

Enhancing LLM Generation with Knowledge Hypergraph for Evidence-Based Medicine

Chengfeng Dou^{1,2}, Ying Zhang³, Zhi Jin^{1,2}, Wenpin Jiao^{1,2}, Haiyan Zhao^{1,2}, Yongqiang Zhao^{1,2}, Zhengwei Tao^{1,2}

¹ School of Computer Science, Peking University;

² Key Laboratory of High Confidence Software Technologies(PKU), MOE, China

³ Beijing Key Lab of Traffic Data Analysis and Mining, Beijing Jiaotong University, Beijing, China

{chengfengdou, zhi jin, jwp, zhhy, sei}@pku.edu.cn

{tttzw, yongqiangzhao}@stu.pku.edu.cn {19112043}@bjtu.edu.cn

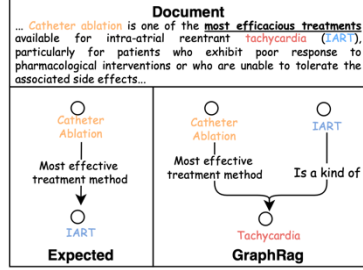


Figure 1: The phenomenon of mis-decomposition of complex relationships. LLMs omit the conditional variable “intra-atrial reentrant” for “tachycardia,” leading to incorrect extraction.

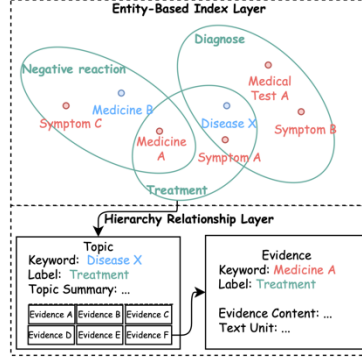


Figure 2: The Schema of EbmkG. The green ellipse denotes the hyperrelations corresponding to topic. The blue entities denote topic keywords, while red entities indicate evidence keywords. Evidence under the same topic has the same label.

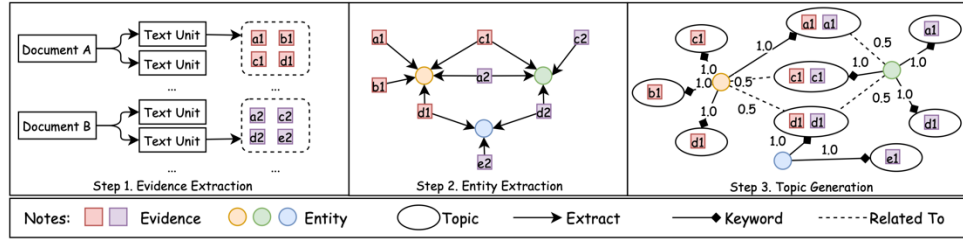


Figure 3: The construction process of EbmkG involves several key elements. Colored squares are utilized to represent evidence, with different colors distinguishing the sources of the evidence and the text within the squares indicating specific label of the evidence. Colored circles are employed to represent entities that are extracted from the evidence. White ovals represent topics, which are derived from evidence with the same aspect words.

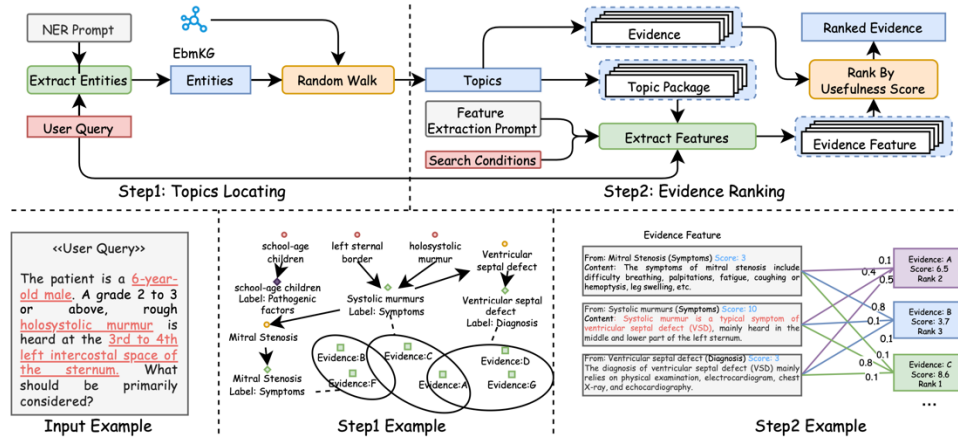


Figure 4: This figure is divided into two sections: the upper part outlines data processing, and the lower part provides an example. In ‘Step 1 Example’, entities in the EbmkG are shown as circles, topics as diamonds, and evidence as rectangles. Only green nodes, filtered by Personal PageRank, advance to Step 2. In ‘Step 2 Example’, an LLM assigns a Usefulness Score to each topic based on predefined Search Conditions, which determines the final evidence score. For NER and Feature Extraction Prompts, see Figure 8 (Search Words Extraction and Evidence Features Extraction). For Search Conditions, refer to Figure 9.

<p>Evidence Extraction</p> <p>Instruction: You are a useful medical AI assistant, please help me extract the information about {keyword}-{label} in the following document.</p> <p>Document: {text_unit}</p>	<p>Topic Summarization</p> <p>Instruction: Please write a brief summary for the chapter {keyword}-{label} based on the following evidence.</p> <p>Evidence: {evidence_list}</p>
<p>Entity Extraction</p> <p>Document: {document}</p> <p>Instruction: Please extract all the medical entities in the above snippet and return them according to the following template:</p> <pre> `json {{{ "name": "...", "type": "...", "mesh headings": [...]}}, ...] ` </pre> <p>Here is an explanation of the fields:</p> <ul style="list-style-type: none"> - name: The description of this entity. - type: the type of the entity, the possible types of the entity are: [Disease, Symptom, Life Habits, Pathogen, Medical Test, Drug, Surgery]. - mesh headings: the corresponding ``Medical Subject Headings`` of the entity, one entity can correspond to multiple headings, please use List to return. If no mesh headings can be found for this entity, please rewrite the entity name according to the mesh naming convention. 	<p>Search Words Extraction</p> <p>Query: {User Query}</p> <p>Task: Please extract the medical entities from the user input that will help in completing the task and try to extract more than 5 medical entities as more entities will help in the subsequent retrieval of evidence.</p> <p>Output Format: [{"name": "...", "type": "...}], ...]</p> <p>Evidence Features Extraction</p> <p>Query: {User Query}</p> <p>Task: To answer the user question, we looked at medical databases to find relevant information. Now after a rough ranking of documents, we have targeted the following documents, which you can consider as absolutely credible.</p> <p>Documents: {Topic Chunk}</p> <p>Please follow the rules below to identify the content in the document that helps answer the patient's question and score it based on relevance.</p> <p>Rules: {Rules}</p> <p>Output Format: [{"reference": "...", "score": ...}]</p>

Figure 8: Prompts used by Graph Construction and Evidence Retrieval.

<p>CMHD</p> <p>Please review the reference materials and indicate whether they contain any content that demonstrates errors in the doctor's response.</p> <p>These errors specifically refer to mistakes in medical knowledge, including incorrect diagnoses, medication errors, improper medical procedures, and inappropriate healthcare advice.</p> <p>Please provide a score for the level of disagreement, with the score ranging from 0 to 10. A higher score indicates a stronger level of disagreement.</p>	<p>MedQA and NLPCC</p> <p>I would like you to evaluate the relevance of a reference to a multiple-choice question on a medical exam. Relevance in this context means that the reference provides evidence from the literature that supports or refutes one of the answer options.</p> <p>Please assign a relevance score between 0 and 10, where a higher score indicates greater relevance.</p>
<p>CMB-Clin</p> <p>I would like you to identify the direct relevance between the reference materials and the doctor's questions. Please consider the following approach:</p> <ol style="list-style-type: none"> 1. If the doctor is asking about diagnostic information, focus on whether the reference materials include diagnostic methods for the specific disease. 2. If the doctor is inquiring about treatment options, concentrate on the treatment plans provided in the reference materials. 3. If the doctor is seeking an overview of a disease, pay attention to the summary and review sections in the reference materials. <p>The scoring range is from 0 to 10, with higher scores indicating that the documents are more likely to address the doctor's questions effectively. If the reference materials do not directly provide the answer the doctor is looking for, please assign a score of 3 or lower.</p>	<p>DDA</p> <p>I am providing you with the abstracts of reference documents. I would like you to evaluate whether these references help in identifying unreasonable decisions made by doctors and assign a score. Please consider the following aspects carefully:</p> <ol style="list-style-type: none"> 1. If the decision is related to diagnosis, assess whether the references can help verify: cases of missed or incorrect diagnoses, and whether there was a failure to collect complete patient information. 2. If the decision is related to lifestyle advice, assess whether the references can help verify: if the advice includes content that could worsen the patient's condition or even endanger their life, such as ignoring contraindications of the disease. 3. If the decision is related to medication advice, assess whether the references can help verify: if there are errors in the knowledge provided, if there is a disregard for drug interactions, or if there is a neglect of potential side effects. 4. If the decision is related to medical care advice, assess whether the references can help verify: if the level of care recommended is appropriate, if emergency measures are suitable, if the recommended department is correct, and if the suggested tests are appropriate. <p>The scoring range is from [0-10] points. The higher the score, the more effectively the documents can highlight the unreasonableness of the doctor's decisions.</p>

Figure 9: Search conditions for different datasets

```
# Task Input
- Input: {input}
- Golden Answer (Key Points): {key_points}
- AI Response: {output}

# Task Description
You are tasked with evaluating the reliability of an AI model's response based on the provided key points in the golden answer. Your evaluation should be formatted as a JSON array, where each object represents the assessment of a specific key point.

# Output Format
The JSON format should be as follows:
```json
[{"keypoint_number": ..., "contradict": ..., "description": ...}, ...]
```

Field Explanations:
- keypoint_number (String): The number or identifier of the key point. Ensure that all key points are accounted for.
- contradict (Boolean):
  - True: if the AI response misses the key point or contradicts the key point.
  - False: if the AI response correctly addresses the key point.
- description (String): A brief explanation of why the contradict field is set to True or False.
```

Figure 12: The prompt of key points evaluation

```
# Task Input
- User Query: {user_query}
- Document: {document}

# Task Description
Given the standard answer to the user's question: {golden_answer}, evaluate how useful the provided document is for answering the user's query.

# Output Format
Return the evaluation in JSON format as follows:
```json
{"score": <int>, "reason": <str>}
```

Field Explanations:
- score (Integer): A score between 0 and 10, indicating the usefulness of the document.
- reason (String): A brief explanation justifying the score.
```

Figure 15: The prompt of Usefulness Evaluation

You are a powerful search engine evaluation assistant. I am using two search engines to retrieve relevant documents for a specific question. The question is as follows:

{question}

The search engines have returned their respective results:

[search_engine1
{docs1}]

[search_engine2
{docs2}]

The known correct answer to the question is: {answer}.

Please help me evaluate which search engine's results are better by considering the following dimensions:

1. **Recall**: Which search engine's results more directly provide the correct answer or are strongly related to the correct answer.
2. **Precision**: Which search engine's results contain fewer irrelevant or misleading pieces of information.

Note: Do not simply assess the quality of the documents based on the amount of information they contain. Instead, consider whether the documents provide direct assistance in answering the question.

Note: When evaluating the accuracy, only check if the content of the documents contains any incorrect information; do not consider whether the documents are relevant to the question.

Please return your results in the following well-structured format:

```
```json
{"recall": {"answer": ..., "exp": ...}, "precision": {"answer": ..., "exp": ...}}
```
```

Here is an explanation of the fields:

- **answer (Str)**: Choose one from ['Search Engine 1 is better', 'Search Engine 2 is better', 'Both are equally good'].
- **exp (Str)**: Please explain why you chose this value.

Figure 13: Prompt for calculating Recall and Precision.

| Test Sample | | Why Graph Rag is Wrong | |
|---|---|---|--|
| <p>L4-5 central disc herniation involves ().</p> <p>A: L4 nerve root.</p> <p>B: L5 nerve root. (Graph Rag)</p> <p>C: S1 nerve root.</p> <p>D: Cauda equina nerve. (Native Rag, EbmKG Rag)</p> <p>E: Femoral nerve.</p> | | <p>An L4/L5 disc herniation does have the potential to compress a nerve root, but exactly which nerve root is compressed depends on the location and direction of the herniation. If the L4/L5 disc herniation is posterolateral, then it is more likely to compress the L5 nerve root. In the case of a centralized disc herniation, it may compress the spinal cord or cauda equina rather than a single nerve root.</p> | |
| <p>Native Rag</p> <p>From: Multiple radiculopathy: spinal stenosis, infectious, cancerous and inflammatory nerve root syndromes</p> <p>Content: ...Cauda equina syndrome - A variant of lumbosacral polyneuropathy is cauda equina syndrome, in which the cauda equina is compressed within the spinal canal. The most common cause is a large centralized disc herniation or centralized spinal stenosis...</p> | <p>Graph Rag</p> <p>Src: L4/L5 RADICULOPATHY</p> <p>Dst: L5 NERVE ROOTS</p> <p>Rel: L4/L5 RADICULOPATHY is a type of lumbar disc herniation that usually compresses the L5 nerve root.</p> <p>-----</p> <p>Src: L5 RADICULOPATHY</p> <p>Dst: Lumbar disc herniation</p> <p>Rel: Lumbar disc herniation is a common cause of L5 radiculopathy</p> | <p>EbmKG Rag</p> <p>From: Back pain in children and adolescents: Overview of causes</p> <p>Keyword: Back pain in children and adolescents</p> <p>Aspect: symptoms</p> <p>Content: Low back pain, sciatica (pain radiating below the knee), limited spinal mobility, limited passive straight leg raising. Severe central disc herniation can lead to cauda equina compression, bladder dysfunction, and anesthesia in the saddle area.</p> | |

Figure 16: Comparisons of results generated by various RAG aRandom Walkoaches.