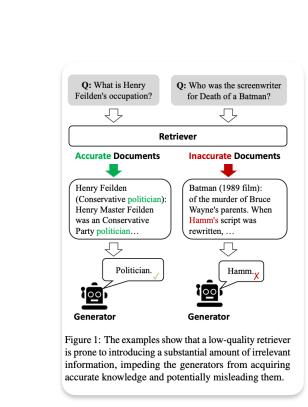


**2401**PDF Document · 668 KB

### • Large Language Models (LLMs) inherently hallucinate due to reliance on parametric knowledge. • Retrieval-Augmented Generation (RAG) is a practical complement but is vulnerable to irrelevant retrieved documents. • Corrective Retrieval Augmented Generation (CRAG) is proposed to enhance generation robustness. • CRAG includes a **lightweight retrieval evaluator** that assesses retrieved document quality for a query, triggering different knowledge retrieval actions based on a confidence degree. • To overcome limitations of static corpora, CRAG utilizes large-scale web searches to augment retrieval results. • A decompose-then-recompose algorithm is designed to focus on key information and filter irrelevant content from retrieved documents. • CRAG is plug-and-play, seamlessly integrating with various RAG-based approaches. • Experiments show CRAG significantly improves the performance of RAG-based methods across different generation tasks.



**LLM Hallucinations:** Large Language Models (LLMs) exhibit impressive abilities but inevitably suffer from **hallucinations** due to **factual errors** and their inability **to guarantee accuracy solely from parametric knowledge.** RAG Limitations:Retrieval-Augmented Generation (RAG) integrates external knowledge but its effectiveness relies heavily on the relevance and accuracy of retrieved documents. Low-quality retrievals can introduce irrelevant information, mislead generators, and cause hallucinations.

Problem with Current RAG: Most conventional RAG approaches indiscriminately incorporate retrieved documents, regardless of relevance. They also treat entire documents as reference knowledge, even though much of the text is often non-essential.

★Introducing CRAG: This paper proposes: Corrective Retrieval-Augmented Generation (CRAG)\*\* to self-correct retriever results and improve document utilization for augmentation, specifically addressing scenarios where the retriever returns inaccurate

### Key points of CRAG:

\* Lightweight Retrieval Evaluator: Assesses the overall quality of retrieved documents for a query, providing a confidence degree (Correct, Incorrect, Ambiguous). • Dynamic Knowledge Retrieval: For "Incorrect" or "Ambiguous" cases, CRAG integrates \*\*large-scale web searches\*\* as an extension to static corpora, broadening the scope and diversity of retrieved information. \* Decompose-then-Recompose Algorithm: Refines retrieved information by selectively focusing on key insights and eliminating redundant/non-essential contexts, optimizing data utilization.

\* Plug-and-Play & Effectiveness: CRAG is designed to be plug-and-play, demonstrated by its implementation into existing RAG frameworks (RAG, Self-RAG). Experiments show it significantly improves RAG-based approach performance.

# **CRAG** aims to **correct** and refine retrieved information to enhance the quality of generated output

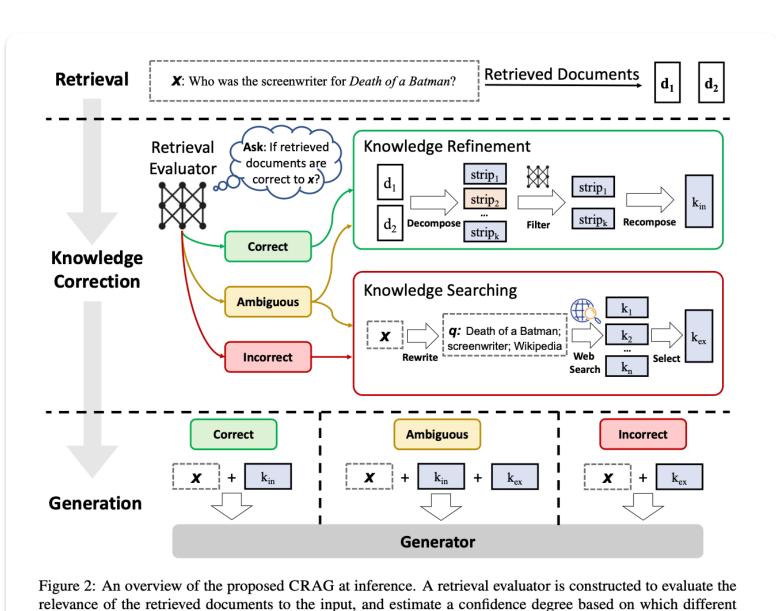
### Core process of CRAG Evaluator -> Action Trigger -> refinement

13 end

14 G predicts y given x and k

**Overview of CRAG** 

**CRAG** 



Algorithm 1: CRAG Inference **Require:** E (Retrieval Evaluator), W (Query Rewriter), G (Generator) **Input** : x (Input question),  $D = \{d_1, d_2, ..., d_k\}$  (Retrieved documents) **Output**: y (Generated response) 1  $score_i = E$  evaluates the relevance of each pair  $(x, d_i), d_i \in D$ 2 Confidence = Calculate and give a final judgment based on  $\{score_1, score_2, ... score_k\}$ // Confidence has 3 optional values: [CORRECT], [INCORRECT] or [AMBIGUOUS] 3 **if** Confidence == [CORRECT] **then** Internal\_Knowledge = Knowledge\_Refine(x, D)s  $k = Internal_Knowledge$ 6 **else if** Confidence == [INCORRECT] **then** External\_Knowledge = Web\_Search(W Rewrites x for searching)  $k = External_Knowledge$ 9 **else if** Confidence == [AMBIGUOUS] **then** Internal\_Knowledge = Knowledge\_Refine(x, D)External\_Knowledge = Web\_Search(W Rewrites x for searching)  $k = Internal_Knowledge + External_Knowledge$ 

knowledge retrieval actions of {Correct, Incorrect, Ambiguous} can be triggered.

### Overview

1 Initial Retrieval: Documents are fetched by another retriever, such as Contriever. Retrieval Evaluator (T5-large based): This component scores the retrieved documents based on their relevance to the Action Trigger: Based on these scores, the system determines whether a document is 'Correct', 'Incorrect', or 'Ambiguous', and then initiates an appropriate corrective action accordingly (e.g., use the document as is, discard and re-retrieve, or perform further verification).

## **Retrieval Evaluator**

CRAG designs corrective strategies to enhance generation robustness during inference. For an input query and retrieved documents, a lightweight retrieval evaluation estimates relevance score

4 Generation: An LLM generates the final answer based on the corrected/selected documents.

**Usefulness**" **Determination** Based on this relevance score, the system determines whether "this search result is useful or not. **Determination = threshold** If the relevance score S is less than the threshold T(S<T), then the system deems the search result "not useful.

To train retrieval evaluator: Base model: T5-large Data Used for fineTuning: Positive Samples: High-quality, relevant passages obtained from the PopQA dataset through the "golden subject wiki title" for each question. These served as the "relevance Negative Samples: Documents randomly sampled from the retrieval results that were, in fact, irrelevant to the input query.

### Scoring

Quantified into three confidence degrees - correct, incorrect and ambiguous three types of actions are designed and triggered accordingly where the upper and lower thresholds are set.

If the confidence score is higher than the upper threshold, the retrieved document is identified as -> implies that no corrective action is needed, and the document can be used as is. Incorrect if below the lower threshold. Otherwise, a more soft and intermediate action, judgment triggers a specific, often more drastic, corrective action (e.g., discarding the document, reretrieving information from scratch).

Ambiguous: Ambiguous is executed. Each retrieved document is conducted individually and integrated triggers a "soft and intermediate action" (e.g., further verification, re-ranking, or using a different generation strategy), as the document is neither clearly correct nor clearly incorrect.

## To evaluate the overall performance of CRAG

2. LLNIS	Evaluated:
•	Public LLMs:
•	LLaMA2-7B, 13B
•	Alpaca-7B, 13B (instruction-tuned)
•	CoVE65B (introduces iterative engineering)
•	Proprietary LLMs (trained with proprietary data):
•	LLaMA2-c13B (chat version)
•	ChatGPT
•	Perplexity.ai (InstructGPT-based production search system)
3. Retrie	val Methods/Baselines:
•	Baselines without retrieval: Standard LLMs as listed above.
•	Standard RAG:
•	LLaMA2-7B, 13B
•	Alpaca-7B, 13B
•	LLaMA2-7B instruction-tuned in Self-RAG
•	The standard RAG approach involves prepending top retrieved
documen	ts to the query before the LM generates output. Google Search API
is used fo	or retrieval.
•	Advanced RAG:
•	SAIL (Luo et al., 2023): Instruction-tuned an LM on Alpaca
data with	retrieved documents inserted before instructions.
•	Self-RAG (Asai et al., 2024): Tuned LLaMA2 on instruction-
tuning da	ta with reflection tokens labeled by GPT-4.
$\mathcal{C}$	Retrieval-augmented baselines trained with private data:

• **Accuracy:** For PopQA, Pub, and ARC datasets.

• **FactScore:** For the Bio dataset.

(Accuracy) (FactScore) (Accuracy) (Accuracy) LMs trained with propriety data 51.8 79.9 52.1 37.9 29.3 71.8 70.1 **75.3 54.7 75.3** 50.8 Perplexity.ai 71.2 Baselines without retrieval 49.8 45.0 29.4 29.4 53.4 Alpaca $_{13B}$ CoVE $_{65B}$ 55.5 54.9 24.4 50.2 - 71.2 - -Baselines with retrieval Self-CRAG 54.9 72.4 Self-CRAG 61.8 86.2

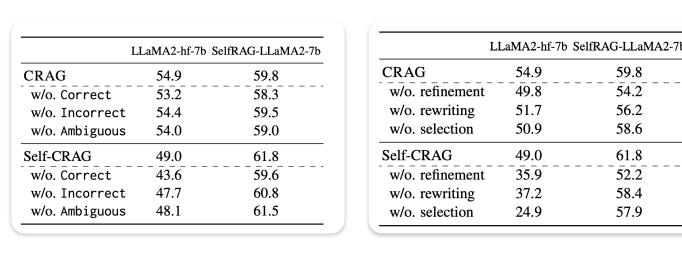
PopQA Bio Pub ARC

RAG (Retrieval-Augmented Generation): External information retrieval + LLM answer • CRAG (Contextualized Retrieval Augmented Generation): In RAG, it evaluates the quality of the retrieved information and adjusts the LLM's answer generation strategy accordingly (smarter utilization of information). • Self-RAG: In RAG, it trains the LLM to "reflect" and judge the necessity of retrieval and the quality of the answer itself (proactive role of the LLM). • Self-CRAG: Aims for the most optimized answer generation by combining CRAG's precise information management with Self-RAG's self-reflection capabilities.

## • CRAG consistently and significantly outperforms standard RAG across all four datasets

- **SelfRAG-LLaMA2-7b base:** Improves PopQA by 7.0%, Bio (FactScore) by 14.9%, PubHealth by 36.6%, and Arc-Challenge by 15.4%. • LLaMA2-hf-7b base: Improves PopQA by 4.4%, Bio (FactScore) by 2.8%, and Arc-Challenge by
- Self-CRAG vs. Self-RAG: Self-CRAG also demonstrates substantial improvements over the state-ofthe-art Self-RAG.
- LLaMA2-hf-7b base: Improves PopQA by 20.0%, Bio (FactScore) by 36.9%, and Arc-Challenge by • SelfRAG-LLaMA2-7b base: Improves PopQA by 6.9%, Bio (FactScore) by 5.0%, and PubHealth by

## **Results**

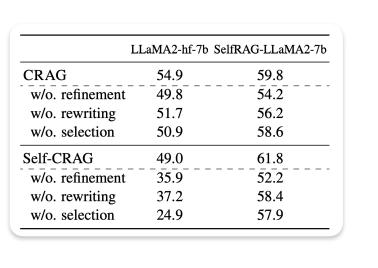


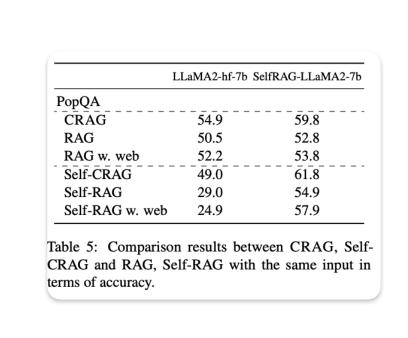
Ret-ChatGPT

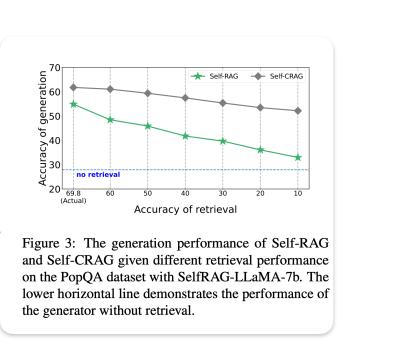
Perplexity.ai

4. Evaluation Metrics:

• Ret-LLaMA-chat







## Process

→ LLM Generation

User Input → Initial Retrieval → Retrieval Evaluator (and optional Web Search) → Decompose (of retrieved documents) → Filtering/Scoring (of decomposed strips)  $\rightarrow$  Recompose (of relevant strips)

## **Step 1:** User Input

Step 2 : Initial Retrieval :Based on the user's query, relevant documents are retrieved from an external database (e.g., a vector DB-FAISS, Pinecone, Vewaviate, Milvus, Chroma..).

**Step 3:** Retrieval Evaluator The quality of the **entire retrieved documents** (relevant/irrelevant) is If the quality of the retrieved documents is deemed low, CRAG can trigger a web search to find better information => corrective part

**Step 4:** Decompose • The **documents** obtained from the initial retrieval or web search are broken down into smaller "knowledge strips."

Step 5: Filtering and Scoring • Each "knowledge strip" is then **scored for relevance**, and strips with low relevance to the user's query are **filtered out**. • A Retrieval Evaluator might be involved here again, or a separate module might perform this role.

Step 6: Recompose • The filtered, i.e., highly relevant, "knowledge strips" are combined again to form the "internal knowledge" that the LLM will

Step 7: LLM Generation

## # Conceptual Pseudocode for Self-CRAG

# Define hypothetical functions and classes

def generate(self, prompt: str, reflection\_tokens=False) -> str: The LLM generates text.

If reflection\_tokens=True, it can include Self-RAG-like reflection tokens. pass # Actual LLM API call logic def assess\_document\_quality(self, query: str, document: str) -> str: Acts as CRAG's Retrieval Assessor. Evaluates document quality as 'High', 'Medium', 'Low', etc. pass # LLM-based document quality assessment logic

def self\_reflect(self, prompt: str) -> dict: Generates Self-RAG's self-reflection tokens. E.g., {'is\_relevant': True, 'is\_supportive': False, 'is\_factual': True} pass # LLM-based reflection token generation logic

### def retrieve\_documents(self, query: str, top\_k: int = 5) -> list[str]: Retrieves relevant documents based on the query.

pass # Vector DB query, web search API call, etc. def run\_self\_crag(user\_query: str, llm: LLM, retriever: Retriever) -> str: Pseudocode representing the overall workflow of Self-CRAG.

if not retrieved\_documents:

\nDocument: {doc}\nQuestion: {user\_query}"

### print(f"User Query: {user\_query}") # 1. Self-RAG's initial judgment: Is retrieval needed? # In practice, the LLM makes this decision after seeing the initial question, simplified here. # Process where LLM analyzes the question and generates Self-RAG Reflection Tokens (initial) initial\_reflection = Ilm.self\_reflect(f"Does the question '{user\_query}' require external knowledge retrieval?") if not initial\_reflection.get('requires\_retrieval', True): # Hypothetical token print("Initial reflection: No retrieval needed. Generating directly...") return Ilm.generate(f"Question: {user\_query}")

print("Initial reflection: Retrieval seems necessary. Proceeding with retrieval...") # 2. Document Retrieval retrieved\_documents = retriever.retrieve\_documents(user\_query)

doc\_quality\_assessment = llm.assess\_document\_quality(user\_query, doc)

print(f"CRAG Document Quality: {doc\_quality\_assessment}")

generated\_text = Ilm.generate(prompt, reflection\_tokens=True)

print("No documents retrieved. Generating based on LLM's internal knowledge.") return llm.generate(f"Question: {user\_query}") final\_answer = "" for doc in retrieved\_documents: print(f"\nProcessing retrieved document: {doc[:100]}...") # Print only a portion of the document #3. CRAG's Document Quality Assessment (Retrieval Assessor)

# 4. Self-RAG's Reflection on Retrieval Results retrieval\_reflection = Ilm.self\_reflect(f"Query: {user\_query}\nDocument: {doc}\nls this document relevant and print(f"Self-RAG Retrieval Reflection: {retrieval\_reflection}") # 5. Determine Generation Strategy by combining CRAG and Self-RAG judgments (Adaptive Generation) generation\_strategy = "internal\_knowledge" # Default strategy

if doc\_quality\_assessment == "High" and retrieval\_reflection.get('is\_relevant') and retrieval\_reflection.get('is\_supportive'): generation\_strategy = "fully\_utilize\_document" print("Strategy: High quality document, fully utilize.") prompt = f"Based on the following document, answer the question:\nDocument: {doc}\nQuestion: {user\_query}" generated\_text = Ilm.generate(prompt, reflection\_tokens=True) # Generate including reflection tokens elif doc\_quality\_assessment == "Medium" and retrieval\_reflection.get('is\_relevant'): generation\_strategy = "combine\_with\_internal" print("Strategy: Medium quality document, combine with internal knowledge.") prompt = f"Considering the following document, answer the question. You may also use your internal knowledge:

else: # Low quality, or not relevant/supportive based on Self-RAG reflection generation\_strategy = "re\_evaluate\_or\_skip" print("Strategy: Low quality/irrelevant document, re-evaluate or skip.") # In practice, this could involve further searching, trying other documents, or generating with LLM's internal prompt = f"You received a document that was not very helpful. Please try to answer the question based on your general knowledge if possible, or indicate if more information is needed:\nQuestion: {user\_query}" generated\_text = llm.generate(prompt, reflection\_tokens=True)

# 6. Self-RAG's Final Reflection on the Generated Answer final\_generation\_reflection = Ilm.self\_reflect(f"Question: {user\_query}\nGenerated Answer: {generated\_text}\nIs this answer factual and complete?") print(f"Self-RAG Final Generation Reflection: {final\_generation\_reflection}") if final\_generation\_reflection.get('is\_factual') and final\_generation\_reflection.get('is\_complete'): final\_answer = generated\_text

print("Final answer deemed factual and complete. Stopping.") break # Stop if a satisfactory answer is obtained print("Answer not yet satisfactory. Trying next document or refining...") # In practice, this would involve logic to search again, try different strategies, or refine the answer. if not final\_answer: print("Could not find a satisfactory answer using retrieved documents. Generating best effort answer.")

return final\_answer # --- Usage Example (Not actually executable) ---# Create hypothetical LLM and Retriever objects (in reality, these would be API calls to OpenAI, HuggingFace, etc.)

 $\# my_IIm = LLM()$ # my\_retriever = Retriever() # answer = run\_self\_crag("What is the capital of France?", my\_llm, my\_retriever) # print(f"\nFinal Self-CRAG Answer: {answer}")

final\_answer = Ilm.generate(f"Please provide the best possible answer to: {user\_query}")

## Conclusion

The CRAG aims to address the issue of LLM producing hallucinations or factually incorrect statements The "decompose-recompse" algorithm plays a crucial role

Decomposing: broken down into smaller, more granular pieces of information, referred to as "knowledge strips Filtering: A retrieval evaluator then assesses the relevance of each individual knowledge strip to the user's query. Irrelevant strips are filtered out, ensuring that only the most pertinent information is retained. Recomposing: The relevant knowledge strips are then concatenated to form a refined set of "internal knowledge". This newly assembled information is what is ultimately used to generate the response

## 1) RAG (Retrieval-Augmented Generatio)

Objective: This technology helps LLMs (Large Language Models) answer questions not by solely relying on their pre-trained knowledge, but by searching for external documents (information) and generating more accurate and factual responses based on that retrieved information.

**Objective:** This technology helps LLMs (Large Language Models) answer questions not by solely relying on their pre-trained knowledge, but by searching for external documents (information) and generating more accurate and factual responses based on that retrieved information. Type to enter text 1 Retrieval: Documents relevant to the user's query are searched from an external database (e.g., web search engines, Wikipedia, academic **Generation:** The retrieved document content and the user's query are fed together into the LLM to generate the answer. Advantages: Reduces hallucination, allows for the incorporation of up-to-date information. **Disadvantages:** The quality of the generated answer can vary significantly depending on the quality of the retrieved documents. Problems can arise if irrelevant documents are retrieved or if the retrieved documents themselves contain errors.

2. CRAG (Contextualized Retrieval Augmented Generation)

1 **Retrieval:** Documents are retrieved in the same way as with RAG.

relevant the document is to the query, whether it contains sufficient information, etc.

picking only the trustworthy ones, or even adopting a completely different strategy."

utilizing and managing the quality of retrieved documents.

answer generation strategy based on that assessment.

or even refusing to generate an answer at all.

knowledge, or may attempt further searches.

filter out "useless information" or process it differently.

**Process (Additions to RAG):** 

3. Self-RAG

improving quality.

4. Self-CRAG

strategy based on that evaluation.

and generation processes.

"This question requires external retrieval."

(special markers within the prompt), such as:

CRAG is a proposed method to improve upon the shortcomings of RAG. Specifically, it focuses on more intelligently

Core Idea: It evaluates how useful retrieved documents are (e.g., informativeness, relevance) and adjusts the LLM's

2 Retrieval Assessor: The "quality" of the retrieved documents is evaluated. For example, it assesses how

3 Adaptive Generation: The LLM changes its answer generation strategy based on the evaluation results:

• **Medium-quality documents:** Refers to the retrieved documents but combines them with the LLM's own

Analogy: If RAG is like "referring to all search results for now," CRAG is like "assessing the reliability of search results,

Self-RAG is an RAG model trained to allow the LLM itself to perform "reflection" or "evaluation" on its retrieval

Core Idea: As the LLM generates an answer, it simultaneously makes meta-judgments (reflection) such as "Does

3 Search Result Evaluation and Answer Generation (Self-Reflection & Generation): The LLM

generates an answer based on the retrieved documents, but at the same time, it is trained to generate specific "tokens"

**Advantages:** The LLM can actively control the retrieval process and autonomously evaluate and improve the quality

Through these "reflection tokens," the LLM provides feedback to itself during the generation process, thereby

Disadvantages: The LLM must be specifically trained (instruction-tuned) to understand and generate these

**Objective:** To create synergy by adding CRAG's ability to "more meticulously manage the quality of

Self-RAG's Retrieval and Reflection: Similar to Self-RAG, the LLM autonomously decides

2 CRAG's Quality Evaluation and Adaptive Generation: CRAG's "search result evaluation"

• For example, even if Self-RAG determines "this search result is relevant," Self-CRAG further

evaluates the "depth" or "accuracy" of that search result in a CRAG-like manner. Instead of just using it, it

reacts more intelligently, for instance, by thinking "this information is insufficient, so I need to search

Advantages: By combining Self-RAG's proactivity with CRAG's meticulous information management

and at the same time, more meticulously checking how trustworthy the materials I referred to were."

Analogy: If Self-RAG is like "asking myself if I did well," Self-CRAG is like "asking myself if I did well,

+----+

| capital of France?"

+----+

+----+ | 2. Call run\_self\_crag |

| Function |

| (Passes User Query) |

+----+

| 3. Execute run\_self\_crag Function Internal Logic (Complex Self-CRAG & CRAG Logic)

| | - Retriever searches for relevant documents from external databases | |

| | 3.3.1. CRAG's Document Quality Assessment (Retrieval Assessor)

| | | - LLM self-judges "Is this document relevant and helpful?" | |

| | 3.3.3. Adaptive Generation Strategy Decision based on CRAG/Self-RAG Judgments |

| | - (e.g., 'Low' quality or 'irrelevant' -> re-evaluate/skip/internal knowledge only) |

| | | - LLM generates answer according to the decided strategy |

| 3.3.5. Self-RAG's Final Reflection on Generated Answer

| | (e.g., 'is\_factual': True, 'is\_complete': True)

| (If satisfactory, end iteration)

| | | - LLM self-judges "Is the generated answer factual and complete?" |

| | <--- (If more documents or not satisfactory, proceed to next document or try another strategy) ---+ |

+----+

| (Optimal Answer Selected) |

+----+

+----+

| 5. Present Answer |

| to User

| "The capital of

| France is Paris." |

+----+

| 4. Return Final Answer |

- (e.g., 'Medium' quality + 'relevant' -> combine with internal knowledge) |

| | | (e.g., 'is\_relevant': True, 'is\_supportive': True)

LLM assesses document informativeness, relevance, reliability |

| | 3.1. LLM's Initial Self-Judgment (Start of Self-RAG)

| | (e.g., 'requires\_retrieval': True)

| (If deemed necessary)

| | (e.g., documents about "Paris")

| (List of retrieved documents)

| | 3.3. Iterative Processing for Each Retrieved Document

| | 3.3.4. Answer Generation

| | 3.3.2. Self-RAG's Reflection on Retrieval Results

| | 3.2. Document Retrieval

| | - "Does this question require external retrieval?"

| 1. User Input

| "What is the

informativeness, reliability, etc., of the retrieved documents, and then fine-tunes the answer generation

Core Idea: It uses Self-RAG's self-reflection mechanism, but like CRAG, it fine-tunes the answer

logic is added here. Beyond simply "relevant/not relevant," it more precisely evaluates the

1 Question Analysis and Retrieval Decision: The LLM analyzes the query and autonomously decides,

this answer require retrieval?", "Is this search result helpful?", or "Is this answer factual?".

**Retrieval:** Based on its decision, it retrieves documents.

• [Is\_Relevant]: Is the search result relevant to the question?

• [Is\_Supportive]: Does the search result support the answer?

• [Is\_Factually\_Correct]: Is the generated answer factually accurate?

"reflection tokens." This means it's difficult to directly apply it to existing, general LLMs.

As the name suggests, Self-CRAG is a combination of CRAG and Self-RAG.

retrieved information" to Self-RAG's ability to "self-evaluate and reflect."

generation strategy based on the quality of the retrieved documents.

the necessity of retrieval and generates reflection tokens after retrieval.

more" or "this information is excellent, so it's enough on its own."

capabilities, it aims for the most powerful and accurate answer generation.

Advantages: Enhances the accuracy of LLM answers by increasing the reliability of the retrieved information. It can

• Low-quality documents: Does not rely on the retrieved documents, instead using the LLM's own knowledge,

• **High-quality documents:** Actively utilizes the retrieved documents to generate the answer.

### Decompose-then-Recompose evaluates the overall quality

of retrieved documents and

to find better information if

can even trigger web searches

the initial retrieval results are

**Decompose:** • Break down a complex and challenging large problem into smaller, more

Decompose-then-Recompose'

manageable, and independent parts (sub-problems or components). • Each sub-problem is much simpler and clearer than the original problem, making it easier to solve or understand individually. • During this process, it's crucial to identify the core and peripheral aspects of the problem, and to understand their interdependencies.

**ReCompose:** 

• Individually solve or optimize each of the decomposed sub-problems or components. • Subsequently, combine these individually solved parts to construct a complete solution for the original large problem. • At this stage, it's important to consider how each part interacts and contributes to the overall system, integrating them harmoniously.