

# Assessing the Effectiveness of Large Language Models in Automating Systematic Literature Reviews: Findings from Recent Studies

MSR155

Sumeyye Samur<sup>1</sup>, Bhakti Mody<sup>1</sup>, Rachael Fleurence<sup>1,2,\*</sup>, Elif Bayraktar<sup>1</sup>, Turgay Ayer<sup>1,3</sup>, Jagpreet Chhatwal<sup>1,4</sup>

<sup>1</sup> Value Analytics Labs, Boston, MA, USA; <sup>2</sup> National Institutes of Health, DC, USA; <sup>3</sup> Georgia Institute of Technology, Atlanta, GA, USA; <sup>4</sup> Massachusetts General Hospital Institute for Technology Assessment, Harvard Medical School, Boston, MA, USA; \*Dr. Fleurence contributed to this article in her personal capacity

## BACKGROUND

- Systematic literature reviews (SLRs) are foundational for evidence synthesis in medical research but are often time-intensive and laborious. On average, completing a traditional SLR can take between 12 and 18 months and typically involves a team of at least three researchers.<sup>1</sup>
- Large Language Models (LLMs) offer potential for automating SLR tasks such as screening, data extraction, and bias assessment.
- However, the feasibility and performance of these models in conducting different steps of SLRs remain insufficiently documented.

## OBJECTIVE

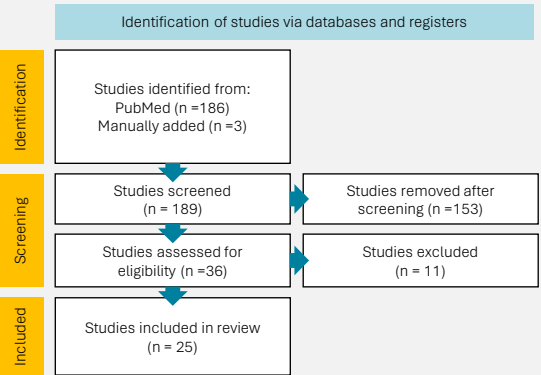
To evaluate the performance of LLMs across key tasks in the SLR process by reviewing recently published studies.

## METHODS

We conducted a targeted literature review of studies applying LLMs in SLR workflows, focusing on performance metrics across multiple review tasks.

**Study Identification:** 25 studies published between Jan 2023 and Jan 2025<sup>2-26</sup> (Figure 1)

**Figure 1: PRISMA Diagram for the Targeted Review of Studies applying LLM in SLR steps**



## KEY FINDINGS

- LLMs show strong potential in automating key components of SLRs, particularly in early screening and data extraction.
- While different LLMs exhibit varying strengths and limitations across tasks, no study has comprehensively evaluated their performance across all SLR steps using AI.
- A hybrid approach that leverages the strengths of multiple LLMs may improve the overall efficiency and accuracy of SLRs.
- Performance of LLMs varies by task, suggesting that some steps in the SLR process may be more suitable for automation than others.
- Human oversight remains essential for conducting AI-assisted SLRs.
- Further research is needed to fully explore the capabilities and limitations of LLMs within SLR workflows.

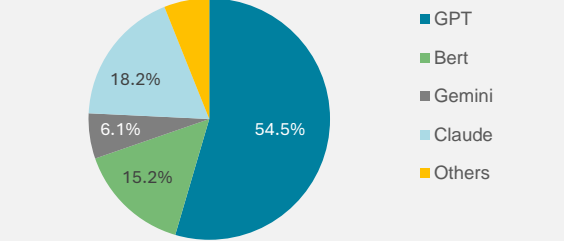
- LLMs Assessed:** GPT, Gemini, Bert, Claude and others such as Google PaLM 2 and Meta Llama 2
- SLR Tasks Analyzed:** Title and abstract screening, full-text screening, data extraction and bias assessment
- Performance Metrics:** Sensitivity, specificity, accuracy, and inter-rater agreement (e.g., kappa score).

## RESULTS

### LLMs Assessed

Among all LLMs, GPT was the most frequently evaluated LLM across studies (Figure 2).

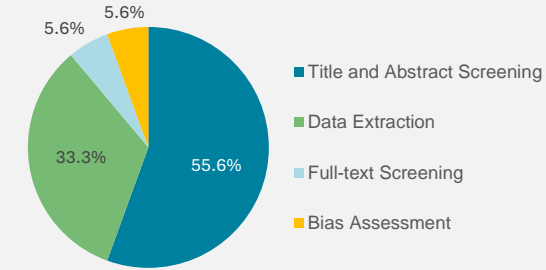
**Figure 2: Distribution of LLMs Evaluated (among 33 independent evaluations)**



### Performance Evaluation

Among all SLR steps, title and abstract screening using LLMs was the most frequently evaluated one across studies (Figure 3).

**Figure 3: Distribution of SLR Tasks Evaluated (among 18 independent evaluations)**



LLMs demonstrated strong performance across various systematic review tasks (Table).

- Title screening:** Sensitivity 94.3%–96.2%; specificity 85.5%–99.6%.
- Abstract screening:** Accuracy 80%–97.5%; sensitivity 62%–95%; specificity 65%–98.7%.

- Full-text screening:** Accuracy 87%; sensitivity 71.4%; specificity 93.8%.
- Data extraction:** Accuracy 67%–96.3%; sensitivity 36%–96.2%; specificity >80%.
- Bias assessment:** Strong agreement with human reviewers (kappa >0.89 for abstracts; 0.65 for full-text).

**Table: Performance of LLMs Across SLR Tasks**

Task	Best Performing Model	Metric	Performance Range
Title Screening	GPT-3.5	Sensitivity	94.3%–96.2% <sup>18</sup>
		Specificity	85.5%–99.6% <sup>18</sup>
Abstract Screening	GPT-4 (Acc, Spec) GPT-3.5 (Sens)	Accuracy	80%–97.5% <sup>22, 26</sup>
		Sensitivity	62%–95% <sup>23, 26</sup>
		Specificity	65%–98.7% <sup>22, 23</sup>
		Accuracy	87% <sup>22</sup>
Full-Text Screening	GPT-4	Sensitivity	71.4% <sup>22</sup>
		Specificity	93.8% <sup>22</sup>
Data Extraction	Claude 2 (Acc, Sens) GPT-4 (Spec)	Accuracy	67%–96.3% <sup>14, 17</sup>
		Sensitivity	36%–96.2% <sup>14, 17</sup>
		Specificity	>80% <sup>17</sup>
		Kappa (Abstracts)	>0.89 <sup>22</sup>
Bias Assessment	GPT-4	Kappa (Full-text)	0.65 <sup>22</sup>

## REFERENCES

- Systematic Reviews. <https://belmont.fda.gov/systematicreviews>. Accessed by April, 2025.
- Tan VT, Gartlehner G, Yaacoub S, et al. Sensitivity and Specificity of Using GPT-3.5 Turbo Models for Title and Abstract Screening in Systematic Reviews and Meta-analyses. *Ann Intern Med* 2024;177:781–789.
- Akismetov O, Jiang X, Palade V. A question-answering framework for automated abstract screening using large language models. *J Am Med Inform Assoc* 2024;31:1939–1952.
- Gorska A, Tacconelli E. Towards Autonomous Living Meta-Analyses: A Framework for Automated Systematic Review and Meta-Analyses. *Stud Health Technol Inform* 2024;316:378–382.
- Chen H, Desampal L, Linoué V, et al. Hallucination Rates and Reference Accuracy of ChatGPT and Bard for Systematic Reviews: Comparative Analysis. *J Med Internet Res* 2024;26:e53164.
- Gwon YN, Kim JH, Chung HS, et al. The Use of Generative AI for Scientific Literature Searches for Systematic Reviews: ChatGPT and Microsoft Bing AI Performance Evaluation. *J Med Inform* 2024;12:e51187.
- Xian M, Ayub U, Nasir SA, et al. Collaborative Large Language Models for Automated Data Extraction in Living Systematic Reviews. *medRxiv* 2024;2023.12.12.23291187.
- Konec A, Thomas L, Gartlehner G, et al. Performance of two large language models for data extraction in evidence synthesis. *Res Synth Methods* 2024;15:618–624.
- Ozmi T, Okada T, Aoki Y, et al. Human-Comparable Sensitivity of Large Language Models in Identifying English Studies Through Title and Abstract Screening: 3-Layer Strategy Using GPT-3.5 and GPT-4 for Systematic Reviews. *J Med Internet Res* 2024;26:e52758.
- Liu M, Sun J, Tan X. Evaluating the effectiveness of large language models in abstract screening: a comparative analysis. *Syst Rev* 2024;13:219.
- Tao K, Omeran ZA, Zhou R, et al. GPT-4 performance on querying scientific publications: reproducibility, accuracy, and impact of an instruction sheet. *BMC Med Res Methodol* 2024;24:139.
- Fabiano N, Gupta A, Shambra N, et al. How to optimize the systematic review process using AI tools. *JCP Adv* 2024;4:e12234.
- Landschulz A, Anselmi D, Macaluso S, et al. Implementation and evaluation of an additional GPT-4-based reviewer in PRISMA-based medical systematic literature reviews. *Int J Med Inform* 2024;189:105531.
- Issay M, Ghannai H, Koliak S, et al. Methodological insights into ChatGPT's screening performance in systematic review. *BMC Med Res Methodol* 2024;24:78.
- Qin X, Liu J, Wang Y, et al. Natural language processing for evidence synthesis: title and abstract screening when updating systematic reviews. *J Clin Epidemiol* 2021;133:121–129.
- Aum S, Choe S. sBERT: automatic article classification model for systematic review using BERT. *Syst Rev* 2021;10:285.
- Kohender Gargan G, Hamedoodi MH, Hajjafarali M, et al. Enhancing title and abstract screening for systematic reviews with GPT-3.5 turbo. *BMC Evol Based Med* 2024;29:69–70.