How Accurate are Large Language Models for Abstract Screening in Systematic Literature Reviews?

Alparslan Sertel¹, <u>Sumeyye Samur</u>¹, Turgay Ayer^{1,2,3}, Jagpreet Chhatwal^{1,4}

¹Value Analytics Labs, Boston, Massachusetts, USA; ²Georgia Institute of Technology, Atlanta, Georgia, USA; ³Emory Medical School, Atlanta, Georgia, USA; ⁴Massachusetts General Hospital Institute for Technology Assessment, Harvard Medical School, Boston, MA, USA.

BACKGROUND

- In healthcare, systematic literature reviews (SLR) are essential for examining specific research topics.
- During screening, human reviewers often assess thousands of articles, which is a time-consuming task prone to errors.
- Recent advancements, like using large language models (LLMs) such as ChatGPT, show potential for automating SLR, especially in text-heavy tasks like screening.
- Prompt engineering can substantially affect the performance of LLMs for conducting SLRs.

OBJECTIVE

Our goal was to evaluate the feasibility of conducting abstract screening using LLMs and explore the accuracy of various prompting techniques in the abstract screening step of SLRs. During screening, human reviewers often assess thousands of articles, which is a time-consuming task prone to errors.

METHODS

 We implemented a binary classification model in Python with the OpenAI GPT-4 API to screen abstracts and categorize them as included or excluded (Figure 1)

Figure 1: Abstract Screening Model Pipeline



KEY FINDINGS

- This study sheds light on the effectiveness of five well-known prompting techniques in a conventional abstract screening step of an SLR
- Few-Shot prompting shows the highest performance in terms of accuracy
- The integration of LLMs holds promise in revolutionizing the landscape of SLR However, further improvements are needed to increase the accuracy of LLMs before they can be fully implemented for automating SLRs
- We tested five prompting techniques: Zero-shot, Few-shot, Chain-of-thought (COT), Zero-shot COT, and Dividing into subtasks (Table 1)

Figure 1: Examples for Used Prompt Techniques



- queries on PubMed and Embase databases.
- Results were compared based on accuracy, defined by percentage of abstracts matching human reviewers' decisions, execution time, and cost.



 The overall accuracy, encompassing the inclusion and exclusion of abstracts, was highest with Few-shot (82%) and lowest with Chain-of-thought (65%) (Figures 2 and 3).





Figure 3: Accuracy for included and excluded papers



 Although various factors such as the number of tokens, server load, network speed can influence the execution time, in our experiments, Zero-shot was the fastest method whereas the Zero-shot COT was the slowest one (Figure 4).





• Because cost is primarily linked to the number of tokens in both input and output prompts, the costliest technique was Few-shot (\$280) with Zero-shot being the least costly one (\$160) (Figure 5).

Figure 5: Cost (\$) by each Prompting Technique



LIMITATIONS

- Performance of the prompting techniques are evaluated by comparing them with human reviewers' decisions, which are presumed as the gold standard.
- The binary classification model is built over GPT-4 without fine-tuning or additional training, only alteration being the prompting techniques.
- We limited our study to abstract screening step of SLR, and future research should evaluate the performance of fulltext review using LLMs.

REFERENCES

https://www.promptingguide.ai/techniques/zeroshot
https://www.promptingguide.ai/techniques/fewshot
Wei et al. (2022)
Kojima et al. (2022)
Ittps://docs.anthropic.com/claude/docs/break-tasks-into-subtasks

Corresponding Author: Jagpreet Chhatwal, PhD; JagChhatwal@mgh.harvard.edu

tics i cabs

MSR44