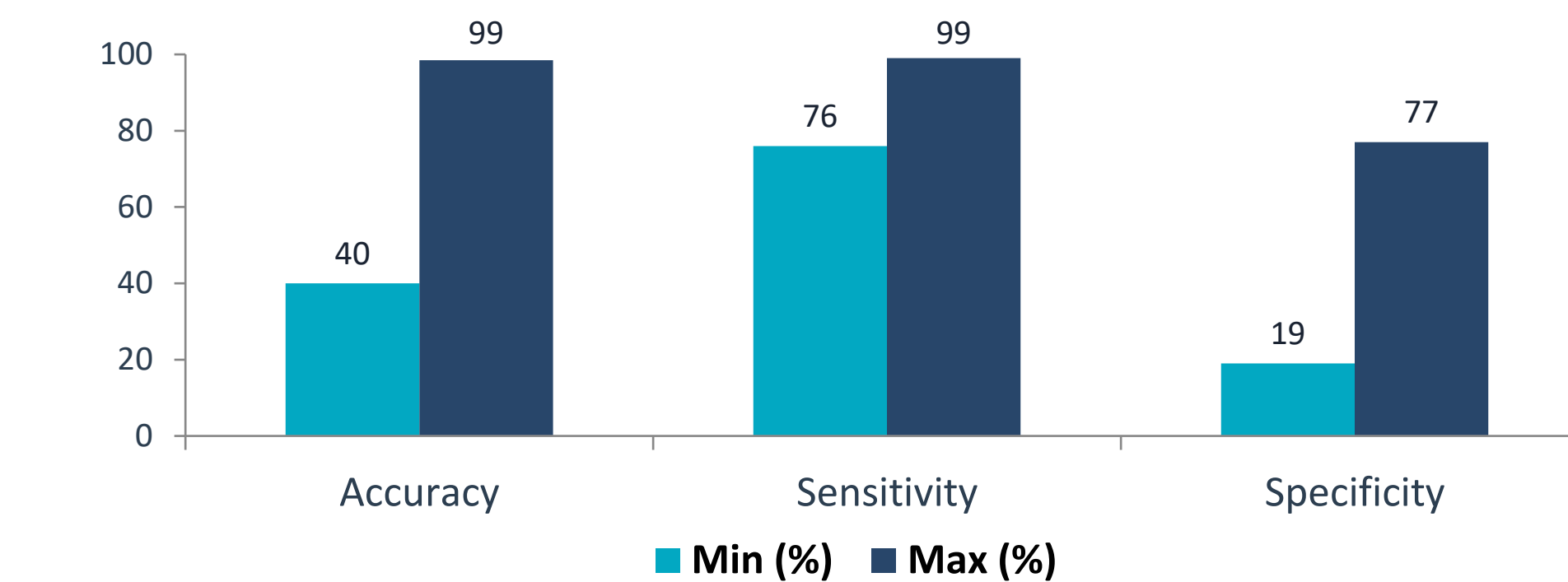
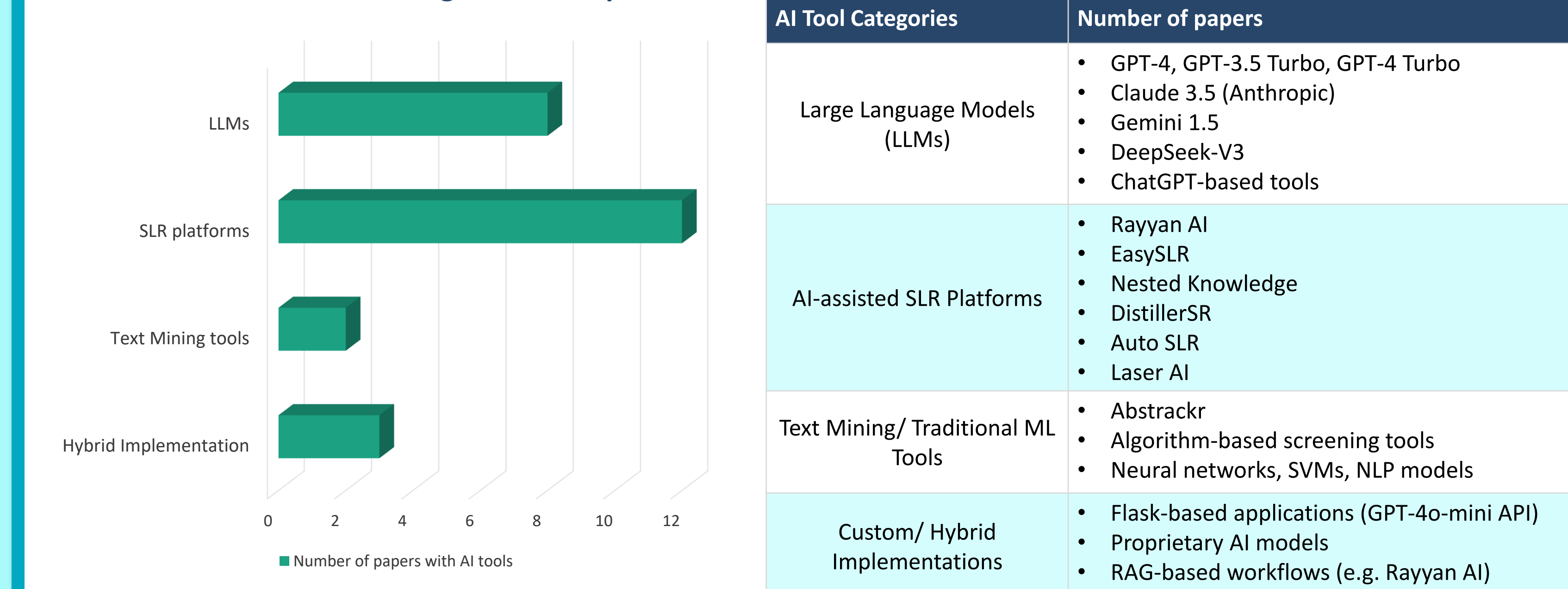


Table 3: AI Tools and Technologies Landscape



- Across 25 studies, diverse AI tools were evaluated spanning multiple technologies and platforms. These included LLMs, AI-assisted SLR platforms, traditional ML models and hybrid implementations
- The increasing adoption of LLM-driven and AI-assisted platforms reflects the rapid evolution of AI in SLR workflows

## CONCLUSIONS

- AI is transforming SLRs by combining efficiency with near human-level performance.
- This study demonstrates that well-configured and validated AI-based tools can achieve comparable sensitivity and specificity to human reviewers, with particularly high sensitivity observed in full-text screening.
- For HTA decision-making, prioritizing sensitivity remains critical to minimize the risk of missed evidence while maintaining acceptable specificity to reduce workload. The variability observed across studies underscores the need for the standardized evaluation of AI tools.
- This study provides the first evidence-informed, HTA-aligned benchmark framework to support consistent validation and confident adoption of AI-assisted SLR workflows within transparent, human-in-the-loop approach.

## REFERENCES

- Li Y, et al. Enhancing systematic literature reviews with generative artificial intelligence: development, applications, and performance evaluation. *J Am Med Inform Assoc.* 2025 Apr 1;32(4):616-625.
- Atkinson, C. F. (2023). Cheap, quick, and rigorous: Artificial Intelligence and the Systematic Literature Review. *Social Science Computer Review*, 42(2), 376–393.
- Carey N, et al. A text-mining tool generated title-abstract screening workload savings: performance evaluation versus single-human screening. *Journal of Clinical Epidemiology*, 2022; 149, 53-59.

Financial Disclosure: Authors are employees of ConnectHEOR Limited. No external funding received. No conflict of interest to declare.

