

Background

Large language models (LLMs) have the potential to increase efficiencies relative to manually conducted systematic reviews; however, caution remains to ensure gold standards are not compromised (1-4).

Objectives

- ✓ Collate studies reporting LLM data extraction of clinical publications
- ✓ Explore performance of LLM extraction according to data domain
- ✓ Identify any factors influencing LLM extraction performance

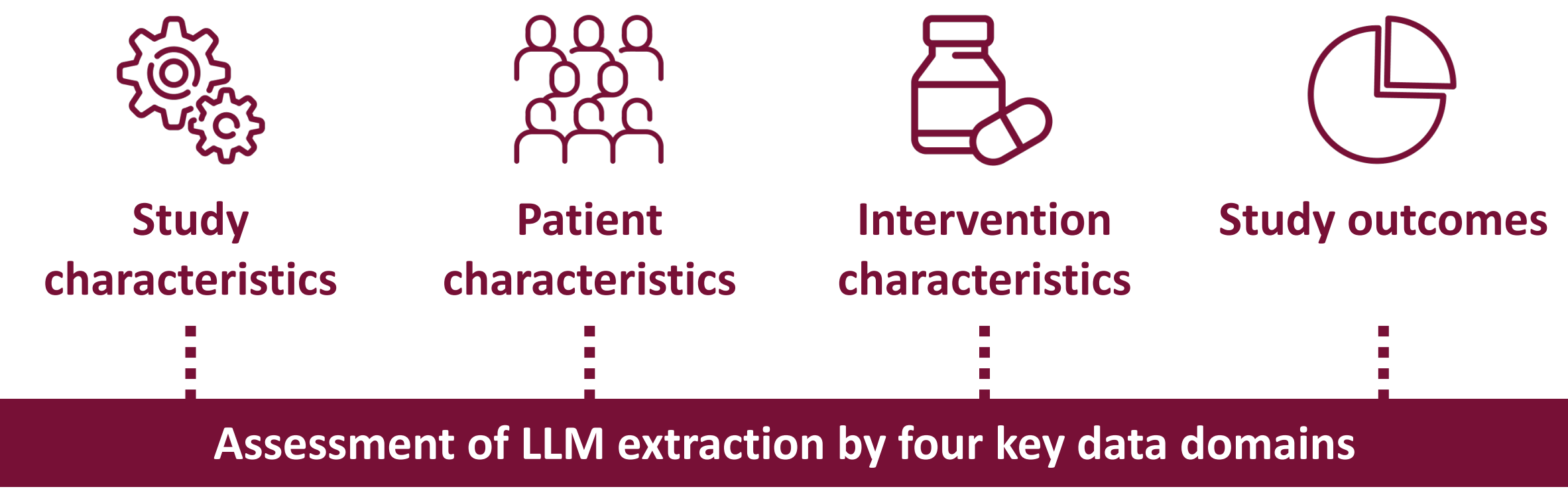
Methods

A rapid systematic review was conducted to identify records reporting LLM data extraction of clinical publications. To be considered for inclusion, clear reporting of the LLM used to perform data extraction was required. A two-step approach was then used for subsequent analyses.



Step 1: Map LLMs for all included records

Step 2: Records reporting quantitative LLM extraction performance by data domain were eligible for further analysis. The data domains of interest were based on those commonly extracted from clinical publications:



Results

A total of 31 records that reported data extraction of clinical publications using an identifiable LLM were identified.
Of these included records, 15 reported the performance of LLM data extraction by data domain.

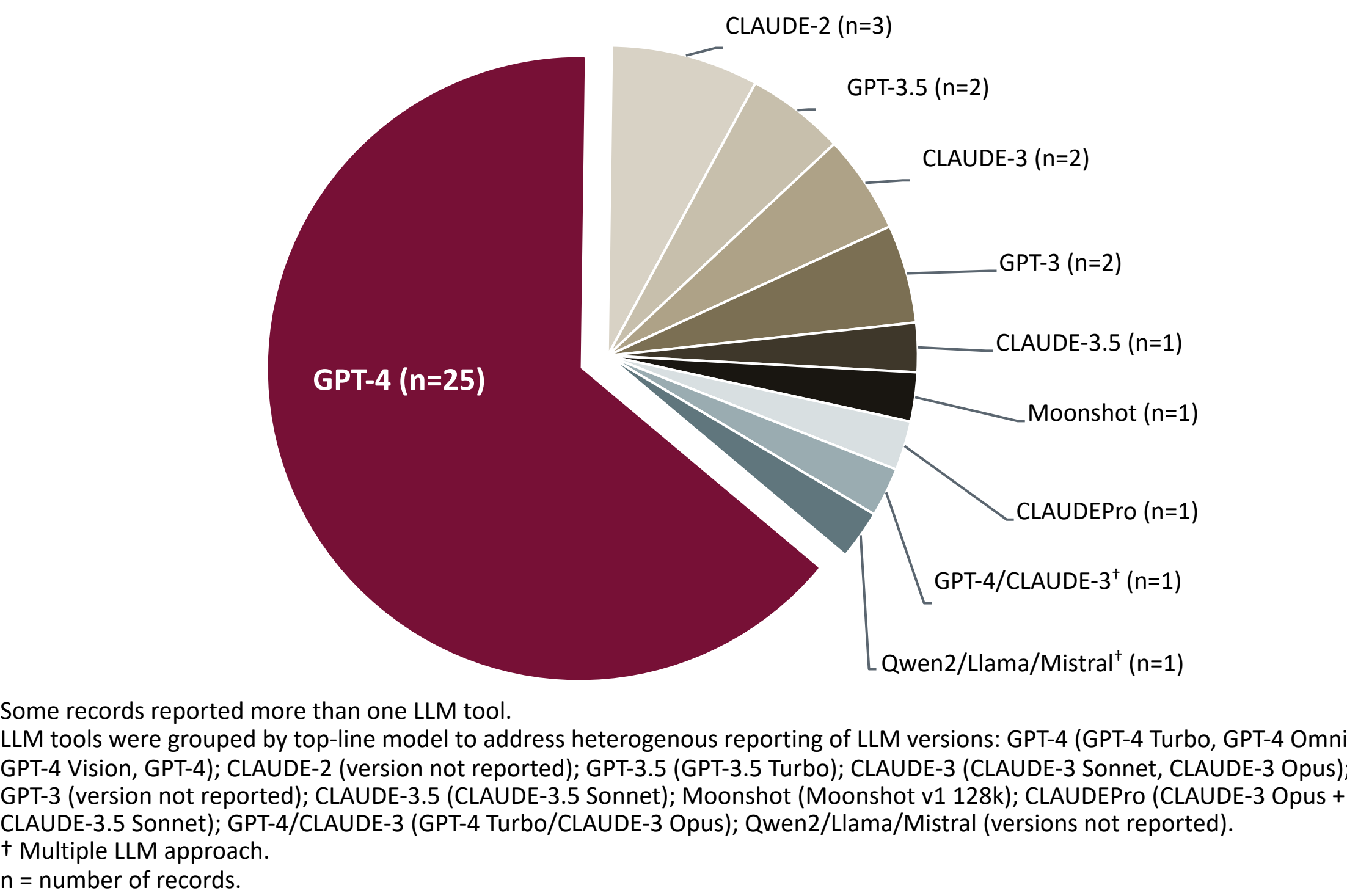


A variety of LLM tools were identified for data extraction

Across the 31 records, 10 different LLM tools were identified.

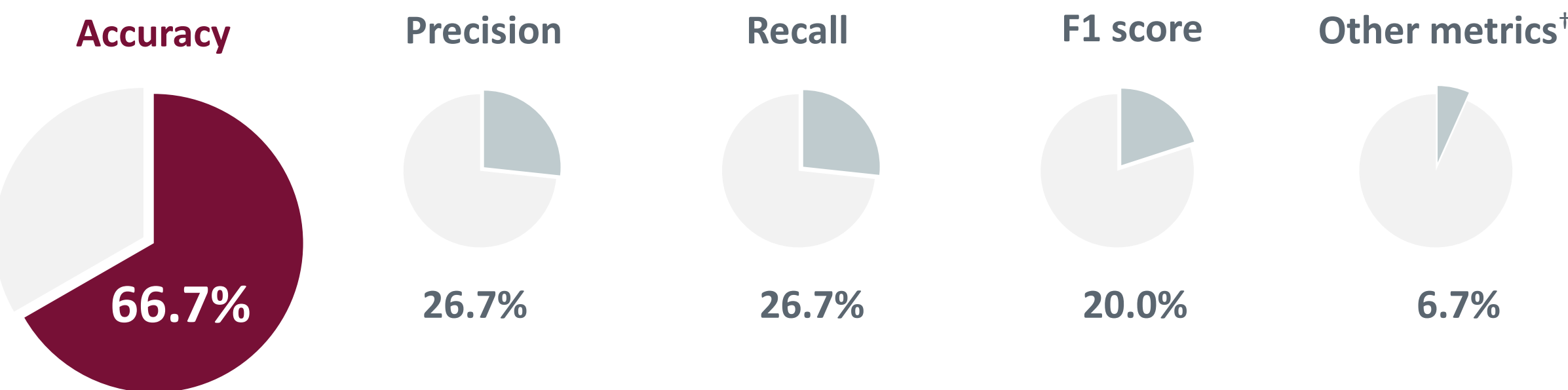
Generative Pre-trained Transformer (GPT)-4 models were the most commonly reported LLM tool (n=25 records; 80.6%), followed by CLAUDE-2 (n=3 records; 9.7%; Figure 1). The remaining LLM tools were reported by ≤2 records. Two records reported extraction of clinical publications using a collaboration of multiple LLM tools: GPT-4/CLAUDE-3 and Qwen2/Llama/Mistral.

Figure 1: Number of records reporting LLMs for data extraction of clinical publications



Heterogenous reporting metrics for LLM performance

Figure 2: Reporting metrics for LLM extraction performance across records reporting performance by data domain (n=15)

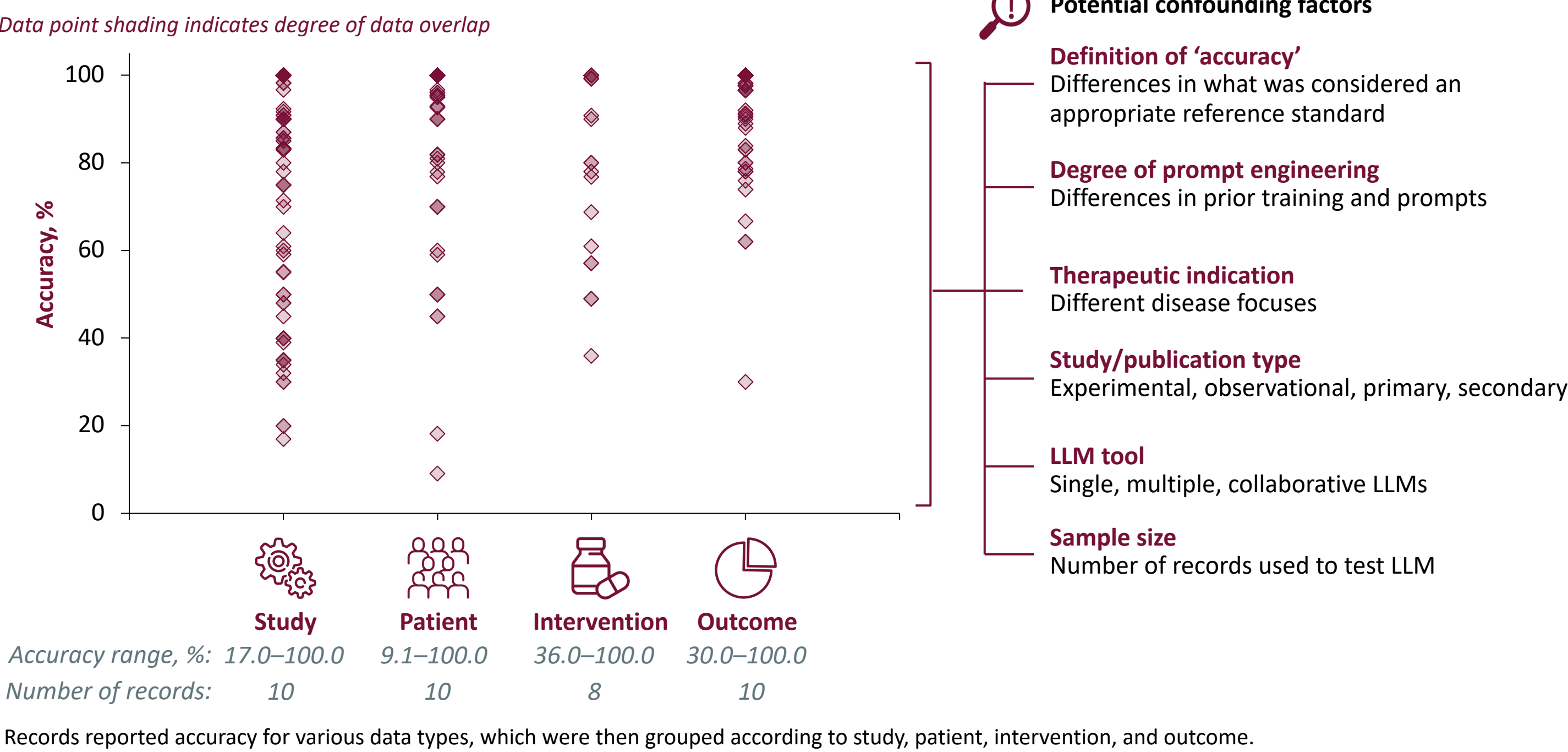


Accuracy was the most common reporting metric (Figure 2), often measured as a “decision match percentage” with a human reviewer.

Reported accuracy was highly variable by data domain

Reported accuracy varied greatly for each data domain (Figure 3). Records differed in relation to various factors, which can be considered potential confounding factors.

Figure 3: Reported accuracy percentages of LLM extraction, by data domain (n=10)



Data type/format, publication language, and use of single or multiple LLMs can impact LLM extraction performance

Key factors impacting LLM extraction performance were identified through a thematic analysis (Table 1).

Table 1: Thematic analysis – Reported factors impacting LLM performance, by LLM tool†

LLM	Key factors impacting LLM performance			
	Data type	Data format	Publication language	Single versus multiple LLM
GPT-4	NR	Free text had better accuracy than table or figure text	NR	NR
GPT-3	NR	Free text had better accuracy than table text	NR	NR
CLAUDEPro	String data had better accuracy/recall than numerical/mixed data	NR	NR	NR
CLAUDE-3.5	NR	NR	English language performed better than Chinese language	NR
GPT-4/CLAUDE-3	NR	NR	NR	Multiple LLMs performed better than individual LLMs

† Thematic analysis was conducted across the 15 records included for further analysis to identify key factors reported to impact LLM performance. Not all records provided insights into factors impacting LLM performance. As the table only presents LLM tools with reported impacting factors, not all records (and associated LLM tools) are represented in the table.

Conclusions

- GPT-4 models were the most commonly reported LLM tool for data extraction of clinical publications
- Substantial heterogeneity exists in the reporting of LLM extraction performance
- LLM data extraction has potential for high accuracy across all four data domains; however, no data domain was found to show reliable accuracy rates
- Standardisation in assessing LLM data extraction performance, and full transparency in its reporting, is required to support future research and guide implementation in evidence synthesis processes

Scan for a video walkthrough



References

- Lieberum JL, et al. Large language models for conducting systematic reviews: on the rise, but not yet ready for use—a scoping review. J Clin Epidemiol. 2025;181:111746.
- Luo X, et al. Potential roles of large language models in the production of systematic reviews and meta-analyses. J Med Internet Res. 2024;26:e56780.
- Yun HS, Marshall U, Trikalinos TA, Wallace BC. Appraising the potential uses and harms of LLMs for medical systematic reviews. arXiv preprint arXiv:2305.11828. 2023.
- Forero DA, Abreu SE, Tovar BE, Oermann MH. Large Language Models and the Analyses of Adherence to Reporting Guidelines in Systematic Reviews and Overviews of Reviews (PRISMA 2020 and PRIOR). J Med Syst. 2025;49(1):80.

Abbreviations

- GPT, Generative Pre-trained Transformer
- LLM, large language model
- NR, not reported
- PRISMA, Preferred Reporting Items for Systematic Reviews and Meta-Analyses