

# Benchmarking Large Language Models in Evidence-Based Medicine

Jin Li , Yiyang Deng , Qi Sun , Junjie Zhu , Yu Tian , Jingsong Li , and Tingting Zhu 

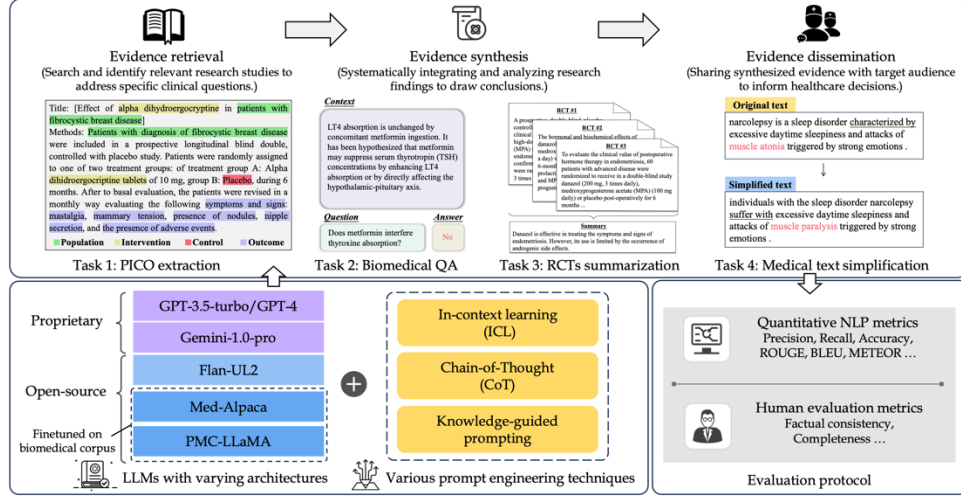


Fig. 1: The workflow of this study.

TABLE VII: An overview of error types.

Category	Source	Example
Factual contradiction	Reference summary	Post extubation atelectasis may be more frequent after <b>nasal intubation</b> , particularly in very low birth weight infants. One route of intubation does not seem to be preferable to the other...
	Output summary	<b>Flan-UL2-20b</b> : The authors compared the incidence of post-extubation atelectasis in patients intubated with oral versus nasal endotracheal tubes. The incidence of atelectasis was significantly higher in patients <b>intubated with oral tubes</b> .
Certainty Illusion	Reference summary	Metronidazole, given as a single dose, is likely to provide a parasitological cure for trichomoniasis, but it is <b>not known</b> whether this treatment will have any effect on pregnancy outcomes...
	Output summary	<b>gpt-3.5-turbo-1106</b> : In summary, the evidence from these studies suggests that treatment of asymptomatic trichomoniasis in pregnant women with metronidazole <b>does not</b> prevent preterm delivery, and the impact of Trichomonas infection on perinatal outcomes may be limited.
Evidence omission	Reference summary	Insufficient evidence exists from randomised trials to confirm the relative effects of <b>open versus closed reduction of intracapsular fractures</b> , or the effects of <b>intra-operative impaction or compression</b> of an intracapsular fracture treated by internal fixation.
	Output summary	<b>Flan-UL2-20b</b> : The results of this study suggest that <b>closed reduction</b> is as effective as <b>open reduction</b> in the treatment of displaced intracapsular fractures of the neck of the femur in young adults.
Semantic distortion	Original text	We observed <b>unfavourable generation shifts in diabetes</b> among men but not among women.
	Simplified text	<b>gpt-3.5-turbo-instruct</b> : We noticed that men were <b>more likely to develop diabetes</b> compared to women.
Fabricated attribute	Original text	Framingham offspring study participants who attended the seventh examination were at least 60 years old and free of dementia and stroke were included.
	Simplified text	<b>gpt-4-0613</b> : The Framingham Offspring Study included participants who were at least 60 years old. These participants did not have dementia or <b>a history of stroke</b> .
Misinterpretation	Original text	<b>Subordinate individuals</b> had lower heart rates than dominants, and consequently they suffered lower rates of body mass loss during periods of restricted pellet food supply.
	Simplified text	<b>Gemini-1.0-pro</b> : When <b>animals were lower in the social hierarchy</b> , their heart rates were slower. This helped them lose less weight when food was scarce.
Hallucination	Original text	Five patients with DM2 and stage II/III HF consumed dark chocolate and a beverage containing approximately 100 mg of EPI per day for 3 months.
	Simplified text	<b>PMC-LLaMA-13b</b> : The patients with DM2 and stage II/III HF who consumed dark chocolate and a beverage containing approximately 100 mg of epicatechin per day for three months <b>experienced a significant improvement in their cardiac function</b> .

**TABLE II:** Task-specific prompt used for querying models. Blue text denotes domain-specific prefixes. Yellow, blue and green highlights represent ICL demonstrations, knowledge-guided prompting, and CoT prompting, respectively. Input variables are denoted by red text enclosed in `<>`.

Task	Prompts
PICO extraction	<p><b>Task Description</b>  You are a skilled medical expert. Your task is to generate an HTML version of an input text, marking up specific entities related to healthcare. The entities to be identified are: "Participant", "Intervention", "Control", and "Outcomes". Use HTML <code>&lt;span&gt;</code> tags to highlight these entities. Each <code>&lt;span&gt;</code> should have a class attribute indicating the type of the entity.</p> <p><b>Markup Format</b>  Use <code>&lt;span class="participant"&gt;</code> to denote a participant entity.  Use <code>&lt;span class="intervention"&gt;</code> to denote an intervention entity.  Use <code>&lt;span class="control"&gt;</code> to denote a control entity.  Use <code>&lt;span class="outcome"&gt;</code> to denote an outcome entity.  Leave the text as it is if no such entities are found.</p> <p><b>Entity Recognition Guide</b>  [entity recognition guide for PICO extraction task]</p> <p><b>Examples (N-shot)</b> [i=1,2,...N]  Input i: <code>&lt;example.input.text.i&gt;</code>  Annotated output i: <code>&lt;example.output.i&gt;</code></p> <p><b>Input Text</b>  <code>&lt;input.text&gt;</code></p>
Biomedical QA	<p><b>Task Description</b>  You are a skilled medical expert. Considering the information from a biomedical study provided in the reference text, is it correct to conclude that "<code>&lt;question&gt;</code>"? Please respond with "yes" or "no". Please first respond with "yes" or "no", followed by a brief explanation of your reasoning process, ensuring that your explanation aligns with the study's findings.</p> <p><b>Examples (N-shot)</b> [i=1,2,...N]  Question i: <code>&lt;example.question.i&gt;</code>  Reference text i: <code>&lt;example.reference.text.i&gt;</code>  Answer i: <code>&lt;example.answer.i&gt;</code></p> <p><b>Input Text</b>  Question: <code>&lt;question&gt;</code>  Reference text: <code>&lt;reference&gt;</code></p>
RCT summarization	<p><b>Task Description</b>  You are a skilled medical expert. Consolidate the information from these randomized controlled trial abstracts into a comprehensive summary.</p> <p><b>Summarization Guide</b>  [summarization guide for RCT summarization task]</p> <p><b>Examples (N-shot)</b> [i=1,2,...N][j=1,2,... # References in Example i]  Reference Text <code>&lt;j&gt;</code> in Example <code>&lt;i&gt;</code>  Title: <code>&lt;example.title.i.j&gt;</code>  Abstract: <code>&lt;example.abstract.i.j&gt;</code>  Summarization <code>&lt;example.summarization.i&gt;</code></p> <p><b>Input Text</b> [j=1,2,... # References in Input]  Reference Text <code>&lt;j&gt;</code>  Title: <code>&lt;title.j&gt;</code>  Abstract: <code>&lt;abstract.j&gt;</code></p>
Medical text simplification	<p><b>Task Description</b>  You are a skilled medical expert. Simplify the given medical text, making complex medical information accessible and relevant to people without a medical background.</p> <p><b>Summarization Guide</b>  [summarization guide for medical text simplification task]</p> <p><b>Examples (N-shot)</b> [i=1,2,...N]  Input i: <code>&lt;example.input.text.i&gt;</code>  Simplified output i: <code>&lt;example.output.i&gt;</code></p> <p><b>Input Text</b>  <code>&lt;input.text&gt;</code></p>