

# 1 Title

2 Large Language Model in Medical Information Extraction from Titles and Abstracts with  
3 Prompt Engineering Strategies: A Comparative Study of GPT-3.5 and GPT-4

## 5 Authors:

### 6 Authors:

7 Yiyi Tang<sup>a,b,#</sup>, Ziyang Xiao<sup>b,c,#</sup>, Xue Li<sup>a,c,d,\*</sup>, Qiwen Fang<sup>c</sup>, Qingpeng Zhang<sup>c,e</sup>, Daniel Yee Tak  
8 Fong<sup>f</sup>, Francisco Tsz Tsun Lai<sup>c,d,g</sup>, Celine Sze Ling Chui<sup>d,f,h</sup>, Esther Wai Yin Chan<sup>c,d</sup>, Ian Chi  
9 Kei Wong<sup>c,d,i,j</sup>, Research Data Collaboration Task Force<sup>k</sup>

10

11 <sup>a</sup> Department of Medicine, School of Clinical Medicine, Li Ka Shing Faculty of Medicine, The  
12 University of Hong Kong, Hong Kong SAR, China

13 <sup>b</sup> Department of Statistics and Actuarial Science, Faculty of Science, The University of Hong  
14 Kong, Hong Kong SAR, China

15 <sup>c</sup> Department of Pharmacology and Pharmacy, Li Ka Shing Faculty of Medicine, The  
16 University of Hong Kong, Hong Kong SAR, China

17 <sup>d</sup> Laboratory of Data Discovery for Health (D<sup>2</sup>4H), Hong Kong Science Park, Hong Kong SAR,  
18 China

19 <sup>e</sup> Musketeers Foundation Institute of Data Science, The University of Hong Kong, Hong Kong  
20 SAR, China

21 <sup>f</sup> School of Nursing, Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong  
22 Kong SAR, China

23 <sup>g</sup> Department of Family Medicine and Primary Care, Li Ka Shing Faculty of Medicine, The  
24 University of Hong Kong, Hong Kong SAR, China

25 <sup>h</sup> School of Public Health, Li Ka Shing Faculty of Medicine, The University of Hong Kong,  
26 Hong Kong SAR, China

27 <sup>i</sup> Research Department of Practice and Policy, School of Pharmacy, University College London,  
28 London, UK

29 <sup>j</sup> Aston Pharmacy School, Aston University, Birmingham, UK

30 <sup>k</sup> Li Ka Shing Faculty of Medicine, The University of Hong Kong, Hong Kong SAR, China

31

32 <sup>#</sup> Co-first author with equal contribution

33 <sup>\*</sup> Corresponding author

34

35 Xue Li

36 Assistant Professor

37 Department of Medicine and Department of Pharmacology & Pharmacy

38 LKS Faculty of Medicine, The University of Hong Kong

1 PB306, 3/F, Professional Block, Queen Mary Hospital  
2 102 Pok Fu Lam Road, Hong Kong  
3 Tel: +852 2255 3319  
4 Email: [sxueli@hku.hk](mailto:sxueli@hku.hk)

5

6 Co-authors:

7 Yiyi Tang: [yiyitang@connect.hku.hk](mailto:yiyitang@connect.hku.hk)

8 Ziyang Xiao: [xiaozy@connect.hku.hk](mailto:xiaozy@connect.hku.hk)

9 Qiwen Fang: [u3594701@connect.hku.hk](mailto:u3594701@connect.hku.hk)

10 Qingpeng Zhang: [qpzhang@hku.hk](mailto:qpzhang@hku.hk)

11 Daniel Yee Tak Fong: [dtyfong@hku.hk](mailto:dtyfong@hku.hk)

12 Francisco Tsz Tsun Lai: [fttlai@hku.hk](mailto:fttlai@hku.hk)

13 Celine Sze Ling Chui: [cslchui@hku.hk](mailto:cslchui@hku.hk)

14 Esther Wai Yin Chan: [ewchan@hku.hk](mailto:ewchan@hku.hk)

15 Ian Chi Kei Wong: [wongick@hku.hk](mailto:wongick@hku.hk)

16 Research Data Collaboration Task Force: [rdctf@hku.hk](mailto:rdctf@hku.hk)

17

18 Author contributions:

19 YT: Study concept and design; Information collection and screening; Data extraction,  
20 analysis, and cross-checking; Drafting of the manuscript;

21 ZX: Study concept and design; Information collection and screening; Data extraction,  
22 analysis, and cross-checking; Drafting of the manuscript;

23 XL: Study concept and design; Drafting of the manuscript; Study supervision;

24 QZ: Study concept and design; Study supervision;

25 EWYC: Study concept and design;

26 ICKW: Study concept and design;

27 All authors: Critical revision of the manuscript of significant intellectual contribution;

28 Data interpretation.

29

## 30 **Abstract**

31 **Background:** While it is believed that large language models (LLMs) have the potential to  
32 facilitate the review of medical literature, their accuracy, stability and prompt strategies in  
33 complex settings have not been adequately investigated. Our study assessed the capabilities  
34 of GPT-3.5 and GPT-4.0 in extracting information from publication abstracts. We also

1 validated the impact of prompt engineering strategies and the effectiveness of evaluating  
2 metrics.

3  
4 **Methodology:** We adopted a stratified sampling method to select 100 publications from  
5 nineteen departments in the LKS Faculty of Medicine, The University of Hong Kong,  
6 published between 2015 and 2023. GPT-3.5 and GPT-4.0 were instructed to extract seven  
7 pieces of information – study design, sample size, data source, patient, intervention,  
8 comparison, and outcomes – from titles and abstracts. The experiment incorporated three  
9 prompt engineering strategies: persona, chain-of-thought and few-shot prompting. Three  
10 metrics were employed to assess the alignment between the GPT output and the ground truth:  
11 ROUGE-1, BERTScore and a self-developed LLM Evaluator with improved capability of  
12 semantic understanding. Finally, we evaluated the proportion of appropriate answers among  
13 different GPT versions and prompt engineering strategies.

14  
15 **Results:** The average accuracy of GPT-4.0, when paired with the optimal prompt engineering  
16 strategy, ranged from 0.736 to 0.978 among the seven items measured by the LLM evaluator.  
17 Sensitivity of GPT is higher than the specificity, with an average sensitivity score of 0.8550  
18 while scoring only 0.7353 in specificity. The GPT version was shown to be a statistically  
19 significant factor impacting accuracy, while prompt engineering strategies did not exhibit  
20 cumulative effects. Additionally, the LLM evaluator outperformed the ROUGE-1 and  
21 BERTScore in assessing the alignment of information.

22  
23 **Conclusion:** Our result confirms the effectiveness and stability of LLMs in extracting  
24 medical information, suggesting their potential as efficient tools for literature review. We  
25 recommend utilizing an advanced version of LLMs and the prompt should be tailored to  
26 specific tasks. Additionally, LLMs show promise as an evaluation tool related for complex  
27 information.

28

## 29 **Introduction**

30 Large language models (LLMs), including the GPT series, have emerged as a promising tool to  
31 revolutionize many practices in medicine [1,2]. LLMs are distinct from traditional natural  
32 language processing (NLP) models in their ability to generate responses that align with users'  
33 requirements [3], without requiring dedicated fine-tuning for specialised tasks [4]. Medical  
34 evidence summarization is one of these areas where GPT shows promise to improve the  
35 traditional process of extracting information from the vast amount of medical research papers  
36 [5-7].

37 Research has also demonstrated the effectiveness and cost efficiency of employing these  
38 automated tools in medical information extraction [8, 9, 10]. For example, one study showed  
39 that text-mining-based single screening reduced workload by over 60% compared to  
40 alternative methods [9]. The advent of Large Language Model has created new possibilities in  
41 automated medical information extraction. Many pioneering experiments in 2024

1 demonstrating considerable enhancement in the functionality and accuracy of automated  
2 medical information extraction [10, 11, 12, 13, 14, 15].

3 Despite the promising potential of LLMs in literature review, there remains a need for  
4 comprehensive empirical research addressing common concerns on applying LLM to medical  
5 information extraction, including accuracy, consistency, adaptability across medical domain,  
6 and the effects of prompt engineering [16, 17, 18]. Take prompt engineering for instance:  
7 although it has been widely reported as a useful strategy to enhance LLM's performance [10,  
8 16, 19,], research also pointed out that over half of the research on effects of prompts failed to  
9 report baseline performance, making the positive gain less credible [20]. Such overlooked facts  
10 also include the simple models may achieve the better performance than models with delicate  
11 prompt design [16]; optimal models and prompt designs diverges among tasks [16,17] . With  
12 the methodology of LLM-related research still being unstandardized [21], many understated  
13 observations are worth of detailed investigation. The sophisticated patterns of LLM's  
14 performance remains unclear, indicating a notable lack of comprehensive research in  
15 addressing this confusion.

16 Therefore, this study designed a series of rigorous experiments to assess the capability of  
17 LLMs in extracting critical information from titles and abstracts of medical research literatures.  
18 Papers were sampled from various medical domains to ensure the generalizability of the results,  
19 rather than previous paradigms focusing on merely one medical domain. It performs  
20 comprehensive statistical analysis on validating the effects of two GPT models and three  
21 common prompt engineering, Persona, Chain of Thought, and Few-shot Prompting. It also  
22 incorporated three automated evaluators to enhance the reliability. By navigating the finer  
23 details of LLMs, this study aims to provide more empirical evidence in uncovering the  
24 complicated nature of LLM models in medical information extraction.

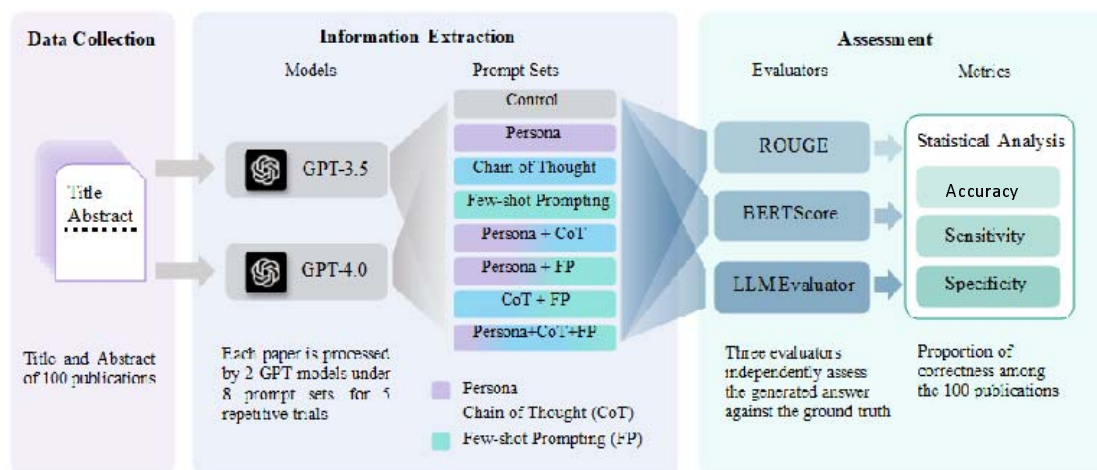
## 25 **Methods**

### 26 **Study design**

27 The scope of this study encompassed 100 research papers capturing a wide spectrum of  
28 subjects randomly selected from the publication pool of the Li Ka Shing Faculty of Medicine,  
29 University of Hong Kong. The selected papers were published between 1, January 2015, and  
30 31, December 2023, with titles and abstracts fully available on Scopus. To ensure  
31 comprehensive coverage, we adopted a stratified sampling method to randomly select papers  
32 from eighteen departments in proportion to the total number of publication records affiliated  
33 with that department. The paper's affiliation is the institution affiliated with the corresponding  
34 author. The departments and their related domains of the paper selected are presented in  
35 Supplementary Material 1.

36 Figure 1 presents the overall study design. All titles and abstracts were obtained from the  
37 Scopus online dataset and pre-processed to remove unreadable characters. Two undergraduate  
38 student researchers with training background on big data and statistics independently labelled

1 the information to be extracted according to pre-defined criteria, to obtain the ground truth. To  
2 ensure accuracy, a third reviewer, a full-time senior research assistant with postgraduate degree  
3 in epidemiology, was employed cross-checking to establish the ground truth.



4

5 Figure 1. Flowchart of overall study design.

6

7 Subsequently, the titles and abstracts were proceeded to GPTs to extract information. We  
8 implemented several prompt sets to compare the effectiveness of prompt engineering. The  
9 assessment of the information extraction performance was based on semantic similarity  
10 between GPT's output and ground truth, measured by several NLP metrics and a  
11 self-developed independent LLM evaluator based on GPT-4.0. Finally, we performed a  
12 statistical analysis on the results.

13 To compare the performance of GPT-3.5 and GPT-4.0, we conducted independent evaluations  
14 using the latest model versions at the time of study: *gpt-3.5-turbo-0125* and  
15 *gpt-4-0125-preview*, referred to as GPT-3.5 and GPT-4.0 in later script. These models represent  
16 the most advanced version of their respective series and are provided by OpenAI through the  
17 API platform. Experiments were executed using Python scripts to interact with the OpenAI  
18 API. Each model received prompts via individual API requests without the conversation  
19 history, maintaining the independence of each interaction and preventing prior context from  
20 influencing the model's performance. All experiments were repeated five times to evaluate  
21 performance stability. Each of the selected papers would go through 80 experiments,  
22 including 2 GPT models (GPT3.5 and GPT4.0), 8 prompt sets, and 5 repetitive trials. In total,  
23 there were 8,000 experiments.

24 This design aims to yield a fair and thorough comparison of the two models, highlighting their  
25 respective strengths and limitations in processing and analysing medical research literature.

## 1 **Information Extraction**

2 Information extraction is a pivotal stage in a literature review. Not only does this facilitate the  
3 identification of related papers, it also has the potential to enhance the transparency of LLM's  
4 decision as an intermediate step in automatic literature screening. In this study, we identified  
5 seven important items in literature screening as representative samples, including sample size,  
6 data source, and PICOS (patient, intervention, comparison, outcomes, and study design). Their  
7 respective definitions are provided in Table 1.

8 We believe these elements to be the basis of efficient and precise literature screening, providing  
9 researchers with a clear, standardised framework for evaluation. Particularly, with the PICO, as  
10 the gold standard for clinical study assessment, this offers a systematic approach to identify  
11 relevant research questions and assess the quality of studies.

## 12 **Validation on the effects of Prompt Engineering**

13 Prompt engineering is an essential mechanism for optimizing the interaction with LLMs,  
14 serving to refine and enhance user queries in order to improve the task performance. In this  
15 section, we will examine and identify the effect of several prompt engineering strategies  
16 discussed in current directions of research, including *Adopting a Persona*, *Chain of Thought*  
17 and *Few-shot Learning*.

18 *Adopting a Persona* [22,23] is often achieved by instructing the LLMs to adopt the role of an  
19 expert in the related field of research. *Chain of Thought* [24,25] asks the model to explain the  
20 reasoning or the rationale behind each step in its problem-solving process. Although our task  
21 may not involve complicated logical reasoning, we are interested in investigating whether  
22 incorporating requests for justification could lead to improved performance and greater  
23 transparency. *Few-shot Learning* [26, 27] refers to the process in which we provide LLMs with  
24 expert output examples for similar tasks, which could serve as a guide for the model's  
25 responses.

26 Table 1. Definition of Information to extract [28] with example derived from of one  
27 publication [29].

Item	Definition	Example
Study design	Type of study	Retrospective study
Sample size	The number of participants involved in the study, and the basic characteristics of the participants.	17 patients received S-1 and 13 patients received SP
Data source	Source of the experimental data, such as databases, previous studies or surveys	Hong Kong cancer registry and tertiary centre
Patient	The patient involved in the experiment with some of their most important characteristics	Patients with metastatic gastric cancer who received either the S-1 or SP regimen as first line treatment

		for metastatic gastric cancer
Intervention	Main intervention, exposure, or prognostic factor in the experiment	S-1 regimen
Comparison	Main alternative group being considered	Compare between S-1 and S-1 plus cisplatin (SP)
Outcomes	The outcome that the experiment is trying to accomplish, measure, improve or affect	Baseline characteristics, overall response rate, median progression-free survival (PFS), overall survival (OS), and toxicity; and clinical benefit rate and 30-day mortality

1  
2 We adopted the following approach for our study. Firstly, we established a standard prompt  
3 without any specialised engineering strategy to serve as a control. This prompt simply asked  
4 the LLM to perform the task without additional instruction or context. We selected three  
5 prompt engineering strategies as mentioned previously. For each strategy, we crafted a series of  
6 prompts that incorporated the specific tactic. Following this, we then systematically removed  
7 one strategy at a time from the prompts, creating various ablated conditions for comparison  
8 against the baseline prompt and each other. For each prompt condition, we evaluated the  
9 LLM's performance using several metrics. Table 2 outlines the specific prompts that have been  
10 designed for each of these prompt engineering strategies.

11 Table 2. Prompt Setting for information extraction

Group	Prompt
Control	<p># Context  [a] You will be provided with titles and abstracts of medical papers, and your task is to parse it into structured data, including Study Design, Data Source, Sample Size, Patient, Intervention, Comparison and Outcomes, and separate them by semicolon.</p> <p># Input  [insert paper title and abstract]</p> <p># Instruction  Please read, extract and concisely report the following key details from the abstract:  Study Design: What type of methodology was employed in the study?  Sample Size: How many participants were included in the study?  Data Source: Where was the data for this study sourced from?  Based on the Study Design, if the paper is a review paper OR a laboratory study, please marks Patient, Intervention, Comparison and Outcomes all as NA.  Else, answer the following PICO question:  - Patients: Who is the study's targeted patient or population group?  - Intervention: What is the key intervention that the study assesses?</p>

	<p>- Comparison: Is there a comparison group or control used, and what does it consist of?</p> <p>- Outcomes: What outcomes are being measured to determine the intervention's success?</p> <p>Answer "NA" if any of the item is not mentioned in the abstract.</p> <p># Output</p> <p>Please [b1] output the structured data separated by semicolon, such as: [b2]</p> <p>Study Design: [output]; Data Source: [output]; Sample Size: [output]; Patient: [output]; Intervention: [output]; Comparison: [output]; Outcomes: [output]; [c]</p>
Strategy 1: Persona	<p># Inserted at [a]</p> <p>Imagine you are an expert in research methodology. Your role is essential in supporting a team of researchers by meticulously extracting critical information from medical paper abstracts. You have been trained to identify and collate specific elements that are crucial for the team's meta-analysis and database entry tasks.</p>
Strategy 2: Chain of thought	<p># Inserted at [b1]</p> <p>present a concise reasoning for each step you take, and how you arrive at the final structured data. Also, please</p> <p># Inserted at [b2]</p> <p>Reasoning: [output];</p> <p># Inserted at [b3]</p> <p>Reasoning: The abstract explicitly indicates that the study is a retrospective cohort study. The sample size is explicitly mentioned, consisting of three distinct groups with their respective counts. The data source is not explicitly named, so we mark it with NA. Since this is a cohort study (an epidemiological study) instead of a review paper or a laboratory study, we proceed with identifying the PICO elements. The patient population is women with PCOS, PCO, and age-matched controls undergoing IVF. The intervention is the IVF treatment itself. The comparison is made between the women with PCOS, those with PCO, and the age-matched controls. The outcomes being measured include various obstetric complications and outcomes such as GDM, GHT, PET, IUGR, gestation at delivery, baby's Apgar scores, and NICU admissions;</p>
Strategy 3: Few-shot Prompting	<p># Inserted at [c]</p> <p>Here is an example for your reference:</p> <p># Input</p> <p>Title: Obstetric outcomes in women with polycystic ovary syndrome and isolated polycystic ovaries undergoing in vitro fertilization: a retrospective cohort analysis</p> <p>Abstract: Objective: This retrospective cohort study evaluated the obstetric</p>



	<p>outcomes in women with polycystic ovary syndrome (PCOS) and isolated polycystic ovaries (PCO) undergoing in vitro fertilization (IVF) treatment. Methods: We studied 104 women with PCOS, 184 with PCO and 576 age-matched controls undergoing the first IVF treatment cycle between 2002 and 2009. Obstetric outcomes and complications including gestational diabetes (GDM), gestational hypertension (GHT), gestational proteinuric hypertension (PET), intrauterine growth restriction (IUGR), gestation at delivery, baby's Apgar scores and admission to the neonatal intensive care unit (NICU) were reviewed. Results: Among the 864 patients undergoing IVF treatment, there were 253 live births in total (25 live births in the PCOS group, 54 in the PCO group and 174 in the control group). The prevalence of obstetric complications (GDM, GHT, PET and IUGR) and the obstetric outcomes (gestation at delivery, birth weight, Apgar scores and NICU admissions) were comparable among the three groups. Adjustments for age and multiple pregnancies were made using multiple logistic regression and we found no statistically significant difference among the three groups. Conclusion: Patients with PCO±PCOS do not have more adverse obstetric outcomes when compared with non-PCO patients undergoing IVF treatment. © 2014 Informa UK Ltd. All rights reserved: reproduction in whole or part not permitted.</p> <p># Output [b3]</p> <p>Study Design: Retrospective cohort study; Sample Size: 864; Data Source: NA; Population: Women with polycystic ovary syndrome (PCOS) and isolated polycystic ovaries (PCO); Intervention: In vitro fertilization (IVF) treatment; Comparison: Age-matched controls; Outcomes: Obstetric complications (GDM, GHT, PET and IUGR) and the obstetric outcomes (gestation at delivery, birth weight, Apgar scores and NICU admissions);</p>
--	--

## 1 **Evaluation**

2 To evaluate the accuracy of the generated outcomes, we employed the established automatic  
3 metrics in NLP, including ROUGE-1 [30] and BERTScore [31]. These metrics were  
4 specifically designed to measure the quality of generated text compared to the reference text  
5 produced by human. ROUGE-N, a metric based on n-gram analysis, examined the overlap of  
6 common words and phrases between the two summaries. On the other hand, BERTScore  
7 encodes both the generated and reference texts using a pre-trained large language model to  
8 produce embeddings that capture the true semantic meaning of each text. The similarity is then  
9 calculated based on these embeddings. A more detailed explanation and relevant formulas are  
10 provided in Supplementary Material 2.

11 Unlike the N-gram (ROUGE-1) method that relies on exact matches, BERTScore can account  
12 for semantic similarities at the word and sentence level. In medical evidence extraction, this is  
13 particularly useful for evaluating complex medical terms and phrases that may have varied  
14 wording but similar meanings. For both metrics, we utilised the F1-score – which is the  
15 harmonic mean of the precision and recall scores that ranges from 0 to 1 – as our final standard

1 for analysis.

2 Noticeably, recent research papers highlighted the inherent challenges in assessing the LLM  
3 responses using traditional automatic metrics in NLP, such as ROUGE and BERTScore, which  
4 may lack sensitivity to nuanced semantic differences. The advanced language understanding  
5 and processing capabilities in LLMs may be beneficial to tackle this challenge. Therefore, we  
6 implemented an independent evaluation mechanism using a separate instance of GPT-4 model,  
7 specifically configured to assess the alignment between the generated responses correspond  
8 and the ground truth. This second GPT model was tailored by setting its temperature to 0,  
9 ensuring deterministic outputs for consistent evaluation, and by designing prompts to evaluate  
10 semantic similarity without relying on prior conversational history. This configuration allows  
11 the model to focus solely on the evaluation task, providing a more context-aware assessment  
12 that better captures subtle semantic alignment than traditional metrics. The detailed prompt  
13 could be found in Supplementary Material 2.

14

15 All three evaluators generated continuous metric ranging from 0 to 1, with distinct mean and  
16 standard deviation according to their different measurement on similarity. Therefore, we  
17 calibrated the evaluators using threshold to enable direct comparison among results by  
18 different evaluators. Specifically, we first created an *accordance dataset* by manually  
19 comparing the extraction results **from the two researchers** for each label in 100 papers. A  
20 score of 1 is assigned if the results matched (indicating agreement), while a score of 0 if they  
21 differed (indicating disagreement). This accordance dataset was solely used to calibrate the  
22 evaluators' threshold values, establishing a basis for measuring agreement consistency without  
23 influencing the final evaluation. During this calibration process, we calculated threshold values  
24 for the metric score produced by the evaluators across different element categories, in order to  
25 define what constitutes an acceptable level of agreement. Specifically, we iterate over the  
26 potential threshold value from 0 to 1 with a step size of 0.01 and assigned a "true" prediction  
27 for metric scores above the threshold, and false for scores below. Then, we determined which  
28 threshold would yield the highest accuracy rate of F1-score across all comparisons between the  
29 evaluators and the accordance ground truth and selected that as the eventual standard.

30 Finally, we used a **separate test set** to testify whether all the three evaluators calibrated on the  
31 accordance dataset is able to measure GPT-generated result with ground truth. We constructed  
32 the test set by randomly selecting 10 pairs of GPT-generated answers and corresponding  
33 ground truth labels across seven elements from various model and prompt combinations  
34 (GPT-3.5, GPT-4.0). We then manually assessed the alignment between each GPT-generated  
35 answer and the ground truth, which will then be served as the "true answer" for the test set.

36 To assess the overall performance of the models, we then applied the evaluators with the  
37 predefined thresholds to calculate the accuracy, sensitivity and specificity in GPT's  
38 information extraction results.

39 We defined the accuracy rate  $p_{\text{correct}}$  as the proportion of GPT's outputs that align with the  
40 ground truth in the five repetitive trials. It is calculated separately across the 100 papers as  
41 follows

1

$$\bar{p}_{correct} = \frac{1}{100} \sum_{i=1}^{100} \frac{\sum_{t=1}^5 1\{s_{t,i} > threshold_s\}}{5}$$

2 where  $s_{t,i}$  is the metric score for the  $i^{\text{th}}$  paper in trial  $t$ , and  $threshold_s$  is the threshold  
3 calculated for the specific element nature. The average  $p_{correct}$  was employed to horizontally  
4 compare the GPT models and prompt engineering strategies.

5 To address the risk of hallucination – producing information not grounded in the source  
6 material – and the possibility that not all elements of interest are present in a given abstract, we  
7 extended our metrics to include sensitivity and specificity. Sensitivity measures the proportion  
8 of correct information being extracted, while specificity measures the proportion of irrelevant  
9 information being discarded. To evaluate the model's performance in handling hallucination,  
10 we expect to see high specificity to avoid any misleading information. High sensitivity is also  
11 important to indicate all necessary information has been involved. For clarity, we categorized  
12 an element as *positive* if it was correctly identified and labelled from the abstract; otherwise, it  
13 was categorized as *negative*. Detailed definitions are found in Supplementary Material 2.

#### 14 **Statistical analysis**

15 For each extracted item evaluated by one metric, a 2-way Analysis of Variance (ANOVA)  
16 model was used to analyse the impact of two factors, GPT versions and prompt engineering  
17 strategies. We summarized all  $P$  values across items and evaluators in one table, to analyse the  
18 significance of the GPT model and prompts effects on the performance. Statistical analysis was  
19 performed using the python package *statsmodels* (version 0.14.1) [31]. All significance levels  
20 were set as 0.05, with all necessary assumptions for ANOVA, including normality and  
21 homogeneity of variances, being assessed, and satisfied.

## 22 **Results**

### 23 **Paper selection and Data source**

24 Figure 2 illustrates the characteristics and distribution of the sampled publications. These  
25 scholarly articles were collected from nineteen departments within the Faculty of Medicine at  
26 the University of Hong Kong, signifying a wide coverage of medical domains. The collection  
27 encompasses various research fields, from broad disciplines, such as surgery, medicine, and  
28 public health, to more specialised areas, such as emergency medicine, Chinese medicine, and  
29 paediatric and adolescent medicine.

30 The selected publications also provided comprehensive coverage across study designs. The  
31 labelled ground truth indicated that the dataset consisted of 22 retrospective studies, 13  
32 laboratory studies, 10 prospective studies, 7 case reports, 5 reviews, 4 randomised controlled  
33 trials, and other types of study design. Among these study types, review and laboratory study  
34 did not include elements like sample size, patient, intervention, comparison, and outcomes.

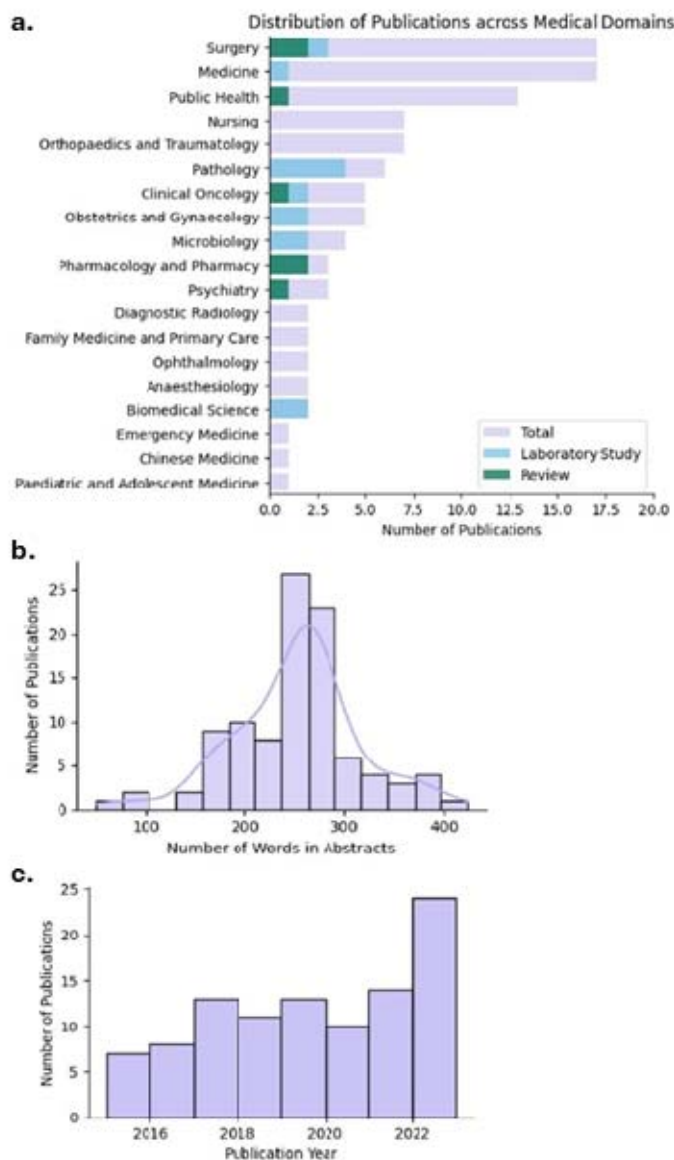
1 The ground truth labelled these elements as not applicable. Figure 2 also indicated the distinct  
2 proportion of these two study designs across medical domains.

3 The distribution of publication year in general exhibits a uniform pattern from 2015 to 2023  
4 with an ascending trend in the recent years. The length of the abstract adheres to a normal  
5 distribution, with a mean length of 252 words.

6

7 Figure 2. Characteristics of sampled publications.

8 (a) An illustration of the number of publications in each medical domains, with proportion of  
9 laboratory study and literature review indicated by different colours. (b) The distribution of the  
10 number of words in abstract as input. (c) The bar plot of number of publications in each year,  
11 from 2015 to 2023.



1

## 2 Evaluator Performance

3 For testing the evaluators, we randomly selected the content and metric scores from three  
4 evaluators of 10 paper \* 7 elements from different combination of models (GPT-3.5, GPT-4.0),  
5 prompt types and trials, and manually marked down the accordance between extracted  
6 information and ground truth. This test set was independent of the accordance dataset used for  
7 threshold calibration. Table 3 presents performance metrics for three different evaluators to  
8 assess the quality of semantic similarity rating, and they are compared based on their accuracy,  
9 precision, and recall. Detailed samples and results are found in Supplementary Material 3. With  
10 all evaluators showing generally good results (accuracy, precision, and recall all above 0.94),  
11 the LLM Evaluator demonstrated the highest score across all metrics, indicating robust  
12 performance in evaluating information alignment.

## 1 Overall Performance

2 In our experiments, GPTs achieved considerable accuracy in extracting information from  
3 papers across medical disciplines. Measured with the ROUGE-1 Score and LLM Evaluator,  
4 GPT-4.0 achieved over 80% correctness in six out of the seven items with the optimal prompt  
5 engineering strategy. Supplementary Material 4 includes a comprehensive table summarising  
6 the average proportions of correctness, covering all 7 items under 8 prompt settings, generated  
7 by GPT-3.5 and GPT-4.0 and measured by the three different metrics.

8 Table 3. Accuracy, precision, and recall of evaluators.

	BERT	ROUGE-1	ChatGPT-4.0
Accuracy	0.94286	0.95714	0.97143
Precision	0.95082	0.98276	0.98305
Recall	0.98305	0.96610	0.98305

9

10 The performance of GPT can be stratified into three levels, corresponding to three distinct  
11 degrees of complexity among the seven information extraction tasks. The first level  
12 encompasses questions where a direct answer can typically be found in the raw text. The  
13 sample size is an example of this level, and both GPT-3.5 and GPT-4.0 achieve accuracy levels  
14 exceeding 0.95 in extracting sample size. The second level pertains to questions requiring  
15 understanding and summarisation skills to extract answers. Most extracted items, including  
16 study design, data source, patient, comparison, and outcomes, belong to this category. Figure 3  
17 shows that GPT-3.5 achieves optimal performance from 0.7 to 0.8 for these items and GPT-4.0  
18 from 0.8 to 0.9. Finally, intervention represents the third level, which demands a high level of  
19 understanding and domain expertise to discern the correct answer accurately from potentially  
20 misleading information. In this regard, GPT-3.5 performed under 0.6 while GPT-4.0  
21 demonstrated accuracy around 0.7.

22 Noticeably, both GPT-3.5 and GPT-4.0 demonstrated stability in information extraction. In  
23 Figure 3, all empirical distributions of  $p_{\text{correct}}$  reveal a bimodal pattern, with performance  
24 clustering at high and low accuracy extremes, indicating that the GPT models are either all  
25 correct or all incorrect in their extractions.

## 26 Sensitivity and Specificity

27 Besides calculating direct average accuracy, we also summarised the overall sensitivity and  
28 specificity scores of each model and prompt strategy types measured by 3 different evaluators  
29 in Supplementary Material 4. Figure 4 is a visual comparison of accuracy, sensitivity and  
30 specificity of GPT across eight prompt designs based on the result from the most reliable  
31 evaluator (LLM Evaluator). GPT-4.0 has an average sensitivity score of 0.8550 while scoring  
32 only 0.7353 in specificity. This difference is more distinct in GPT-3.5, with sensitivity 0.8147  
33 and specificity of 0.5671.

## 1 Comparing the Performance of GPT-3.5 and GPT-4.0

2 ANOVA analysis results that the GPT version is a statistically significant factor influencing  
 3 model performance. As presented in Table 4, the ANOVA analysis revealed that 19 out of the *P*  
 4 values assessing the impact of GPT were significantly lower than 0.05. The only two  
 5 exceptions on *P* value were associated with Study Design and Sample Size measured by  
 6 BERTScore, which may relate to the low accuracy of BERTScore mentioned above.

7 Table 4. Summary of *P* values in ANOVA analysis.

Evaluator	Factor	Study Design	Sample Size	Data Source	Patient	Intervention	Comparison	Outcomes
ROUGE	GPT	< .001 <sup>a</sup>	.01	.008	< .001	< .001	< .05	< .001
ROUGE	Prompt	.073	< .001	< .001	.446	.943	< .001	.995
ROUGE	Interaction	.007	.006	.899	.753	.991	.815	.997
BERT	GPT	.342	.890	.001	< .001	< .001	< .001	< .001
BERT	Prompt	.02	< .001	< .001	.991	.947	.001	.656
BERT	Interaction	.002	< .001	.924	.995	.968	.651	.194
GPT	GPT	.004	< .001	.002	< .001	< .001	.002	< .001
GPT	Prompt	.01	< .001	< .001	.331	.869	.126	.122
GPT	Interaction	.087	< .001	.892	.716	.999	.497	.233

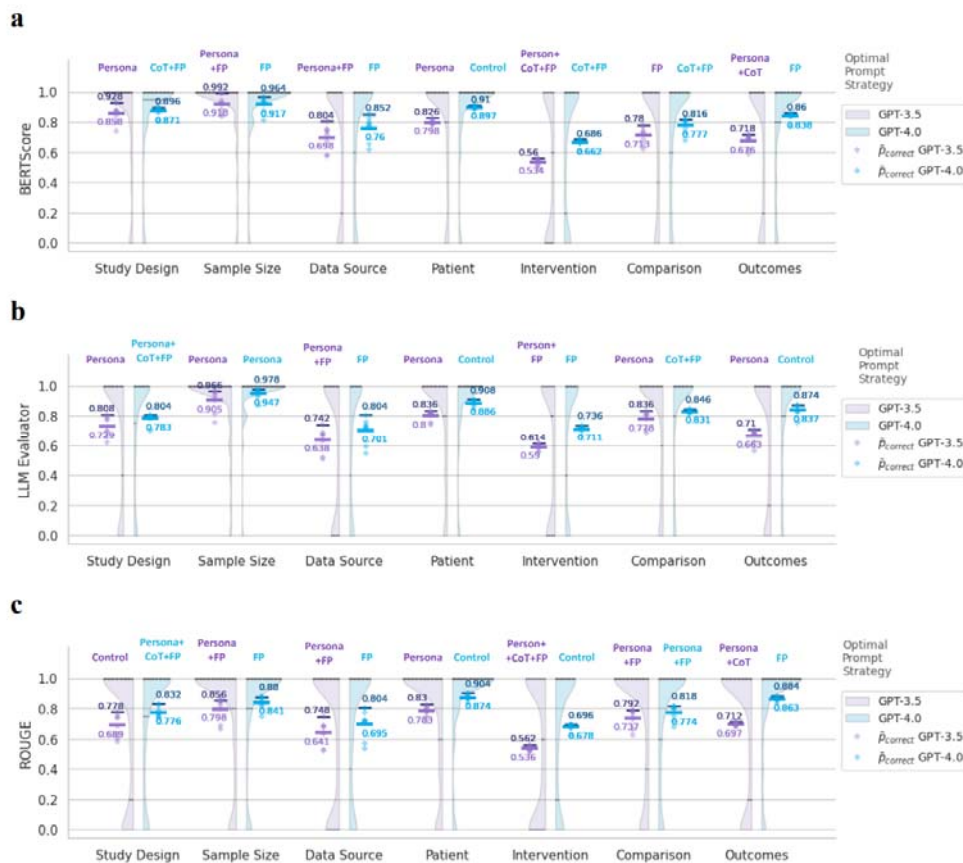
8 a: For example, the first cell represents the *P* value corresponding to the factor “GPT versions”, utilizing study  
 9 design data evaluated by the ROUGE metric as input data for the ANOVA analysis.

10

## 11 Effects of Prompt Engineering Strategies

12 Prompt engineering strategies are likely to influence model performance positively. As  
 13 presented in Table 4, the ANOVA analysis revealed that the impact of the GPT prompt was  
 14 statistically significant for two extracted items, Sample Size, Data Source, measured by all  
 15 three evaluators. There needs to be more evidence for other items to prove the impact of  
 16 prompt engineering strategies. It is also noticeable that prompt engineering strategies may not  
 17 have additive effects with each other. For example, in Figure 4, the combination Persona +  
 18 Chain-of-Thought did not perform as well as either Persona or Beta. Combined strategies, such  
 19 as Persona + Chain-of-Thought + Few-shot Prompting, could lead to inferior results compared  
 20 to a single strategy.

21 The effects of GPT versions and prompt engineering strategies will likely interact. In ANOVA  
 22 analysis, the interaction between the GPT version and prompt engineering strategies was  
 23 statistically significant based on the Sample Size extraction, as assessed by all three  
 24 evaluators (ROUGE, *P* < .001; BERTScore, *P* < .001; LLM Evaluator, *P* < .001). However, for  
 25 other items, interaction may exist but needs more statistical strength. Figure 3 indicates that  
 26 GPT-3.5 tended to favour the Persona strategy, persona, while GPT-4.0 tended to prefer the  
 27 few-shot prompting. Chain of thought, was relatively less effective in the information  
 28 extraction task.

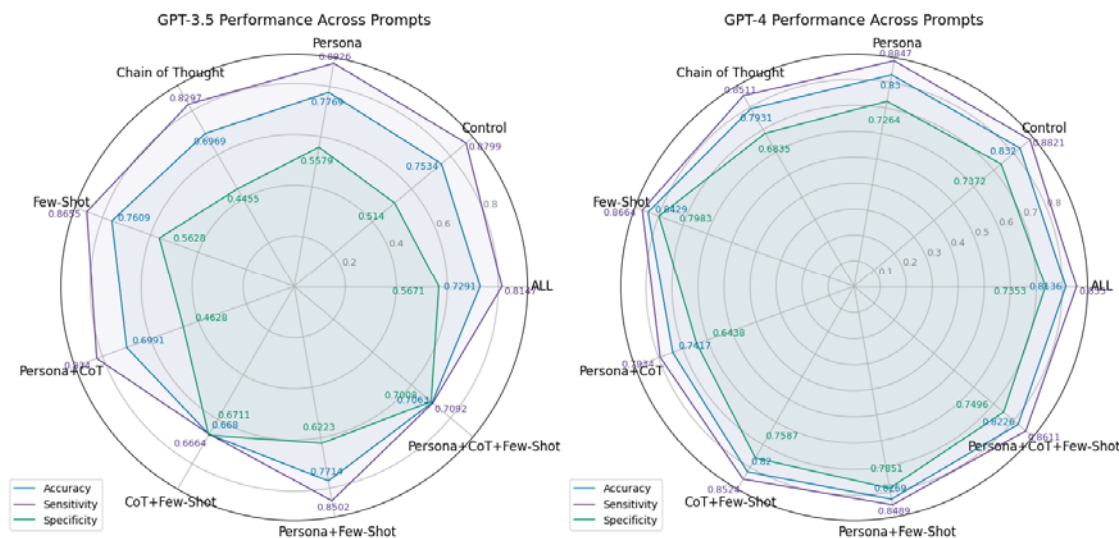


1  
2  
3  
4  
5  
6  
7  
8  
9  
10  
11  
12  
13

Figure 3. Distribution of Information Extraction Accuracy.

Violin plots illustrate the empirical distribution of  $p_{correct}$  across seven medical items, and three prompt engineering strategies, Persona, Chain of Thought (CoT), Few-shot Prompting (FP).  $p_{correct}$  denotes the proportion of correct answers from five repeated trials per paper. Each distribution aggregates 800  $p_{correct}$  values via kernel density estimation. Diamond markers represent  $\bar{p}_{correct}$  for each prompt strategy, where  $\bar{p}_{correct}$  denotes the mean of  $p_{correct}$  among the 100 publications. Furthermore, the highest  $\bar{p}_{correct}$  are marked in dark color with the corresponding optimal prompt strategy highlighted above each column. Both mean and maximum  $\bar{p}_{correct}$  are depicted using bar marker.





1

2 Figure 4. Accuracy, Sensitivity, and Specificity of Information Extraction across Prompt  
3 Strategies.

4 “ALL” represents the mean value across all prompt designs.

5 All values are measured by the LLM Evaluators.

## 6 Discussion

7 Our research pioneers the exploration of a new generation of LLMs in medical evidence  
8 summarisation and offers potential applications in various scenarios. It provides empirical  
9 evidence to support the development of credible automatic tools for medical literature  
10 screening and review. With critical information extracted, automatic tools can strike a balance  
11 between efficiency and transparency. To ensure comprehensive coverage across various  
12 medical domains, a stratified sampling method was adopted for paper selection from almost all  
13 affiliated medical schools and departments of a university. Furthermore, we employed multiple  
14 evaluators, repetitive trials, and experiments on prompt engineering strategies in the  
15 experiment to enhance the integrity of results. Our findings demonstrated that GPTs can  
16 effectively extract or summarise information described in the abstracts. Notably, GPT-4.0  
17 exhibits robust performance in providing thorough answers and understanding and  
18 summarising abstracts. However, there is still room for improvement in accurately discerning  
19 information that requires sophisticated understanding and domain expertise. When combined  
20 with appropriate prompt engineering strategies, the accuracy level achieves over 0.8 in  
21 extracting information related to study design, sample size, data source, patient, comparison,  
22 and outcomes.

23 We observed that GPT has displayed different levels of performance in extracting information  
24 across the seven items. This may be due to the varying complexities involved in the  
25 information extraction tasks. The first level encompasses questions where a direct answer can  
26 typically be found in the raw text. The sample size is an example of this level, and both  
27 GPT-3.5 and GPT-4.0 achieve accuracy levels exceeding 0.95 in extracting sample size. The  
28 second level pertains to questions requiring understanding and summarisation skills to extract

1 answers. Most extracted items, including study design, data source, patient, comparison, and  
2 outcomes, belong to this category. Finally, intervention represents the third level, which  
3 demands a high level of understanding and domain expertise to discern the correct answer  
4 accurately from potentially misleading information. In this regard, GPT-3.5 performed under  
5 0.6 while GPT-4.0 demonstrated accuracy around 0.7.

6 The field of large language models is rapidly advancing. Our investigations reveal that the  
7 effect of the GPT version on the accuracy of information extraction is significant (Table 4).  
8 GPT-4.0 presents a more robust performance in summarising complex information that may  
9 not be readily apparent in the raw text, such as the PICOS. The increase in accuracy is mainly  
10 driven by a significant improvement in specificity, the ability to discard irrelevant information,  
11 which align with observations in previous research [15] On the other hand, the drawback of  
12 GPT-4.0 compared to its predecessor is associated with time and cost. According to the  
13 OpenAI website, by March 2024, the price of GPT-3.5 Turbo was one-twentieth that of  
14 GPT-4.0 Turbo [33]. In our experiment, we found that the time required for GPT-3.5 to label  
15 100 papers is approximately one-tenth of the time taken by GPT-4.0. This significant difference  
16 may be attributed to the rate limits imposed by the API, as noted on OpenAI's website, the rate  
17 limit for GPT-4-turbo is 500 RPM (Requests Per Minute) for Tier 1 users, while GPT-3.5-turbo  
18 offers a higher rate limit of 3500 RPM [34]. Both the two models mark an improvement in  
19 efficiency compared to human labour, by reducing 8 to 10 hours of labelling to around 5  
20 minutes (GPT-3.5) or 40 minutes (GPT-4.0) in our experiments.

21 Prompt engineering strategies play an essential role in enhancing LLMs' performance. We  
22 found that the optimal prompt engineering strategies vary depending on the extraction tasks  
23 and GPT versions employed. Overall, two useful strategies are recommended to attempt:  
24 persona and few-shot prompting. Although the chain of thought strategy might help guide  
25 multi-step tasks, it might not be effective in straightforward tasks like the information  
26 extraction in this study. Further, the few-shot prompting strategy may improve the overall  
27 accuracy by raising the specificity scores. This is because the incorporation of examples  
28 labelled as 'NA' in the prompts can likely guide the GPT model in recognising and categorising  
29 non-applicable instances more accurately, leveraging the model's predictive nature to enhance  
30 overall accuracy in information extraction tasks. Interestingly, it is worth noting that the  
31 combination of prompt engineering strategies may not yield additive effects on the final results.  
32 Considering the cost associated with input tokens, a conservative approach is recommended to  
33 employ prompt engineering strategy in solving simple medical information extraction.

34 Noteworthy, we also identified the overall higher sensitivity score contrast to specificity score,  
35 as recorded in Figure 4 and Table 3 of Supplementary Material 4. Specificity, the ability to  
36 avoid hallucination, is the weakest compared to the other metrics, sensitivity and accuracy,  
37 representing the ability of extracting correct information, and overall accuracy. GPT-4.0  
38 outperforms GPT-3.5 significantly in reducing the risk of hallucination. However, this  
39 consistently higher sensitivity might be a result of the imbalance of dataset. Since the dataset  
40 have more element identified in the abstract and less labelled as Not Applicable (NA), naturally  
41 there will be fewer number in the negative class. Therefore, any misclassification will have a  
42 disproportionately large impact on the specificity measurement, making the metric highly  
43 sensitive to the model's performance on a small number of cases. Also, the nature of the GPT

1 model might also play a role in this result. As a generative transformer, GPTs generate text  
2 based on the probability of the next word or phrase. In tasks that require extraction from texts,  
3 this nature might make them inherently more inclusive in their responses, favouring sensitivity.  
4 In brief, we suggest GPT performs better in tasks focus more on reducing false positive. A  
5 more cautious attitude is recommended when applying GPT to tasks that are vulnerable to  
6 hallucination, in particular with the older version of GPT.

7 Moreover, this study extensively examines and compares the performance of evaluators  
8 utilised in the experiment, including two well-established NLP metrics, ROUGE-1 and  
9 BERTScore, and one newly developed LLM evaluator. Overall, the three evaluators provide  
10 consistent performance evaluation across various extraction items and prompt engineering  
11 strategies. Our study also revealed an interesting observation regarding the potential of GPT as  
12 a promising and unique tool to assess the accuracy of generated text compared to the ground  
13 truth. Notably, LLM evaluators can leverage their pre-trained knowledge base to evaluate text  
14 based not only on lexical similarity but also on semantic similarity. This ability effectively  
15 addresses some significant limitations of existing NLP metrics.

16 Our study also has limitations. First, while we attempted to cover a wide range of medical  
17 domains within a hundred papers, each specific medical domain might be under-sampled.  
18 Moreover, when the targeted literature focuses on one area, domain knowledge can be  
19 provided as contextual information to enhance performance. Thus, future research could  
20 validate GPT's performance practically on one specific medical domain. Another limitation of  
21 this study is that we solely tested GPT from the abstracts. Given the proliferating capability of  
22 LLMs in handling long text, figures, and tables, it is recommended that future researchers  
23 extend the GPT tools to operate on full text or the PDF level. This expansion would extract  
24 more valuable information sources and open up broader possibilities for GPT to facilitate  
25 medical research.

## 26 **Conclusion**

27 GPT has been demonstrate notable accuracy in clinical text summarization [13], our study  
28 further showcases that GPT can be a stable and reliable tool for information extraction from  
29 titles and abstracts of literature across multiple medical domains. Both GPT versions and  
30 prompt engineering strategies will impact the accuracy of GPT's output. Conservative prompt  
31 strategy is recommended for simple information extraction tasks, and latest versions of GPT  
32 for tasks that are vulnerable to hallucination. Further investigation is needed to assess and  
33 improve LLM's performance in extracting complex or professional information. We  
34 encourage more research and studies to continue refining and advancing this tool, unlocking  
35 the potential of the new generation of technology in medical research.

## 36 **Data Availability**

37 All the data and codes of this study will be available for open access after publication.

## 1 **Acknowledgement**

2 We extend our heartfelt thanks to Professor Wanling Yang, Dr. Ching Lung Cheung, Dr.  
3 Joshua Wing Kei Ho, and Dr. Eric Yuk Fai Wan for their insights and perspectives on this  
4 study. We also thank Ms Lisa Lam for proof-reading of this paper.

5 This study received no funding.

## 6 **Declaration of Interest**

7 Xue Li received research grants from the Hong Kong Health and Medical Research Fund  
8 (HMRF, HMRF Fellowship Scheme, HKSAR), Research Grants Council Early Career Scheme  
9 (RGC/ECS, HKSAR), Janssen, and Pfizer; internal funding from the University of Hong Kong;  
10 and consultancy fees from Merck Sharp & Dohme and Pfizer; she is also a non-executive  
11 director of Advanced Data Analytics for Medical Science (ADAMS) Limited Hong Kong; all  
12 are unrelated to this work.

## 13 **Reference**

- 14 1. Li J, Dada A, Puladi B, Kleesiek J, Egger J. ChatGPT in healthcare: A taxonomy and  
15 systematic review. *Comput Methods Programs Biomed.* 2024 Mar;245:108013. doi:  
16 10.1016/j.cmpb.2024.108013. Epub 2024 Jan 15. PMID: 38262126.
- 17 2. Lim ZW, Pushpanathan K, Yew SME, Lai Y, Sun CH, Lam JSH, Chen DZ, Goh JHL, Tan  
18 MCJ, Sheng B, Cheng CY, Koh VTC, Tham YC. Benchmarking large language models'  
19 performances for myopia care: a comparative analysis of ChatGPT-3.5, ChatGPT-4.0, and  
20 Google Bard. *EBioMedicine.* 2023 Sep;95:104770. doi: 10.1016/j.ebiom.2023.104770. Epub  
21 2023 Aug 23. PMID: 37625267; PMCID: PMC10470220.
- 22 3. Shah NH, Entwistle D, Pfeffer MA. Creation and Adoption of Large Language Models in  
23 Medicine. *JAMA.* 2023 Sep 5;330(9):866-869. doi: 10.1001/jama.2023.14217. PMID:  
24 37548965.
- 25 4. Nori H, King N, McKinney SM, Carignan D, Horvitz E. Capabilities of gpt-4 on medical  
26 challenge problems. *arXiv preprint arXiv:230313375.* 2023.4
- 27 5. Tang L, Sun Z, Idnay B, Nestor JG, Soroush A, Elias PA, Xu Z, Ding Y, Durrett G, Rousseau  
28 JF, Weng C, Peng Y. Evaluating large language models on medical evidence summarization.  
29 *NPJ Digit Med.* 2023 Aug 24;6(1):158. doi: 10.1038/s41746-023-00896-7. PMID: 37620423;  
30 PMCID: PMC10449915.
- 31 6. Shaib C, Li ML, Joseph S, Marshall IJ, Li JJ, Wallace BC. Summarizing, simplifying, and  
32 synthesizing medical evidence using GPT-3 (with varying success). *arXiv preprint*  
33 *arXiv:230506299.* 2023.
- 34 7. Tian S, Jin Q, Yeganova L, Lai PT, Zhu Q, Chen X, Yang Y, Chen Q, Kim W, Comeau DC,  
35 Islamaj R, Kapoor A, Gao X, Lu Z. Opportunities and challenges for ChatGPT and large

- 1 language models in biomedicine and health. *Brief Bioinform.* 2023 Nov 22;25(1):bbad493.  
2 doi: 10.1093/bib/bbad493. PMID: 38168838; PMCID: PMC10762511.
- 3 8. O'Mara-Eves A, Thomas J, McNaught J, Miwa M, Ananiadou S. Using text mining for study  
4 identification in systematic reviews: a systematic review of current approaches. *Syst Rev.*  
5 2015 Jan 14;4(1):5. doi: 10.1186/2046-4053-4-5. Erratum in: *Syst Rev.* 2015 Apr 28;4:59. doi:  
6 10.1186/s13643-015-0031-5. PMID: 25588314; PMCID: PMC4320539.
- 7 9. Shemilt I, Khan N, Park S, Thomas J. Use of cost-effectiveness analysis to compare the  
8 efficiency of study identification methods in systematic reviews. *Syst Rev.* 2016 Aug  
9 17;5(1):140. doi: 10.1186/s13643-016-0315-4. PMID: 27535658; PMCID: PMC4989498.
- 10
- 11 10. Matsui K, Utsumi T, Aoki Y, Maruki T, Takeshima M, Yoshikazu T. Large language model  
12 demonstrates human-comparable sensitivity in initial screening of systematic reviews: A  
13 semi-automated strategy using gpt-3.5. Available at SSRN: <https://ssrn.com/abstract=4520426>.  
14 4520426. 2023.
- 15 11. Mahuli SA, Rai A, Mahuli AV, Kumar A. Application ChatGPT in conducting systematic  
16 reviews and meta-analyses. *Br Dent J.* 2023 Jul;235(2):90-92. doi:  
17 10.1038/s41415-023-6132-y. PMID: 37500847.
- 18 12. Hill JE, Harris C, Clegg A. Methods for using Bing's AI-powered search engine for data  
19 extraction for a systematic review. *Res Synth Methods.* 2024 Mar;15(2):347-353. doi:  
20 10.1002/jrsm.1689. Epub 2023 Dec 8. PMID: 38066713.
- 21 13. Van Veen, D., Van Uden, C., Blankemeier, L., Delbrouck, J. B., Aali, A., Bluethgen, C., ... &  
22 Chaudhari, A. S. (2024). Adapted large language models can outperform medical experts in  
23 clinical text summarization. *Nature medicine*, 30(4), 1134-1142.
- 24 14. Landschaft, A., Antweiler, D., Mackay, S., Kugler, S., Rüping, S., Wrobel, S., ... &  
25 Allende-Cid, H. (2024). Implementation and evaluation of an additional GPT-4-based  
26 reviewer in PRISMA-based medical systematic literature reviews. *International Journal of*  
27 *Medical Informatics*, 189, 105531.
- 28 15. Yang, J., Walker, K. C., Bekar-Cesaretli, A. A., Hao, B., Bhadelia, N., Joseph-McCarthy, D.,  
29 & Paschalidis, I. C. (2024). Automating biomedical literature review for rapid drug discovery:  
30 Leveraging GPT-4 to expedite pandemic response. *International Journal of Medical*  
31 *Informatics*,
- 32 16. Sivarajkumar, S., Kelley, M., Samolyk-Mazzanti, A., Visweswaran, S., & Wang, Y.  
33 (2024). An empirical evaluation of prompting strategies for large language models in  
34 zero-shot clinical natural language processing: algorithm development and validation  
35 study. *JMIR Medical Informatics*, 12, e55318.
- 36 17. Akyon, S. H., Akyon, F. C., Camyar, A. S., Hızlı, F., Sari, T., & Hızlı, Ş. (2024).  
37 Evaluating the Capabilities of Generative AI Tools in Understanding Medical Papers:  
38 Qualitative Study. *JMIR Medical Informatics*, 12(1), e59258.
- 39
- 40 18. Gilbert S, Harvey H, Melvin T, Vollebregt E, Wicks P. Large language model AI chatbots

- 1 require approval as medical devices. *Nat Med.* 2023 Oct;29(10):2396-2398. doi:  
2 10.1038/s41591-023-02412-6. PMID: 37391665.
- 3 19. Meskó, B. (2023). Prompt engineering as an important emerging skill for medical  
4 professionals: tutorial. *Journal of medical Internet research*, 25, e50638.
- 5 20. Zagher, J., Naguib, M., Bjelogrić, M., Névéol, A., Tannier, X., & Lovis, C. (2024). Prompt  
6 Engineering Paradigms for Medical Applications: Scoping Review. *Journal of Medical*  
7 *Internet Research*, 26, e60501.
- 8 21. Wei, Q., Yao, Z., Cui, Y., Wei, B., Jin, Z., & Xu, X. (2024). Evaluation of ChatGPT-generated  
9 medical responses: a systematic review and meta-analysis. *Journal of Biomedical Informatics*,  
10 104620.
- 11 22. Grabb D. The impact of prompt engineering in large language model performance: a  
12 psychiatric example. *Journal of Medical Artificial Intelligence.* 2023;6.
- 13 23. Kojima T, Gu SS, Reid M, Matsuo Y, Iwasawa Y. Large language models are zero-shot  
14 reasoners. *Advances in neural information processing systems.* 2022;35:22199-213.
- 15 24. Wei J, Wang X, Schuurmans D, Bosma M, Xia F, Chi E, Le Q, Zhou D. Chain-of-thought  
16 prompting elicits reasoning in large language models. *Advances in neural information*  
17 *processing systems.* 2022;35:24824-37.
- 18 25. Wang X, Wei J, Schuurmans D, Le Q, Chi E, Narang S, Chowdhery A, Zhou D.  
19 Self-consistency improves chain of thought reasoning in language models. *arXiv preprint*  
20 *arXiv:220311171.* 2022.
- 21 26. Zhao Z, Wallace E, Feng S, Klein D, Singh S. Calibrate Before Use: Improving Few-shot  
22 Performance of Language Models. In: Marina M, Tong Z, editors. *Proceedings of the 38th*  
23 *International Conference on Machine Learning; Proceedings of Machine Learning Research:*  
24 *PMLR;* 2021. p. 12697--706.
- 25 27. Brown T, Mann B, Ryder N, Subbiah M, Kaplan JD, Dhariwal P, Neelakantan A, Shyam P,  
26 Sastry G, Askell A, Agarwal S, Herbert-Voss A, Krueger G, Henighan T, Child R, Ramesh A,  
27 Ziegler D, Wu J, Winter C, Hesse C, Chen M, Sigler E, Litwin M, Gray S, Chess B, Clark J,  
28 Berner C, McCandlish S, Radford A, Sutskever I, Amodei D. Language models are few-shot  
29 learners. *Advances in neural information processing systems.* 2020;33:1877-901.
- 30 28. Duke University. *LibGuides: Evidence-Based Practice: PICO: Duke University; 2019* [cited  
31 2024 April 20]. Available from: <https://guides.mclibrary.duke.edu/ebm/pico>.
- 32 29. Lau KS, Lam KO, Chan WLW, Lee VHF, Kwong DLW, Leung TW. S-1 Versus S-1 Plus  
33 Cisplatin as First-line Treatment for Metastatic Gastric Cancer. *Hong Kong Journal of*  
34 *Radiology.* 2017;20:318. doi: 10.12809/hkjr1716810.
- 35 30. Lin C. ROUGE: A Package for Automatic Evaluation of Summaries. *Text Summarization*  
36 *Branches Out: Association for Computational Linguistics;* 2004. p. 74-81.
- 37 31. Zhang T, Kishore V, Wu F, Weinberger KQ, Artzi Y. Bertscore: Evaluating text generation  
38 with bert. *arXiv preprint arXiv:190409675.* 2019.

- 1        32. Seabold S, Perktold J. Statsmodels: Econometric and Statistical Modeling with Python.  
2            Proceedings of the 9th Python in Science Conference. 2010;2010.
- 3        33. OpenAI. Pricing: OpenAI; 2024 [cited 2024 April 20]. Available from:  
4            <https://openai.com/pricing>.
- 5        34. OpenAI. Rate Limits: OpenAI; 2024 [cited 2024 April 20]. Available from:  
6            <https://platform.openai.com/docs/guides/rate-limits/usage-tiers?context=tier-one>.
- 7
- 8
- 9
- 10
- 11

## 1 **Abbreviation**

2 LLM: Large language model

3 NLP: Natural language processing

4 PICOS: patient, intervention, comparison, outcomes, and study design



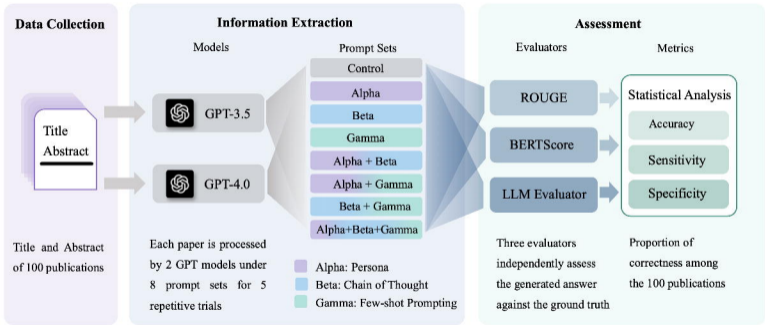


Figure 2. Characteristics of sampled publications.

(a) An illustration of the number of publications in each medical domains, with proportion of laboratory study and literature review indicated by different colours. (b) The distribution of the number of words in abstract as input. (c) The bar plot of number of publications in each year, from 2015 to 2023.

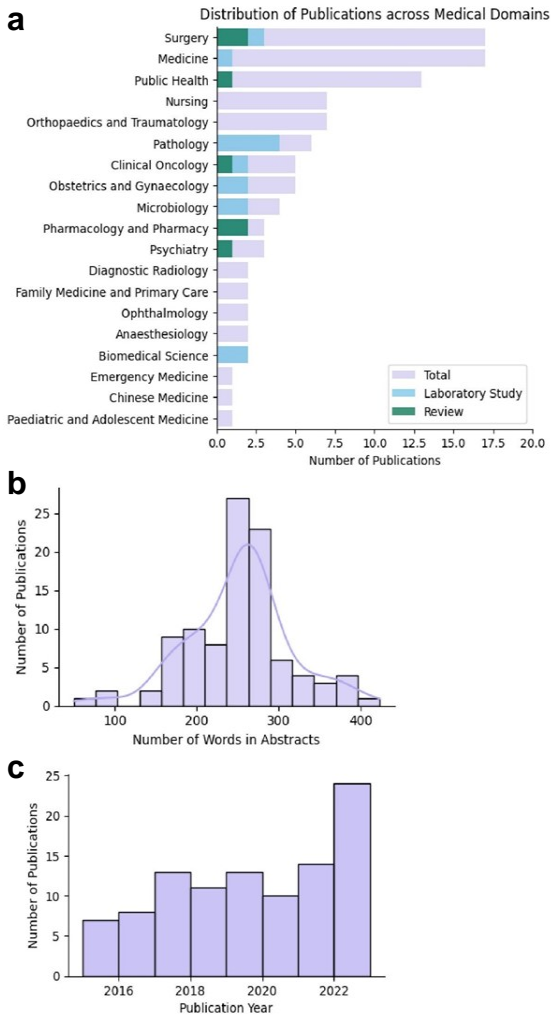
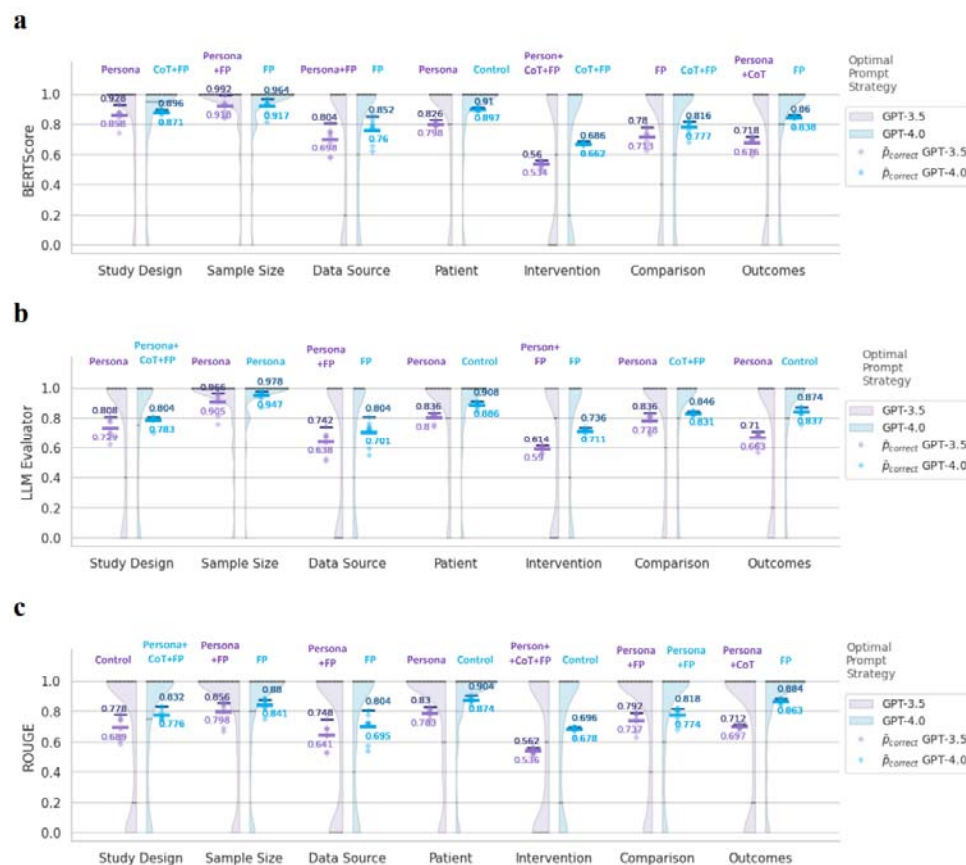
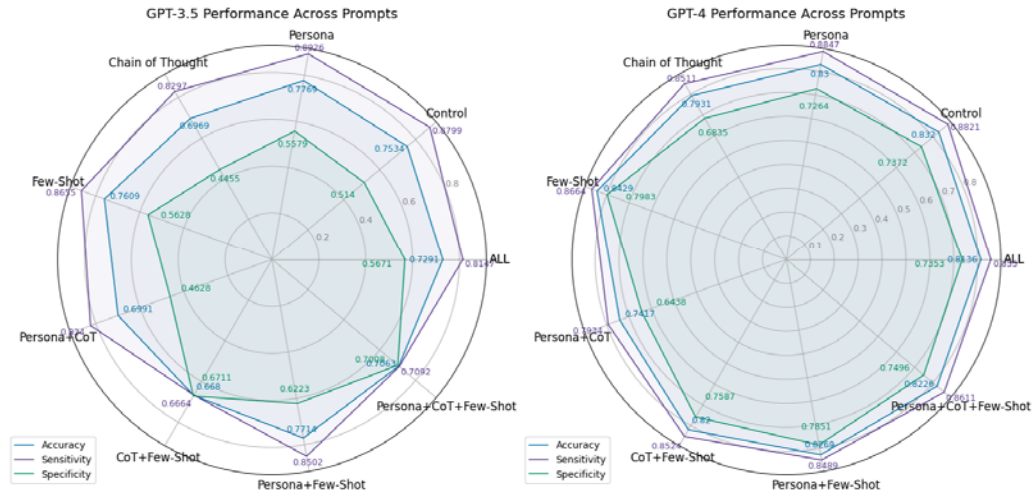


Figure 3. Distribution of information extraction accuracy



Violin plots illustrate the empirical distribution of  $p_{correct}$  for 100 papers across seven medical items, and three prompt engineering strategies, Persona, Chain of Thought (CoT), Few-shot Prompting (FP).  $p_{correct}$  denotes the proportion of correct answers from five repeated trials per paper. Each distribution aggregates 800  $p_{correct}$  values via kernel density estimation. Diamond markers represent  $\bar{p}_{correct}$  for each prompt strategies, where  $\bar{p}_{correct}$  denotes the mean of  $p_{correct}$  among the 100 publications. Furthermore, the highest  $\bar{p}_{correct}$  are marked in dark color with the corresponding optimal prompt strategy highlighted above each column. Both mean and maximum  $\bar{p}_{correct}$  are depicted using bar marker.

**Figure 4.** Accuracy, Sensitivity and Specificity of Information Extraction across Prompt Strategies



“ALL” represents the mean value across all prompt design.

All values are measured by the LLM Evaluators.