

# Iterative Utility-Focused Evidence Refinement for Retrieval-Augmented Generation

Anonymous ACL submission

## Abstract

Relevance and utility are two frequently used measures to evaluate the effectiveness of an information retrieval (IR) system. Relevance emphasizes the aboutness of a result to a query, while utility refers to the result’s usefulness or value to an information seeker. In retrieval-augmented generation (RAG), high-utility results should be prioritized to feed to LLMs due to their limited input bandwidth. Re-examining RAG’s three core components—relevance ranking derived from retrieval models, utility judgments, and answer generation—aligns with Schutz’s philosophical system of relevances, which encompasses three types of relevance representing different levels of human cognition that enhance each other. These three RAG components also reflect three cognitive levels for LLMs in question-answering. Therefore, we propose an Iterative uTility-focused Evidence refineMent (ITEM) to promote each step in RAG. We conducted extensive experiments on retrieval (TREC DL, WebAP), utility judgment task (GTI-NQ), and factoid question-answering (NQ) datasets. Experimental results demonstrate improvements of ITEM in utility judgments, ranking, and answer generation upon representative baselines<sup>1</sup>.

## 1 Introduction

Relevance and utility are two frequently used measures to evaluate Information Retrieval (IR) performance (Saracevic, 1996, 1975; Saracevic et al., 1988). Relevance usually refers to *topical relevance* that measures the degree of fit between the subjects of a query and retrieved items, and the criteria of “aboutness” is used (Saracevic et al., 1988; Schamber and Eisenberg, 1988). In contrast, *utility* refers to the “usefulness” of retrieval items to an information seeker, of which the criterion is their “value” to the user (Saracevic, 1996; Saracevic et al.,

<sup>1</sup>Our code and benchmark can be found at <https://anonymous.4open.science/r/ITEM-B486/>.

Question: How does granulation tissue start?
[1] Healthy granulation tissue is granular and uneven in texture; it does not bleed easily and is pink / red in colour. The colour and condition of the granulation tissue is often an indicator of how the wound is healing. Dark granulation tissue can be indicative of poor perfusion, ischaemia and / or infection. (Related)
[2] Granulation tissue is collagen-rich tissue which forms at the site of an injury. As the body heals, this tissue fills in the injury, and may eventually scar over. The scar may fade over time, especially if the wound is small. In some cases, the body produces too much granulation tissue, in a condition known as proud flesh, in which case medical treatment may be required to halt the overproduction. The appearance of granulation tissue is a good sign. When a wound starts granulating, it means that the body is starting to rebuild after the injury. This highly fibrous tissue is usually pink because the body produces numerous small blood vessels to provide a supply of oxygen and nutrients to remove waste. (Highly relevant)
[3] Granulation Tissue. in Pathology. After tissue damage, repair process starts. It can begin as early as 24 hours. Fibroblasts and endothelial cells begin proliferating to form a specialized type of tissue that is the hallmark of healing called granulation tissue. The term derives from its pink, soft, granular appearance on surface of wound but its histological features are: it can begin as early as 24 hours. Fibroblasts and endothelial cells begin proliferating to form a specialized type of tissue that is the hallmark of healing called granulation tissue. The term derives from its pink, soft, granular appearance on surface of wound but its histological... (Perfectly relevant -> Utility)

Figure 1: An example between utility and relevance from TREC DL dataset.

1988). As the example from the TREC DL dataset shown in Figure 1, topical relevance does not necessarily mean utility, while utility indicates a higher standard of relevance. Since topical relevance is relatively easy to observe and measure (Schamber et al., 1990), the studies of IR models have been primarily focused on improving relevance for a long time (Bruce, 1994).

In the modern LLM era, retrieval-augmented generation (RAG) has become a hot research topic that equips LLMs with external knowledge (Xie et al., 2023; Shi et al., 2023; Izacard et al., 2023; Su et al., 2024; Glass et al., 2022). Given the constrained bandwidth of LLM inputs, it is essential to prioritize high-value results to guide LLMs. Consequently, utility needs to be emphasized more than topical relevance in RAG. More recently, Zhang et al. (2024) highlighted the use of LLMs for utility judgments. In this paper, we aim to further promote the utility judgment performance of LLMs so that RAG can be enhanced by high-utility references. **Schutz’s Philosophical Theory of Relevance.** Relevance is foundational in information retrieval (IR) and remains widely debated (Mizzaro, 1998). Saracevic (1996) discussed the nature of relevance in the IR system as the effectiveness of interactive exchange on different levels, and they are non-independent interdependencies, which are primarily influenced by Schutz’s philosophical theory of relevance. Schutz considered relevance as the property that determines the connections and relations

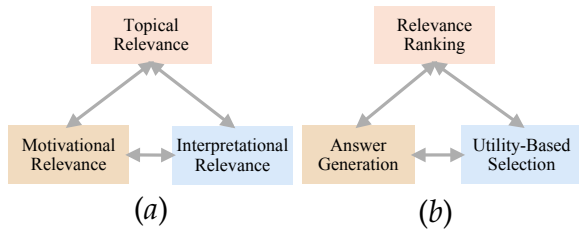


Figure 2: (a): Schutz’s “system of relevancies”, (b): the relation of each relevance to the components in RAG. The same color in the two frameworks is the corresponding connection.

in our lifeworld. He identified three types of basic and interdependent relevance that interact dynamically within a “system of relevancies” (Saracevic, 1996; Schutz, 1970): (i) Topical relevance, which refers to the perception of what is separated from one’s experience to form one’s present object of concentration; (ii) Interpretational relevance, which involves the past experiences in understanding the currently concerned object; and (iii) Motivational relevance, which refers to the course of action to be adapted based on the interpretations. The motivational relevance, in turn, helps obtain additional materials to become a user’s new experience, which further facilitates topical and interpretational relevance. Schutz posited that one’s perception of the world may be enhanced by this dynamic interaction, as shown in Figure 2. By incorporating utility judgments into RAG, we can re-examine its three components: topical relevance or relevance ranking derived from retrieval models, utility judgments, and answer generation. Topical relevance is an emerging focus on a topic, utility is the deeper understanding of the topic, and answers indicate the final solution based on the interpretations and will guide users’ actions. *Therefore, topic relevance, utility, and derived answer also reflect three cognitive levels for LLMs in question-answering, from low to high, i.e., aboutness, the value for deriving an answer, and the derived answer.*

**Iterative uTility-focused Evidence refineMent (ITEM).** Inspired by the philosophical theory of relevance, we believe the dynamic interactions between the three components in RAG can promote the performance of each step. To verify the idea, we leverage LLMs to perform each step in RAG shown in Figure 2, i.e., relevance ranking, utility judgments (classification), and answer generation. We propose an Iterative uTility-focused Evidence refineMent (ITEM) to enhance the utility judgment and QA performance of LLMs by interactions between the steps. ITEM has two variants depending

on whether relevance ranking is involved in the iterations. We are curious to see which option will be better for the tasks: fewer iterations with more components in an iteration, more iterations with fewer components in an iteration, or more iterations with more components.

We experiment on various information-seeking tasks, i.e., multi-grade passage retrieval on TREC DL (Craswell et al., 2020), multi-grade non-factoid answer passage retrieval on WebAP (Yang et al., 2016), utility judgments benchmark on GTI-NQ (Zhang et al., 2024), and factoid QA on NQ (Kwiatkowski et al., 2019). Experimental results have demonstrated that ITEM can outperform competitive baselines, including various single-shot judgment approaches in terms of utility judgments, topical relevance ranking, and answer generation, which confirms the viability of the adaptation of Schutz’s viewpoint of the relevance system into RAG. We also find that: 1) for difficult tasks (i.e., utility judgments of non-factoid answer passages in WebAP) and complicated candidate passage list (i.e., GTI-NQ), more components in the iteration and multiple iterations are usually more beneficial; 2) our ITEM achieves performance comparable to the long reasoning mode while requiring very lower computational cost, thereby offering a more efficient and practical solution for evidence refinement; 3) for factoid QA tasks, more iterations with fewer components performs the best, indicating that more components and more iterations are not always needed, especially for simpler tasks.

## 2 Related Work

**Multi-dimensional relevance.** The concept of “relevance” is central to information retrieval theory. Researchers have extensively debated its definition and measurement (Mizzaro, 1997). Early approaches primarily defined and assessed relevance through exact term matching (Vickery, 1959) or logical entailment (Hillman, 1964). However, subsequent empirical studies revealed the limitations of system-oriented relevance analysis, prompting diverse perspectives on relevance (Saracevic, 1975; Swanson, 1986; Saracevic, 1996; Lancaster, 1968; Goffman and Newill, 1964; Kemp, 1974; Bruce, 1994). For example, Cooper (1971) introduced logical relevance and utility. Saracevic (1996) summarized five frameworks for information science: systems, communication, situational, psychological, and interaction frameworks, and categorized

five distinct types of relevance, i.e., 1) system or algorithmic relevance; 2) topical or subject relevance; 3) cognitive relevance or pertinence; 4) situational relevance or utility; and 5) motivational or affective relevance. Bruce (1994) explored cognitive dimensions of relevance. Over time, scholarly consensus has coalesced around two primary perspectives: the system view and user view, with topical relevance and utility serving as their respective representative frameworks.

**Utility-Focused Information Retrieval.** Utility is a distinct measure of relevance compared to topical relevance (Zhao et al., 2024; Saracevic et al., 1988; Saracevic, 1975, 1996; Ji et al., 2024; Zhang et al., 2023), and more recently, Zhang et al. (2024) highlighted the use of LLMs for utility judgments. However, Zhang et al. (2024) only conducted a preliminary exploration of LLMs in utility judgments. Our work aims to further explore how to improve the performance of utility judgments for LLMs.

**Retrieval-Augmented Generation (RAG).** RAG approaches are widely employed to mitigate the hallucination issues in large language models (LLMs) (Xie et al., 2023; Zhou et al., 2024; Su et al., 2024). The current RAG approaches are categorized as follows: (i) single-round retrieval (Borgeaud et al., 2022; Lewis et al., 2020; Glass et al., 2022; Izacard et al., 2023; Shi et al., 2023), which involves using the initial input as a query to retrieve information from an external corpus and then the information is incorporated as part of the input for the model; and (ii) multi-round retrieval (Su et al., 2024; Jiang et al., 2023b; Ram et al., 2023; Khandelwal et al., 2020; Trivedi et al., 2023), which needs multi-round retrieval based on feedback from LLMs.

**Iterative Relevance Feedback via LLMs.** Recent works (Li et al., 2023; Shao et al., 2023) have achieved great success in using LLMs to obtain the information needs of the question as pseudo-relevance feedback for iterative retrieval. They posit that a single retrieval may not yield comprehensive information, thus requiring multiple retrievals. In contrast, our methodology involves making iterative utility judgments on the results obtained from a single retrieval.

### 3 Utility Judgments (UJ) via LLMs

Drawing inspiration from Schutz’s framework, we re-examine the three core components of RAG—namely, relevance ranking from retrieval

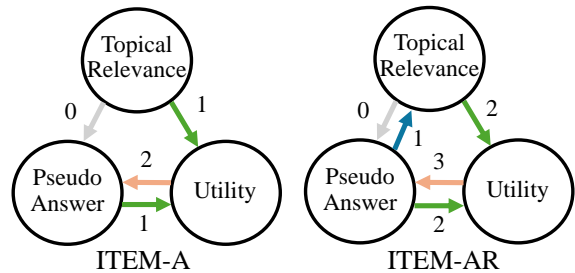


Figure 3: Flowchart illustrating the first iteration of ITEM. For ITEM-A, the process involves pseudo-answer generation followed by utility judgments and pseudo-answer generation. For ITEM-AR, the process includes pseudo-answer generation, relevance ranking, utility judgments, and pseudo-answer generation.

models, utility judgments from the selector, and answer generation—and observe that they closely align with Schutz’s system. Specifically, topical relevance, utility, and the generated answer correspond to three cognitive levels in LLM-based question answering: aboutness, the value for answer derivation, and the final answer itself, respectively, moving from lower to higher cognitive processing. Inspired by Schutz’s theory, we presume they can also interact with each other and enhance each other. Therefore, we propose an Iterative uTility-focused Evidence refineMent (ITEM) for utility judgments.

#### 3.1 Utility Judgments Definition

Following Zhang et al. (2024), given a question  $q$  and a list of retrieved passages  $\mathcal{D} = [p_1, p_2, \dots, p_n]$ , the goal of utility judgments for LLMs is to identify a set of passages  $U = \{p_1, \dots, p_m\}$ ,  $m$  is the number of passage with utility selected by LLMs. There are two typical input approaches for LLMs to select the set  $U$  from  $\mathcal{D}$ : pointwise and listwise. The pointwise approach independently evaluates the utility of individual passages as a binary classification task, while the listwise method assesses multiple passages simultaneously using the entire list as input.

#### 3.2 Single-Shot Utility Judgments

The most common approach to judge utility for the LLM is to perform a single-shot utility judgment, i.e.,  $U = LLM(q, \mathcal{D}, I)$ , where  $I$  is the instruction. To improve the accuracy of utility judgments, Zhang et al. (2024) proposed a unified framework where the LLM generates a pseudo-answer  $a$  and performs the utility judgments in a single output:  $a, U = LLM(q, \mathcal{D}, I)$ , i.e., **UJ-ExpA** in our baselines.

Answer generation instruction
<b>Implicit answer:</b> To answer the question, output what information is necessary to answer the question based on the references.
<b>Explicit answer:</b> Answer the following question based on the given information with one or few words/sentences.

Figure 4:  $I_a$  instruction contains the *implicit answer* and *explicit answer*.

Utility judgments instruction
<b>Listwise:</b> Directly output the passages you selected that have utility in generating the reference answer to the question.
<b>Pointwise:</b> Directly output whether the passage has utility in generating the reference answer to the question or not.
The <b>requirements</b> for judging whether a passage has utility in answering the question are: The passage has utility in answering the question, meaning that the passage not only be relevant to the question, but also be useful in generating a correct, reasonable and perfect answer to the question. Directly output the passages you selected that have utility in generating the reference answer to the question.

Figure 5:  $I_u$  instruction contains listwise and pointwise approaches.

### 3.3 Iterative uTility-focused Evidence refineMent (ITEM)

Drawing on Schutz’s theory of relevance, we introduce an Iterative uTility-focused Evidence refineMent (ITEM) framework for RAG. The framework iteratively refines evidence through cycles of topical relevance ranking, pseudo-answer generation, and utility judgments. To manage LLM inference cost, two variants are proposed: ITEM -A and ITEM -AR, which implement iterative loops between two or three RAG components, respectively, as shown in Figure 3.

#### ITEM with Answering in the Loop (ITEM-A).

Formally, at each iteration  $t$  ( $t \geq 1$ ), given the pseudo answer  $a_t$  generated based on the utility judgment result  $U_{t-1}$  from the previous iteration, we perform utility judgments on the candidate passages list  $\mathcal{D}$  to obtain a set of passages with utility  $U_t$ :

$$a_t = LLM(q, U_{t-1}, I_a), \quad (1)$$

$$U_t = LLM(q, \mathcal{D}, a_t, I_u), \quad (2)$$

where  $I_a$  represents the answer prompts for LLMs (as detailed in Figure 4),  $a_t$  can be in two forms: (i) *explicit answer* to the question  $q$ ; (ii) *implicit answer* that specifies the necessary information to answer the question  $q$ .  $I_u$  denotes the utility judgment prompts for LLMs (as detailed in Figure 5). We consider  $U_0 = \mathcal{D}$  as the initial candidate set, where  $\mathcal{D}$  represents the initial results ranked by a retriever such as BM25 (Robertson et al., 2009).

**ITEM with both Answering and Ranking of Topical Relevance in the Loop (ITEM-AR).** In the ITEM-A framework, topical relevance is not up-

Relevance ranking instruction
Reference answer: {answer}
Rank the {num} passages above based on their relevance to the query. The passages should be listed in descending order using identifiers.

Figure 6:  $I_r$  instruction

dated during the iteration process. To incorporate dynamic updating of topical relevance, we integrate a relevance ranking task into the ITEM framework, ensuring that all three tasks are executed in a loop. Formally, at iteration  $t$  ( $t \geq 1$ ), the answer  $a_t$  is generated based on the judging result  $U_{t-1}$  from the previous iteration. Subsequently, given  $a_t$ , the passage list  $R_{t-1}$  from the previous iteration is ranked based on the relevance to the question, yielding a new ranked list  $R_t$ . Finally, the judging result  $U_t$  is derived using the ranked list  $R_t$  and the answer  $a_t$ :

$$a_t = LLM(q, U_{t-1}, I_a), \quad (3)$$

$$R_t = LLM(q, R_{t-1}, a_t, I_r), \quad (4)$$

$$U_t = LLM(q, R_t, a_t, I_u), \quad (5)$$

where  $I_r$  is the relevance ranking prompt for LLMs (as detailed in Figure 6), respectively.

**Overall.** At iteration  $t$ , we have two ways to produce the set  $U_t$ : (i) Set-based approach: Asking LLMs to identify the set of passages that have utility using listwise and pointwise input forms, which called ITEM-A<sub>s</sub> or ITEM-AR<sub>s</sub> variants; (ii) Rank-based approach: Requesting LLMs to provide a ranked passage list based on utility (utility ranking prompt is shown in Appendix H.2) using the listwise input approach and considering the top- $k$  passages in the list to build  $U_t$ , which called ITEM-A<sub>r</sub> or ITEM-AR<sub>r</sub> variants. We set  $k = 5$  and more details of  $k$  are shown in Appendix A.3. We find that ITEM-AR<sub>r</sub> does not improve ranking performance as well as ITEM-A<sub>r</sub> (see Appendix A.4 for experimental analysis), so we do not employ ITEM-AR<sub>r</sub> in the ranking experiment. The rank-based approach has poor performance on the utility judgment task (details can be found in Appendix A.4), so we only employ the set-based approach on the utility judgment task. We stop the iteration when at most  $m$  ( $m=3$  in our paper) iterations are reached or the set of selected passages does not change, i.e.,  $t = m$  or  $U_t = U_{t-1}$ . Full details of all prompts can be found in Appendix H.

## 4 Experimental Setup

### 4.1 Datasets

Our experiments are conducted on four benchmark datasets, including two retrieval datasets,

i.e., TREC DL (Craswell et al., 2020) and WebAP (Yang et al., 2016), a utility judgment dataset, i.e., GTI-NQ (Zhang et al., 2024), and an open-domain question answer (ODQA) dataset, i.e., NQ (Kwiatkowski et al., 2019). Detailed statistics of the experimental datasets are shown in Table Appendix D. We use two representative retrievers to gather candidate passages in  $\mathcal{D}$  for utility judgments on TREC DL, WebAP, and NQ datasets. Construction details can be found in Appendix G.

**TREC DL.** We use the TREC-DL19 and TREC-DL20 datasets (Craswell et al., 2020). Judgments of TREC DL are on a four-point scale, i.e., “perfectly relevant”, “highly relevant”, “related”, and “irrelevant”. We consider the passages that are “perfectly relevant” to have utility. We filter questions of two datasets that contain the passages labeled “perfectly relevant: The passage is dedicated to the query and contains the exact answer. ” and combine them to form a whole dataset, i.e., the TREC DL. The annotation instruction for different points is shown in Appendix D from the original paper.

**WebAP.** WebAP (Yang et al., 2016) is a non-factoid answer passage collection built on Gov2. Non-factoid questions usually require longer answers, such as sentence-level or passage-level (Keikha et al., 2014a; Yang et al., 2016; Keikha et al., 2014b). Relevant passages are annotated and categorized as “perfect”, “excel”, “good”, and “fair”. The annotation instruction is similar to TREC DL. So we considered the “perfect” passages to have utility.

**NQ.** Natural Questions (NQ) consist of factoid questions issued to the Google search engine (Kwiatkowski et al., 2019). Each question is annotated with a long answer (typically a paragraph) and a short answer (one or more entities). Following Zhang et al. (2024), we use the questions that have long answers in our experiments.

**GTI-NQ.** Ground-truth inclusion (GTI) benchmark is constructed by Zhang et al. (2024) for utility judgment task. The GTI-NQ constructs a candidate passage set of 10 passages for each query sourced from the NQ dataset, comprising the long answer (designated as the utility passage), highly relevant noisy passages, weakly relevant noisy passages, and counterfactual passages.

## 4.2 Evaluation metrics

For the utility judgments task, we evaluate the results of judgments using **Precision**, **Recall**, and **F1**. For the ranking task, we use the normalized

discounted cumulative gain (**NDCG**) (Järvelin and Kekäläinen, 2017) to evaluate the ranking performance. For the answer generation task, we use the standard exact match (EM) metric and F1.

## 4.3 LLMs

We conduct our experiments using several representative LLMs, i.e., (i) ChatGPT (OpenAI, 2022) (we use the gpt-3.5-turbo-1106 version), (ii) Mistral (Jiang et al., 2023a) (the Mistral-7B-Instruct-v0.2 version), and (iii) Llama 3 (Meta, 2024) (the Meta-Llama-3-8B-Instruct version). (iv) Qwen 3 (Yang et al., 2025) (Qwen3-8B version), supporting seamless switching between thinking mode and non-thinking mode within a single model. To ensure the reproducibility of the experiments, the temperature for all experiments is set to 0.

## 4.4 Baselines

We utilize the following baselines on the utility judgments task and question answering performance based on the utility judgment results:

**Single-shot utility judgments.** (i) **Vanilla:** Ask LLMs to provide utility judgments based on the instruction directly. (ii) **UJ-ExpA:** Utility judgments and provide explicit answers simultaneously through a single output, which is shown to be effective in Zhang et al. (2024). (iii) **UJ-ImpA:** Utility judgments and provide implicit answers that are necessary to answer the question through a single output.

**$k$ -sampling.** (Zhang et al., 2024) proposed  $k$ -sampling to alleviate the sensitivity of LLMs to input order. Specifically, the  $k$ -sampling method randomizes the order of the input passage list  $k$  times in addition to the original input and aggregates the  $k + 1$  utility judgment results through voting. For fair comparison, we use the  $k = 5$ , more details are in Appendix F.

To evaluate the effectiveness of the proposed ITEM framework in ranking tasks, we are using a verbalized ranking. Therefore, we also employ another verbalized ranking method, i.e., **RankGPT** (Sun et al., 2023) as our main baseline, which uses the LLMs to directly rank input passages based on their relevance to the query.

## 5 Experimental Results of LLMs

This section will present the performance of each task within our ITEM framework. By default, the pseudo answer is the *explicit answer* in all experiments, if not specified otherwise.

Method	WebAP						TREC DL					
	Listwise			Pointwise			Listwise			Pointwise		
	M	L	C	M	L	C	M	L	C	M	L	C
Vanilla	20.79	21.79	28.43	23.05	25.09	26.85	45.67	49.39	55.19	45.11	47.64	49.84
UJ-ExpA	27.94	26.99	30.50	25.27	29.25	27.44	54.10	52.83	57.49	43.53	53.73	48.09
UJ-ImpA	25.06	26.22	29.89	28.35	25.29	26.32	48.29	48.22	56.18	48.31	50.20	48.83
5-sampling	30.16	28.97	31.49	-	-	-	52.31	52.68	60.49	-	-	-
ITEM-A <sub>s</sub> w. ExpA (1)	29.76	27.50	36.89	29.10	31.08	<u>32.02</u>	53.78	<u>53.66</u>	<u>62.52</u>	49.44	<u>52.09</u>	53.61
ITEM-A <sub>s</sub> w. ImpA (1)	26.06	25.59	34.97	28.28	30.53	29.34	49.39	53.73	58.11	46.01	53.68	<u>54.61</u>
ITEM-AR <sub>s</sub> w. ExpA(1)	<u>35.50</u>	<b>31.44</b>	36.58	-	-	-	52.34	48.97	62.00	-	-	-
ITEM-A <sub>s</sub> w. ExpA (3)	31.65	29.32	39.57	<b>30.50</b>	<b>32.67</b>	31.43	54.86	<b>56.03</b>	<b>63.18</b>	<b>51.74</b>	52.46	<b>55.74</b>
ITEM-A <sub>s</sub> w. ImpA (3)	28.36	26.10	<b>40.78</b>	30.13	29.64	<u>32.54</u>	52.05	55.14	60.56	46.59	<u>53.76</u>	54.90
ITEM-AR <sub>s</sub> w. ExpA(3)	<b>37.06</b>	29.08	38.58	-	-	-	<u>56.27</u>	52.10	61.37	-	-	-

Table 1: The F1 performance (%) of utility judgments with different LLMs on the different datasets (the numbers in parentheses represent  $m$ -values). “-” indicates no experiments are performed under the pointwise approach because of that the  $k$ -sampling method and our ITEM-AR<sub>s</sub> require listwise input. **bold** indicates the best performance. Underline means the best performance among all variants of our ITEM with the same  $m$  value. “M”, “L”, and “C” mean “Mistral”, “Llama 3” and “ChatGPT”, respectively.

Method	Llama3-8B		ChatGPT	
	Listwise	Pointwise	Listwise	Pointwise
Vanilla	43.38	28.55	59.37	35.31
UJ-ExpA	47.07	39.32	66.13	37.17
UJ-ImpA	43.31	38.72	57.40	37.29
k-sampling	49.20	-	71.17	-
ITEM-As-ExpA (1)	49.26	47.52	72.44	54.89
ITEM-As-Imp (1)	47.47	37.98	68.92	43.17
ITEM-ARs-ExpA (1)	<u>50.77</u>	-	<u>74.43</u>	-
ITEM-As-ExpA (3)	49.73	<b>48.90</b>	73.55	<b>55.45</b>
ITEM-As-Imp (3)	48.03	38.34	69.68	43.58
ITEM-ARs-ExpA (3)	<b>51.22</b>	-	<b>76.34</b>	-

Table 2: The F1 performance (%) of utility judgments with different LLMs on the GTI-NQ dataset. **bold** and Underline are defined in Table 1.

## 5.1 Utility Judgment Results

Table 1 shows the F1 performance on the TREC DL and WebAP datasets using three LLMs. Further, we utilize a better-performing open-source LLM, i.e., Llama-3.1-8B, and a closed LLM, i.e., ChatGPT, to conduct experiments on GTI-NQ, as shown in Table 2. We assess the long reasoning process of the popular LLM Qwen3 (Thinking) against baseline models and our own model without thinking, using the larger GTI-NQ dataset. Results are shown in Table 3. Since ITEM-A<sub>r</sub> and ITEM-AR<sub>r</sub> have poor F1 performance in utility judgments (refer to Table 13 in Appendix A.4 for details), we restrict our experiments to ITEM-A<sub>s</sub> and ITEM-AR<sub>s</sub> in this section.

**ITEM with a Single Iteration vs. Baselines.** All LLMs using our ITEM with a single iteration generally outperform the single-shot utility judgments on three datasets and may even surpass the  $k$ -sampling method. For example, ChatGPT on the TREC DL dataset using our ITEM-A<sub>s</sub> w. ExpA and ImpA in the listwise approach improve the F1 performance by 8.7% and 3.4% over UJ-ExpA and UJ-

Method	Qwen3		Time (s) / Query	
	List	Point	List	Point
Vanilla (w/o Think)	41.82	42.30	0.4	0.7
UJ-ExpA (w/o Think)	47.30	43.65	1.5	3.2
UJ-ImpA (w/o Think)	43.96	38.92	0.9	1.6
k-sampling (w/o Think)	49.28	-	9.2	-
Vanilla (w. Think)	55.96	-	31.2	-
( $m=1$ , w/o Think)				
ITEM-As-ExpA	50.53	52.56	2.0	4.3
ITEM-As-Imp	51.58	52.43	1.4	1.9
ITEM-ARs-ExpA	<u>55.72</u>	-	3.3	-
( $m=3$ , w/o Think)				
ITEM-As-ExpA	51.69	<b>53.33</b>	3.4	10.3
ITEM-As-Imp	52.38	52.34	2.4	3.6
ITEM-ARs-ExpA	<b>56.02</b>	-	7.1	-

Table 3: The F1 performance (%) of utility judgments with Qwen 3-8B on the GTI-NQ dataset. **bold** and Underline are defined in Table 1. “w/ thinking” and “w/o thinking” refer to the model generating with its thinking function enabled and disabled, respectively. Due to inference cost constraints, our evaluation of the thinking function was conducted under the vanilla listwise input setting. The terms “List” and “Point” refer to the “Listwise” and “Pointwise” approaches.

ImpA, respectively. Explicit generation of pseudo-answers by LLMs enhances their performance in utility judgment tasks, highlighting the importance of task interaction. Moreover, concurrent execution of answer generation and utility judgment within a single inference cycle yields inferior performance compared to sequential task execution through separate reasoning phases.

**ITEM with Multiple Iterations vs. ITEM with Single Iteration.** All LLMs using our ITEM-A and ITEM-AR generally demonstrate improved performance with multiple iterations compared to single iterations on all three datasets. For instance, on the WebAP dataset, Mistral, Llama 3, and ChatGPT (using our ITEM-A w. ExpA) improved their F1 scores in the listwise approach by 6.4%, 6.6%,

and 7.3%, respectively, after multiple iterations. Moreover, our method achieves state-of-the-art performance compared to all baselines by leveraging the iterative framework. The performance improvement from multiple iterations underscores the importance of iterative interaction and further supports Schutz’s interactive framework. Moreover, ChatGPT outperforms other LLMs on all datasets using both input approaches.

**ITEM-A<sub>s</sub> vs. ITEM-AR<sub>s</sub>.** In our utility-emphasized iterative RAG framework, ITEM-A<sub>s</sub> and ITEM-AR<sub>s</sub> are the two major methods we propose. From Table 1, we find that ITEM-AR<sub>s</sub> works better than ITEM-A<sub>s</sub> most of the time for complex questions (WebAP, all the questions are non-factoid) and the complex candidate passage list (GTI-NQ, containing different types of passage), indicating complicated questions or passage lists need more components in the loop. For TREC DL, which contains factoid questions, we find that ITEM-AR<sub>s</sub> is worse than ITEM-A<sub>s</sub> most times on TREC DL. This is reasonable since factoid questions are relatively easier to answer and may not need more components involved in the iteration.

**Listwise vs. Pointwise.** The general performance of utility judgments for LLMs is better with the listwise approach than with the pointwise approach. The primary rationale lying in the listwise approach exposes the LLM to broader contextual information, thereby facilitating more effective interaction during the LLMs in judging the passages’ utility.

**Thinking vs Non-thinking.** As shown in Table 3, the long reasoning mode demonstrates superior performance compared to other single-shot baselines and the  $k$ -sampling method. Our proposed ITEM framework achieves comparable performance to this long reasoning method but at a significantly lower computational cost (only about 23% of long reasoning). This indicates that the ITEM framework successfully balances efficiency with effectiveness, offering a more practical and resource-efficient solution for RAG.

## 5.2 Ranking Performance

We also assess whether the ranking performance has been improved within ITEM on retrieval datasets (Table 4) and utility judgments benchmark (Table 5). In terms of ranking performance, we consider two rankings: relevance ranking (ITEM-AR<sub>s</sub>) and utility ranking (ITEM-A<sub>r</sub>). We can observe that: (i) Our ITEM with a single iteration significantly improves the ranking of topical relevance

Method	Mistral		Llama 3		ChatGPT	
	TREC	WebAP	TREC	WebAP	TREC	WebAP
$D$	58.69	21.89	58.69	21.89	58.69	21.89
RankGPT	69.81	29.34	75.61	41.73	80.56	42.49
ITEM-A <sub>r</sub> (1)	70.57	37.11	73.95	40.89	80.79	50.30
ITEM-AR <sub>s</sub> (1)	<u>71.29</u>	<u>37.48</u>	<u>77.22</u>	43.80	81.38	48.42
ITEM-A <sub>r</sub> (3)	<b>74.27</b>	43.80	<b>77.34</b>	<b>45.88</b>	<b>83.12</b>	<b>51.61</b>
ITEM-AR <sub>s</sub> (3)	73.24	<u>45.45</u>	74.80	44.87	82.89	48.80

Table 4: The NDCG@5 performance (%) of the ranking using different LLMs on the different datasets. **Bold** and Underline are defined in Table 1.

Method	Llama 3		ChatGPT	
	@5	@10	@5	@10
$D$	29.46	45.26	29.46	45.26
RankGPT	71.50	74.05	77.27	78.64
ITEM-A <sub>r</sub> (1)	74.36	76.91	85.99	87.26
ITEM-AR <sub>s</sub> (1)	<u>75.46</u>	<u>77.75</u>	84.54	85.14
ITEM-A <sub>r</sub> (3)	75.95	78.18	<b>87.48</b>	<b>88.47</b>
ITEM-AR <sub>s</sub> (3)	<b>76.38</b>	<b>78.56</b>	85.95	86.39

Table 5: The NDCG performance (%) of the ranking using different LLMs on the GTI-NQ dataset. **Bold** and Underline are defined in Table 1.

performance compared to the RankGPT. For instance, relevance ranking outperforms RankGPT in NDCG@5 by 2.1% on the TREC dataset. The performance improvement may stem from the interaction between tasks. (ii) After iterations, relevance and utility ranking performance have been improved on all datasets and all LLMs. The ranking benefits from our dynamic iterative framework, confirming Schutz’s theory of dynamic iterative interaction. (iii) From Tables 4&5, we can find that ITEM-AR<sub>s</sub> is generally better than ITEM-A<sub>r</sub> when  $m = 1$ . However, when  $m = 3$ , it may have the opposite performance. The possible reason is that when  $m$  is small, the answer is not very good, and utility is more dependent on the answer than on relevance. As iterations proceed, we observe improved answer quality and utility performance. In contrast, relevance does not show as marked an improvement.

## 5.3 Results of Answer Generation

In the answer generation task, the results of utility judgments are fed to LLMs for answer generation. We use the factoid QA dataset (i.e., NQ) for answer generation evaluation, as shown in Table 6. From Tables 1&2, we find that the listwise approach generally outperforms the pointwise approach for utility judgments. Consequently, our answer generation experiments utilize only the listwise utility judgments. The following observations can be made from Table 6: (i) ITEM outperforms

References	Mistral		ChatGPT	
	EM	F1	EM	F1
Golden	46.09	62.59	66.40	76.86
<i>D</i>	31.58	47.69	46.54	57.00
Vanilla	31.16	47.43	48.52	58.64
UJ-ExpA	32.76	48.46	47.72	58.01
UJ-ImpA	30.67	46.83	49.01	59.30
5-sampling	33.24	48.84	48.90	58.97
ITEM-A <sub>s</sub> (1)	32.98	49.00	49.38	59.78
ITEM-AR <sub>s</sub> (1)	<u>33.30</u>	<u>49.26</u>	<u>49.52</u>	<u>59.64</u>
ITEM-A <sub>s</sub> (3)	<b>33.73</b>	<b>49.63</b>	<b>49.69</b>	<b>60.18</b>
ITEM-AR <sub>s</sub> (3)	33.40	49.27	49.06	59.67

Table 6: The answer generation performance (%) of all LLMs on the NQ dataset using reference passages collected from different methods. **Bold** means the best performance except for the answer generation with golden evidence. Underline is defined in Table 1.

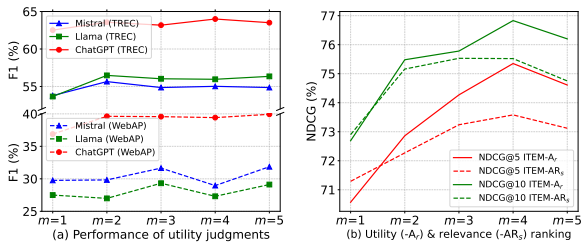


Figure 7: (a): utility judgments performance (%) in terms of  $m$  values in ITEM-A<sub>s</sub>. (b): relevance ranking (ITEM-AR<sub>s</sub>) and utility ranking (ITEM-A<sub>r</sub>) performance (%) of Mistral on the TREC DL dataset.

baselines across all metrics on all LLMs (except for the EM score of Llama 3), indicating that ITEM can help the LLMs to find better evidence for generating answers. (ii) Similar to Tables 1&2, when the  $m = 1$ , ITEM-AR<sub>s</sub> performs better than ITEM-A<sub>s</sub>, which shows the importance of relevance reranking in ITEM. However, as the number of iterations increases, ITEM-A<sub>s</sub> performs better than ITEM-AR<sub>s</sub>. We are keen to discern the optimal choice for different tasks: 1) More components and more iterations are not always needed, especially for simpler tasks; 2) Fewer iterations with numerous components, or increased iterations with few components.

## 6 Further Analyses

**Iteration Rounds.** Figure 7 shows the performance of (a): utility judgments under ITEM-A<sub>s</sub> and (b): ranking with varying maximum iteration rounds  $m$ . We observe the following: 1) Varying the value of  $m$  affects the performance of utility judgments and ranking. 2) Based on empirical observations balancing the cost and performance,  $m$  was operationally configured with distinct values for different question types on utility judgments ( $m=3$  in our paper on all experiments for fair comparison):  $m=2$  for factoid questions, whereas  $m=3$  is better for non-factoid questions in practical appli-

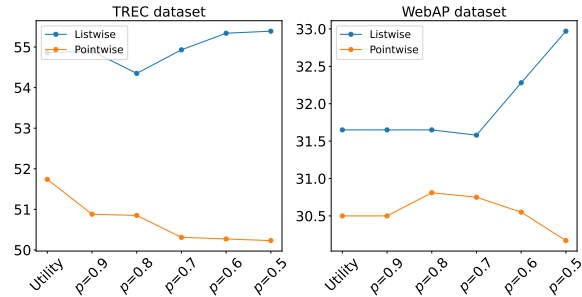


Figure 8: The utility judgments F1 performance (%) of Mistral in different iteration stop conditions ( $m=3$ ) under ITEM-A<sub>s</sub>.

cations. 3) Utility ranking generally outperforms relevance ranking, which confirms the effectiveness of utility in the ranking task.

**Iteration Stop Conditions.** In our experiment, we employ utility-based selection results as the stopping criterion—that is, the iteration halts when the selection outcomes remain identical across two consecutive rounds. In this section, we also evaluate an alternative stopping condition based on the pseudo-answer generation performance of ITEM. Specifically, we calculate the ROUGE-L score (Lin, 2004) of the answer in two iterations and stop the iteration if the ROUGE-L of  $a_t$  and  $a_{t-1}$  is greater than  $p$ . The utility judgment performance of different iteration stop conditions is shown in Figure 8. The results show that using different stopping conditions affects the performance of utility judgments. However, using the answer as a stopping condition, different LLMs on different datasets may need to look for different  $p$ , which is not very flexible.

## 7 Conclusion

In this paper, we propose an Iterative uTility-focused Evidence refineMent (ITEM) to enhance the utility judgment and QA performance of LLMs by interactions between the steps, inspired by Schutz’s philosophical discussion of relevance. This is a unified framework of iterative RAG with an emphasis on utility. Our framework achieves state-of-the-art performance in zero-shot scenarios, outperforming previous methods in utility judgments, ranking of topical relevance, and answer generation tasks, indicating that the cognitive process of LLMs on a specific topic can also be improved by a similar process. Our experiments also highlight the significance of dynamic interaction in achieving high performance and stability. Future directions include developing better fine-tuning strategies for utility judgments and creating end-to-end solutions for RAG.

## 611 Limitations

612 There are two primary limitations that should be  
613 acknowledged: (i) Our methods are applied in  
614 zero-shot scenarios without any training. The  
615 zero-shot approach itself does not enhance the  
616 LLMs’s inherent capability in utility judgments but  
617 rather employs strategies to improve performance  
618 on utility judgment tasks. Future research should  
619 explore designing more effective training methods,  
620 e.g., utilizing our iterative framework with self-evo-  
621 lution techniques (Singh et al., 2023), to genuinely  
622 enhance the LLMs’s ability in utility judgments  
623 through training. (ii) The number of candidate pas-  
624 sages in the search scenario is much larger than  
625 20. The number of search results we assumed is  
626 too small. We need to continue to study utility  
627 judgments in large-scale scenarios in the future.

## 628 8 Ethics Statement

629 Our research does not rely on personally identi-  
630 fiable information. All datasets and models used  
631 in our paper are publicly available and have been  
632 widely adopted by researchers. We firmly believe  
633 in the principles of open research and the scientific  
634 value of reproducibility. To this end, we have made  
635 all data, and code associated with our paper pub-  
636 licly available on GitHub. This transparency not  
637 only facilitates the verification of our findings by  
638 the community but also encourages the application  
639 of our methods in other contexts.

## 640 References

641 Keping Bi, Qingyao Ai, and W Bruce Croft. 2019. *Iter-  
642 ative relevance feedback for answer passage retrieval  
643 with passage-level semantic match*. In *Advances  
644 in Information Retrieval: 41st European Conference  
645 on IR Research, ECIR 2019, Cologne, Germany,  
646 April 14–18, 2019, Proceedings, Part I 41*, pages  
647 558–572. Springer.

648 Sebastian Borgeaud, Arthur Mensch, Jordan Hoff-  
649 mann, Trevor Cai, Eliza Rutherford, Katie Millin-  
650 can, George Bm Van Den Driessche, Jean-Baptiste  
651 Lespiau, Bogdan Damoc, Aidan Clark, and 1 oth-  
652 ers. 2022. *Improving language models by retrieving  
653 from trillions of tokens*. In *International conference  
654 on machine learning*, pages 2206–2240. PMLR.

655 Harry W Bruce. 1994. A cognitive view of the sit-  
656 uational dynamism of user-centered relevance es-  
657 timation. *Journal of the American Society for  
658 Information Science*, 45(3):142–148.

659 Daniel Cohen, Liu Yang, and W Bruce Croft. 2018.  
660 *Wikipassageqa: A benchmark collection for research*

*on non-factoid answer passage retrieval*. In *The 41st  
661 international ACM SIGIR conference on research &  
662 development in information retrieval*, pages 1165–  
663 1168.

664 William S Cooper. 1971. A definition of relevance  
665 for information retrieval. *Information storage and  
666 retrieval*, 7(1):19–37.

667 Nick Craswell, Bhaskar Mitra, Emine Yilmaz, Daniel  
668 Campos, and Ellen M Voorhees. 2020. *Overview  
669 of the trec 2019 deep learning track*. *arXiv preprint  
670 arXiv:2003.07820*.

671 Michael Glass, Gaetano Rossiello, Md Faisal Mahub  
672 Chowdhury, Ankita Naik, Pengshan Cai, and Al-  
673 fio Gliozzo. 2022. *Re2G: Retrieve, rerank, gen-  
674 erate*. In *Proceedings of the 2022 Conference  
675 of the North American Chapter of the Association  
676 for Computational Linguistics: Human Language  
677 Technologies*, pages 2701–2715, Seattle, United  
678 States. Association for Computational Linguistics.

679 William Goffman and Vaun A Newill. 1964.  
680 *Methodology for test and evaluation of information  
681 retrieval systems*. Center for Documentation and  
682 Communication Research, School of Library . . . .  
683

684 Helia Hashemi, Mohammad Aliannejadi, Hamed Za-  
685 mani, and W Bruce Croft. 2020. *Antique: A  
686 non-factoid question answering benchmark*. In  
687 *Advances in Information Retrieval: 42nd European  
688 Conference on IR Research, ECIR 2020, Lisbon,  
689 Portugal, April 14–17, 2020, Proceedings, Part II 42*,  
690 pages 166–173. Springer.

691 Helia Hashemi, Hamed Zamani, and W Bruce Croft.  
692 2019. *Performance prediction for non-factoid  
693 question answering*. In *Proceedings of the 2019  
694 ACM SIGIR International Conference on Theory of  
695 Information Retrieval*, pages 55–58.

696 Donald J Hillman. 1964. The notion of relevance (i).  
697 *American Documentation*, 15(1):26–34.

698 Gautier Izacard, Patrick Lewis, Maria Lomeli, Lucas  
699 Hosseini, Fabio Petroni, Timo Schick, Jane Dwivedi-  
700 Yu, Armand Joulin, Sebastian Riedel, and Edouard  
701 Grave. 2023. *Atlas: Few-shot learning with retrieval  
702 augmented language models*. *Journal of Machine  
703 Learning Research*, 24(251):1–43.

704 Kalervo Järvelin and Jaana Kekäläinen. 2017. *Ir eval-  
705 uation methods for retrieving highly relevant docu-  
706 ments*. In *ACM SIGIR Forum*, volume 51, pages  
707 243–250. ACM New York, NY, USA.

708 Kaixin Ji, Danula Hettiachchi, Flora D. Salim, Falk  
709 Scholer, and Damiano Spina. 2024. *Characterizing  
710 information seeking processes with multiple physio-  
711 logical signals*. *SIGIR*.

712 Albert Q Jiang, Alexandre Sablayrolles, Arthur Men-  
713 sch, Chris Bamford, Devendra Singh Chaplot, Diego  
714 de las Casas, Florian Bressand, Gianna Lengyel, Guil-  
715 laume Lample, Lucile Saulnier, and 1 others. 2023a.  
716 *Mistral 7b*. *arXiv preprint arXiv:2310.06825*.

717	Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun,	Stefano Mizzaro. 1998. How many relevances in in-	770
718	Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie	formation retrieval? <u>Interacting with computers</u> ,	771
719	Callan, and Graham Neubig. 2023b. <u>Active re-</u>	10(3):303–320.	772
720	<u>trieval augmented generation</u> . In <u>Proceedings of the</u>		
721	<u>2023 Conference on Empirical Methods in Natural</u>	OpenAI. 2022. <u>Introducing chatgpt</u> .	773
722	<u>Language Processing</u> , pages 7969–7992, Singapore.		
723	Association for Computational Linguistics.		
724	Mostafa Keikha, Jae Hyun Park, and W Bruce Croft.	Ori Ram, Yoav Levine, Itay Dalmedigos, Dor Muhlgay,	774
725	2014a. <u>Evaluating answer passages using sum-</u>	Amnon Shashua, Kevin Leyton-Brown, and Yoav	775
726	<u>marization measures</u> . In <u>Proceedings of the 37th</u>	Shoham. 2023. <u>In-context retrieval-augmented lan-</u>	776
727	<u>international ACM SIGIR conference on Research</u>	<u>guage models</u> . <u>Transactions of the Association for</u>	777
728	<u>&amp; development in information retrieval</u> , pages 963–	<u>Computational Linguistics</u> , 11:1316–1331.	778
729	966.		
730	Mostafa Keikha, Jae Hyun Park, W Bruce Croft,	Ruiyang Ren, Yingqi Qu, Jing Liu, Wayne Xin Zhao,	779
731	and Mark Sanderson. 2014b. <u>Retrieving passages</u>	Qiaoqiao She, Hua Wu, Haifeng Wang, and Ji-Rong	780
732	<u>and finding answers</u> . In <u>Proceedings of the 19th</u>	Wen. 2021. <u>Rocketqav2: A joint training method</u>	781
733	<u>Australasian Document Computing Symposium</u> ,	<u>for dense passage retrieval and passage re-ranking</u> .	782
734	pages 81–84.	<u>EMNLP</u> .	783
735	DA Kemp. 1974. Relevance, pertinence and informa-	Stephen Robertson, Hugo Zaragoza, and 1 others. 2009.	784
736	tion system development. <u>Information Storage and</u>	<u>The probabilistic relevance framework: Bm25 and</u>	785
737	<u>Retrieval</u> , 10(2):37–47.	<u>beyond</u> . <u>Foundations and Trends® in Information</u>	786
738		<u>Retrieval</u> , 3(4):333–389.	787
739	Urvashi Khandelwal, Omer Levy, Dan Jurafsky, Luke	Tefko Saracevic. 1975. <u>Relevance: A review of and</u>	788
740	Zettlemoyer, and Mike Lewis. 2020. <u>Generalization</u>	<u>a framework for the thinking on the notion in infor-</u>	789
741	<u>through memorization: Nearest neighbor language</u>	<u>mation science</u> . <u>Journal of the American Society for</u>	790
742	<u>models</u> . <u>ICLR</u> .	<u>information science</u> , 26(6):321–343.	791
743	Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-	Tefko Saracevic. 1996. Relevance reconsidered.	792
744	field, Michael Collins, Ankur Parikh, Chris Alberti,	In <u>Proceedings of the second conference on</u>	793
745	Danielle Epstein, Illia Polosukhin, Jacob Devlin,	<u>conceptions of library and information science</u>	794
746	Kenton Lee, and 1 others. 2019. <u>Natural ques-</u>	<u>(CoLIS 2)</u> , pages 201–218.	795
747	<u>tions: a benchmark for question answering research</u> .		
748	<u>Transactions of the Association for Computational</u>	Tefko Saracevic, Paul Kantor, Alice Y Chamis, and	796
749	<u>Linguistics</u> , 7:453–466.	Donna Trivison. 1988. <u>A study of information seek-</u>	797
750	F Wilfrid Lancaster. 1968. Information retrieval sys-	<u>ing and retrieving. i. background and methodology</u> .	798
751	tems: Characteristics, testing, and evaluation. ( <u>No</u>	<u>Journal of the American Society for Information</u>	799
752	<u>Title</u> ).	<u>science</u> , 39(3):161–176.	800
753	Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio	Linda Schamber and Michael Eisenberg. 1988. Rele-	801
754	Petroni, Vladimir Karpukhin, Naman Goyal, Hein-	vance: The search for a definition.	802
755	rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-	Linda Schamber, Michael B Eisenberg, and Michael S	803
756	täschel, and 1 others. 2020. <u>Retrieval-augmented gen-</u>	Nilan. 1990. <u>A re-examination of relevance: to-</u>	804
757	<u>eration for knowledge-intensive nlp tasks</u> . <u>Advances</u>	<u>ward a dynamic, situational definition</u> . <u>Information</u>	805
758	<u>in Neural Information Processing Systems</u> , 33:9459–	<u>processing &amp; management</u> , 26(6):755–776.	806
759	9474.	Alfred Schutz. 1970. <u>Reflections on the Problem of</u>	807
760	Xiaonan Li, Changtai Zhu, Linyang Li, Zhangyue Yin,	<u>Relevance</u> . Greenwood Press, Westport, Conn.	808
761	Tianxiang Sun, and Xipeng Qiu. 2023. <u>Llatrieval:</u>	Zhihong Shao, Yeyun Gong, Yelong Shen, Minlie	809
762	<u>Llm-verified retrieval for verifiable generation</u> . <u>arXiv</u>	Huang, Nan Duan, and Weizhu Chen. 2023. <u>Enhanc-</u>	810
763	<u>preprint arXiv:2311.07838</u> .	<u>ing retrieval-augmented large language models with</u>	811
764	Chin-Yew Lin. 2004. <u>Rouge: A package for automatic</u>	<u>iterative retrieval-generation synergy</u> . In <u>Findings</u>	812
765	<u>evaluation of summaries</u> . In <u>Text summarization</u>	<u>of the Association for Computational Linguistics:</u>	813
766	<u>branches out</u> , pages 74–81.	<u>EMNLP 2023</u> , pages 9248–9274, Singapore. Associ-	814
767	Meta. 2024. <u>Welcome llama 3 - meta's new open llm</u> .	ation for Computational Linguistics.	815
768	Stefano Mizzaro. 1997. Relevance: The whole history.	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	816
769	<u>Journal of the American society for information</u>	joon Seo, Rich James, Mike Lewis, Luke Zettle-	817
	<u>science</u> , 48(9):810–832.	moyer, and Wen-tau Yih. 2023. <u>Replug: Retrieval-</u>	818
		<u>augmented black-box language models</u> . <u>arXiv</u>	819
		<u>preprint arXiv:2301.12652</u> .	820

821	Avi Singh, John D Co-Reyes, Rishabh Agarwal, Ankesh Anand, Piyush Patil, Peter J Liu, James Harrison, Jaehoon Lee, Kelvin Xu, Aaron Parisi, and 1 others. 2023. <a href="#">Beyond human data: Scaling self-training for problem-solving with language models</a> . <a href="#">arXiv preprint arXiv:2312.06585</a> .	878
822		879
823		880
824		881
825		882
826		883
827	Weihsiang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. <a href="#">Dragin: Dynamic retrieval augmented generation based on the real-time information needs of large language models</a> . <a href="#">arXiv preprint arXiv:2403.10081</a> .	884
828		885
829		886
830		887
831		888
832	Weiwei Sun, Lingyong Yan, Xinyu Ma, Pengjie Ren, Dawei Yin, and Zhaochun Ren. 2023. <a href="#">Is chatgpt good at search? investigating large language models as re-ranking agent</a> . <a href="#">EMNLP</a> .	889
833		890
834		
835		
836	Don R Swanson. 1986. Subjective versus objective relevance in bibliographic retrieval systems. <a href="#">The library quarterly</a> , 56(4):389–398.	
837		
838		
839	Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. <a href="#">Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions</a> . In <a href="#">Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)</a> , ACL 2023, Toronto, Canada, July 9-14, 2023, pages 10014–10037. Association for Computational Linguistics.	
840		
841		
842		
843		
844		
845		
846		
847		
848	Brian C Vickery. 1959. The structure of information retrieval systems. In <a href="#">Proceedings of the International Conference on Scientific Information</a> , volume 2, pages 1275–1290.	
849		
850		
851		
852	Ellen M Voorhees and 1 others. 2003. <a href="#">Overview of the trec 2003 robust retrieval track</a> . In <a href="#">Trec</a> , pages 69–77.	
853		
854		
855	Xuezhi Wang, Jason Wei, Dale Schuurmans, Quoc Le, Ed Chi, Sharan Narang, Aakanksha Chowdhery, and Denny Zhou. 2022. <a href="#">Self-consistency improves chain of thought reasoning in language models</a> . <a href="#">arXiv preprint arXiv:2203.11171</a> .	
856		
857		
858		
859		
860	Jian Xie, Kai Zhang, Jiangjie Chen, Renze Lou, and Yu Su. 2023. <a href="#">Adaptive chameleon or stubborn sloth: Revealing the behavior of large language models in knowledge conflicts</a> . In <a href="#">The Twelfth International Conference on Learning Representations</a> .	
861		
862		
863		
864		
865	An Yang, Anfeng Li, Baosong Yang, Beichen Zhang, Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen Huang, Chenxu Lv, and 1 others. 2025. <a href="#">Qwen3 technical report</a> . <a href="#">arXiv preprint arXiv:2505.09388</a> .	
866		
867		
868		
869		
870	Liu Yang, Qingyao Ai, Damiano Spina, Ruey-Cheng Chen, Liang Pang, W Bruce Croft, Jiafeng Guo, and Falk Scholer. 2016. <a href="#">Beyond factoid qa: effective methods for non-factoid answer sentence retrieval</a> . In <a href="#">Advances in Information Retrieval: 38th European Conference on IR Research, ECIR 2016, Padua, Italy, March 20–23, 2016. Proceedings 38</a> , pages 115–128. Springer.	
871		
872		
873		
874		
875		
876		
877		
	Hengran Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2023. <a href="#">From relevance to utility: Evidence retrieval with feedback for fact verification</a> . In <a href="#">Findings of the Association for Computational Linguistics: EMNLP 2023</a> , pages 6373–6384, Singapore. Association for Computational Linguistics.	878
		879
		880
		881
		882
		883
		884
	Hengran Zhang, Ruqing Zhang, Jiafeng Guo, Maarten de Rijke, Yixing Fan, and Xueqi Cheng. 2024. <a href="#">Are large language models good at utility judgments? Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR '24)</a> .	885
		886
		887
		888
		889
		890
	Xinran Zhao, Tong Chen, Sihao Chen, Hongming Zhang, and Tongshuang Wu. 2024. <a href="#">Beyond relevance: Evaluate and improve retrievers on perspective awareness</a> . <a href="#">arXiv preprint arXiv:2405.02714</a> .	891
		892
		893
		894
	Yujia Zhou, Zheng Liu, Jiajie Jin, Jian-Yun Nie, and Zhicheng Dou. 2024. <a href="#">Metacognitive retrieval-augmented large language models</a> . In <a href="#">Proceedings of the ACM on Web Conference 2024, WWW '24</a> , page 1453–1463, New York, NY, USA. Association for Computing Machinery.	895
		896
		897
		898
		899
		900

## A Experiment Details

### A.1 Effect of Iteration Numbers

The precision, recall, and F1 performance of different LLMs on different datasets with different iteration numbers is shown in Table 7, Table 8, Table 9, and Table 10.

Method	TREC						WebAP					
	listwise			pointwise			listwise			pointwise		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Vanilla	36.82	60.13	45.67	29.92	<b>91.61</b>	45.11	13.07	50.83	20.79	13.30	86.29	23.05
UJ-ExpA	48.51	61.15	54.10	28.12	96.27	43.53	18.83	54.16	27.94	14.65	<b>91.82</b>	25.27
UJ-ImpA	40.16	60.53	48.29	33.95	83.76	48.31	16.46	52.45	25.06	17.56	73.55	28.35
5-sampling	46.64	59.56	52.31	-	-	-	20.61	56.22	30.16	-	-	-
ITEM-A <sub>s</sub> w. ExpA (m=1)	48.07	61.04	53.78	34.21	89.11	49.44	20.57	53.81	29.76	17.86	78.41	29.10
ITEM-A <sub>s</sub> w. ExpA (m=2)	50.58	61.86	55.65	35.87	88.73	51.09	21.11	50.85	29.83	18.27	82.00	29.88
ITEM-A <sub>s</sub> w. ExpA (m=3)	50.61	59.88	54.86	36.23	90.46	51.74	23.57	48.14	31.65	18.73	81.96	30.50
ITEM-A <sub>s</sub> w. ExpA (m=4)	50.01	61.15	55.02	<b>36.41</b>	90.36	<b>51.90</b>	21.44	44.62	28.96	<b>19.19</b>	80.59	<b>31.00</b>
ITEM-A <sub>s</sub> w. ExpA (m=5)	50.61	59.88	54.86	36.14	90.46	51.65	24.07	47.09	31.86	19.17	78.94	30.85
ITEM-A <sub>s</sub> w. ImpA (m=1)	39.97	64.62	49.39	30.98	89.38	46.01	16.88	57.13	26.06	17.10	81.65	28.28
ITEM-A <sub>s</sub> w. ImpA (m=2)	43.14	61.52	50.72	30.90	87.00	45.60	19.41	54.82	28.67	18.88	78.06	30.40
ITEM-A <sub>s</sub> w. ImpA (m=3)	44.43	62.82	52.05	31.68	87.99	46.59	19.21	54.20	28.36	18.69	77.77	30.13
ITEM-A <sub>s</sub> w. ImpA (m=4)	44.72	61.29	51.71	31.66	87.40	46.49	17.44	47.11	25.46	18.95	78.06	30.50
ITEM-A <sub>s</sub> w. ImpA (m=5)	44.63	60.98	51.54	31.80	89.32	46.91	18.98	48.88	27.35	19.05	76.69	30.52
ITEM-AR <sub>s</sub> (m=1)	43.65	65.34	52.34	-	-	-	25.04	<b>60.99</b>	35.50	-	-	-
ITEM-AR <sub>s</sub> (m=2)	45.10	65.46	53.40	-	-	-	24.42	51.97	33.23	-	-	-
ITEM-AR <sub>s</sub> (m=3)	49.07	<b>65.96</b>	56.27	-	-	-	<b>27.70</b>	55.95	<b>37.06</b>	-	-	-
ITEM-AR <sub>s</sub> (m=4)	50.96	62.32	56.07	-	-	-	23.77	53.40	32.90	-	-	-
ITEM-AR <sub>s</sub> (m=5)	<b>53.01</b>	63.60	<b>57.82</b>	-	-	-	25.85	47.56	33.50	-	-	-

Table 7: The utility judgments performance (%) of Mistral on retrieval datasets (Numbers in parentheses represent  $m$ -values). Numbers in bold indicate the best performance.

### A.2 Quality of Utility Judgments

The relevance labels of TREC DL are of a four-point scale, and we consider the highest level as having utility. To see the utility judgment performance when we consider lower grades to have utility, we measure the precision of utility judgments of Mistral on TREC DL when passages of different grades are treated as positive in Table 11. We can see that almost 70% of the results of positive utility judgments are highly relevant to the question.

### A.3 $k$ values in ITEM-A<sub>r</sub>

Different ranking performance of  $k$  values in ITEM-A<sub>r</sub> is shown in Table 12. Considering the performance of utility ranking and utility judgments, we set  $k=5$ .

### A.4 ITEM-AR<sub>r</sub>

We evaluate two ranking performances of ITEM-AR<sub>r</sub> during the same loop, with the experimental results shown in Table 13. We find that under the ITEM-AR<sub>r</sub> framework, relevance ranking and utility ranking are both improved, and utility ranking performance is generally better than relevance ranking. However, as seen in Table 4 and Table 13, performing ranking twice in the same iteration may not yield better ranking results than performing utility ranking once in the iteration. All experiments in the work are zero-shot. For the reproducibility of the experiments, the temperature is set to 0, and thus all experiments were run once.

## B More Baseline

For methods that involve multiple calls, k-sampling(our baseline) randomly shuffles the input and aggregates multiple results to obtain the final utility judgments. We reproduced the Self-consistency method (Wang et al., 2022), which also leverages multiple LLM calls, on the utility judgments task. That is, we sampled multiple results from the same input (The prompt is the same as UJ-ExpA) and

Method	TREC						WebAP					
	listwise			pointwise			listwise			pointwise		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Vanilla	34.67	<b>85.80</b>	49.39	31.42	<b>98.47</b>	47.64	12.69	77.15	21.79	14.65	<b>87.36</b>	25.09
UJ-ExpA	39.21	80.98	52.83	38.27	90.15	53.73	16.32	77.92	26.99	18.04	77.15	29.25
UJ-ImpA	33.92	83.36	48.22	38.68	71.47	50.20	15.57	<b>82.79</b>	26.22	17.22	47.61	25.29
5-sampling	39.04	80.98	52.68	-	-	-	17.52	83.49	28.97	-	-	-
ITEM-A <sub>s</sub> w. ExpA ( $m=1$ )	39.68	82.88	53.66	37.58	84.84	52.09	17.54	63.67	27.50	19.65	74.31	31.08
ITEM-A <sub>s</sub> w. ExpA ( $m=2$ )	<b>42.35</b>	84.77	<b>56.48</b>	38.25	84.58	52.68	17.39	60.25	26.99	20.23	73.01	31.68
ITEM-A <sub>s</sub> w. ExpA ( $m=3$ )	42.00	84.15	56.03	37.84	85.50	52.46	19.12	62.87	29.32	<b>20.91</b>	74.63	32.67
ITEM-A <sub>s</sub> w. ExpA ( $m=4$ )	41.85	84.41	55.96	38.12	85.16	52.67	17.53	61.85	27.31	20.44	73.83	32.02
ITEM-A <sub>s</sub> w. ExpA ( $m=5$ )	42.36	84.15	56.35	37.35	84.69	51.84	18.94	62.87	29.12	20.88	75.45	<b>32.71</b>
ITEM-A <sub>s</sub> w. ImpA ( $m=1$ )	39.63	83.42	53.73	39.70	82.87	53.68	15.48	73.66	25.59	20.04	64.06	30.53
ITEM-A <sub>s</sub> w. ImpA ( $m=2$ )	38.75	85.63	53.35	38.15	82.36	52.14	15.50	76.47	25.77	18.54	62.69	28.62
ITEM-A <sub>s</sub> w. ImpA ( $m=3$ )	40.84	84.86	55.14	40.58	79.64	53.76	15.99	70.99	26.10	19.54	61.32	29.64
ITEM-A <sub>s</sub> w. ImpA ( $m=4$ )	38.88	82.74	52.90	39.34	81.74	53.12	15.03	74.41	25.01	19.72	59.95	29.68
ITEM-A <sub>s</sub> w. ImpA ( $m=5$ )	41.26	84.61	55.47	<b>40.92</b>	82.14	<b>54.63</b>	15.49	68.93	25.29	19.84	57.21	29.46
ITEM-AR <sub>s</sub> ( $m=1$ )	34.53	84.17	48.97	-	-	-	<b>20.05</b>	72.88	<b>31.44</b>	-	-	-
ITEM-AR <sub>s</sub> ( $m=2$ )	36.27	83.19	50.51	-	-	-	15.92	79.01	26.50	-	-	-
ITEM-AR <sub>s</sub> ( $m=3$ )	38.04	82.68	52.10	-	-	-	17.93	76.87	29.08	-	-	-
ITEM-AR <sub>s</sub> ( $m=4$ )	37.28	83.70	51.58	-	-	-	16.60	78.81	27.42	-	-	-
ITEM-AR <sub>s</sub> ( $m=5$ )	40.25	81.37	53.86	-	-	-	17.04	74.83	27.75	-	-	-

Table 8: The utility judgments performance (%) of Llama 3 on retrieval datasets (Numbers in parentheses represent  $m$ -values). Numbers in bold indicate the best performance.

Method	TREC						WebAP					
	listwise			pointwise			listwise			pointwise		
	P	R	F1	P	R	F1	P	R	F1	P	R	F1
Vanilla	42.13	<b>79.98</b>	55.19	33.86	94.40	49.84	17.13	<b>83.45</b>	28.43	15.80	<b>89.42</b>	26.85
UJ-ExpA	45.74	77.36	57.49	32.06	<b>96.19</b>	48.09	19.51	69.86	30.50	16.23	88.74	27.44
UJ-ImpA	44.19	77.11	56.18	33.45	90.36	48.83	18.37	80.14	29.89	15.58	84.51	26.32
5-sampling	50.78	74.77	60.49	-	-	-	20.70	65.83	31.49	-	-	-
ITEM-A <sub>s</sub> w. ExpA ( $m=1$ )	55.55	71.48	62.52	37.83	91.94	53.61	26.74	59.45	36.89	19.73	84.95	32.02
ITEM-A <sub>s</sub> w. ExpA ( $m=2$ )	57.95	70.40	63.57	40.74	93.04	56.67	29.43	60.58	39.62	19.62	78.62	31.40
ITEM-A <sub>s</sub> w. ExpA ( $m=3$ )	58.36	68.88	63.18	40.00	91.88	55.74	29.30	60.91	39.57	19.80	76.20	31.43
ITEM-A <sub>s</sub> w. ExpA ( $m=4$ )	<b>58.48</b>	70.67	<b>64.00</b>	40.25	93.38	56.25	29.11	61.03	39.42	20.48	79.63	32.58
ITEM-A <sub>s</sub> w. ExpA ( $m=5$ )	58.34	69.69	63.51	39.29	92.16	55.09	29.76	60.68	<b>39.93</b>	20.58	80.42	<b>32.77</b>
ITEM-A <sub>s</sub> w. ImpA ( $m=1$ )	54.36	65.08	59.24	40.89	82.20	54.61	24.79	64.37	35.80	18.78	67.00	29.34
ITEM-A <sub>s</sub> w. ImpA ( $m=2$ )	55.88	63.11	59.27	43.32	83.13	<b>56.96</b>	27.68	62.03	38.28	20.70	70.54	32.00
ITEM-A <sub>s</sub> w. ImpA ( $m=3$ )	57.33	64.17	60.56	41.66	80.48	54.90	<b>30.01</b>	63.60	40.78	21.51	66.77	32.54
ITEM-A <sub>s</sub> w. ImpA ( $m=4$ )	55.98	62.24	58.95	<b>42.34</b>	80.65	55.53	28.43	60.11	38.60	20.60	65.63	31.36
ITEM-A <sub>s</sub> w. ImpA ( $m=5$ )	56.63	62.19	59.28	41.49	83.57	55.45	29.05	60.66	39.29	<b>21.51</b>	68.03	32.68
ITEM-AR <sub>s</sub> ( $m=1$ )	51.94	76.90	62.00	-	-	-	25.32	65.84	36.58	-	-	-
ITEM-AR <sub>s</sub> ( $m=2$ )	53.77	76.19	63.05	-	-	-	25.55	59.26	35.70	-	-	-
ITEM-AR <sub>s</sub> ( $m=3$ )	52.41	74.04	61.37	-	-	-	27.61	63.96	38.58	-	-	-
ITEM-AR <sub>s</sub> ( $m=4$ )	52.75	73.78	61.52	-	-	-	28.84	61.85	39.34	-	-	-
ITEM-AR <sub>s</sub> ( $m=5$ )	52.77	76.28	62.39	-	-	-	28.76	62.54	39.40	-	-	-

Table 9: The utility judgments performance of ChatGPT on retrieval datasets (Numbers in parentheses represent  $m$ -values). Numbers in bold indicate the best performance.

References of RAG	Mistral		Llama 3		ChatGPT	
	EM	F1	EM	F1	EM	F1
Golden Evidence	46.09	62.59	64.45	76.64	66.40	76.86
RocketQAv2	31.58	47.69	<b>50.96</b>	62.01	46.54	57.00
Vanilla	31.16	47.43	49.09	60.56	48.52	58.64
UJ-ExpA	32.76	48.46	49.63	61.10	47.72	58.01
UJ-ImpA	30.67	46.83	48.88	60.26	49.01	59.30
5-sampling	33.24	48.84	48.72	60.71	48.90	58.97
ITEM- $A_s$ w. ExpA ( $m=1$ )	32.98	49.00	50.16	61.88	49.38	59.78
ITEM- $A_s$ w. ExpA ( $m=2$ )	<b>34.31</b>	<b>50.08</b>	50.48	<b>62.32</b>	49.22	59.99
ITEM- $A_s$ w. ExpA ( $m=3$ )	33.73	49.63	50.27	62.09	<b>49.69</b>	<b>60.18</b>
ITEM- $A_s$ w. ExpA ( $m=4$ )	34.21	50.07	50.43	62.20	-	-
ITEM- $A_s$ w. ExpA ( $m=5$ )	33.78	49.63	50.27	62.07	-	-
ITEM- $A_s$ w. ImpA ( $m=1$ )	32.17	48.51	50.37	61.89	48.75	58.99
ITEM- $A_s$ w. ImpA ( $m=2$ )	32.49	48.67	49.63	61.16	49.11	59.14
ITEM- $A_s$ w. ImpA ( $m=3$ )	32.39	48.47	49.68	61.48	48.69	58.94
ITEM- $A_s$ w. ImpA ( $m=4$ )	32.71	48.84	49.41	61.03	-	-
ITEM- $A_s$ w. ImpA ( $m=5$ )	32.33	48.44	49.73	61.42	-	-
ITEM- $AR_s$ ( $m=1$ )	33.30	49.26	50.27	61.69	49.52	59.64
ITEM- $AR_s$ ( $m=2$ )	33.57	49.16	50.70	61.92	49.01	59.75
ITEM- $AR_s$ ( $m=3$ )	33.40	49.27	49.36	60.97	49.06	59.67
ITEM- $AR_s$ ( $m=4$ )	33.46	49.24	49.84	61.54	-	-
ITEM- $AR_s$ ( $m=5$ )	33.89	49.58	49.20	60.84	-	-

Table 10: The answer generation performance (%) of all LLMs in the listwise approach. Numbers in bold indicate the best performance except the answer performance using golden evidence. Due to the high cost of using ChatGPT, we only tested with  $m=1,2,3$  on ChatGPT.

$m$	label $\geq 1$	label $\geq 2$	label $\geq 3$
$m=1$	82.08	68.34	48.07
$m=2$	83.86	69.53	50.58
$m=3$	84.23	<b>71.06</b>	<b>50.61</b>
$m=4$	<b>84.63</b>	70.18	50.01
$m=5$	84.52	70.69	50.61

Table 11: The precision scores (%) of utility judgments using Mistral in different  $m$  (iteration) values. “label” is the manual annotation in the original dataset, i.e., [3]: Perfectly relevant; [2]: Highly relevant; [1]: Related; [0]: Irrelevant.

$k, m$	Ranking					Utility judgments		
	N@1	N@3	N@5	N@10	N@20	P	R	F1
$k=1, m=1$	72.76	71.27	70.57	72.69	84.08	53.66	24.09	33.25
$k=1, m=2$	76.02	71.54	71.38	73.66	84.78	58.54	28.73	38.54
$k=1, m=3$	77.24	72.83	71.83	73.87	85.20	<b>59.76</b>	28.84	38.90
$k=1, m=4$	77.24	73.04	71.91	73.90	85.25	59.76	28.84	38.90
$k=1, m=5$	76.02	72.11	71.42	73.45	84.98	58.54	28.71	38.53
$k=5, m=1$	72.76	71.27	70.57	72.69	84.08	33.17	57.31	42.02
$k=5, m=2$	78.46	73.74	72.86	75.48	86.09	32.93	58.37	42.10
$k=5, m=3$	79.27	75.00	74.27	75.78	86.80	34.15	62.57	44.18
$k=5, m=4$	79.67	75.92	<b>75.35</b>	<b>76.83</b>	<b>87.23</b>	35.12	61.40	<b>44.68</b>
$k=5, m=5$	79.67	75.32	74.61	76.20	86.82	34.63	61.25	44.25
$k=10, m=1$	72.76	71.27	70.57	72.69	84.08	22.56	68.03	33.88
$k=10, m=2$	78.05	72.64	72.90	75.48	85.74	23.66	75.47	36.02
$k=10, m=3$	<b>80.89</b>	<b>76.58</b>	74.54	76.30	86.94	23.78	<b>75.65</b>	36.19
$k=10, m=4$	78.05	74.70	72.85	75.12	85.72	24.51	74.17	36.85
$k=10, m=5$	79.67	75.60	74.84	76.54	86.88	23.66	74.42	35.90

Table 12: The utility ranking performance and utility judgments performance of Mistral on TREC DL dataset in ITEM-AR. “N@k” means “NDCG@k”. Numbers in bold indicate the best performance.

$m$	NDCG@5	NDCG@10	NDCG@20	Utility-F1
1	<u>71.29</u> / <u>72.77</u>	<u>72.90</u> / <u>74.96</u>	<u>84.56</u> / <u>85.75</u>	43.13
2	<b>72.54</b> / 70.99	<u>74.81</u> / <u>73.76</u>	<u>85.77</u> / <u>85.28</u>	40.21
3	<u>72.07</u> / <b>74.14</b>	<u>74.14</u> / <b>76.63</b>	<u>85.53</u> / <b>86.57</b>	<b>45.67</b>
4	<u>71.02</u> / <u>71.06</u>	<u>74.30</u> / <u>74.03</u>	<u>85.09</u> / <u>85.16</u>	43.82
5	<u>72.26</u> / <u>70.12</u>	<b>75.83</b> / <u>72.59</u>	<b>85.88</b> / <u>84.77</u>	44.10

Table 13: Ranking of topical relevance and utility judgments performance (%) of ITEM-AR, using Mistral on the TREC DL dataset. “a/b” means “relevance ranking performance / utility ranking performance”. Numbers with underline mean better performance among all variants of ITEM with the same  $m$ .

aggregated them to get the final utility judgments. To ensure a fair comparison, we used 5 generated results and aggregated them in a way similar to k-sampling to obtain the final utility judgments. The Least-to-Most approach is designed for complex reasoning, but the utility judgment task doesn’t require problem decomposition. Llama 3 8B’s results using the listwise method on the TREC DL and WebAP datasets are shown in Table 14:

## C Case Study

Figure 9 presents two cases from the TREC DL dataset using Mistral under ITEM-A<sub>3</sub>. For the first question in Figure 9, the first pseudo-answer, though relatively correct, includes irrelevant information, leading to a misjudgment of “Passage-2” as “utility”. Based on the results of the first round of utility judgments, the second round of the pseudo-answer is more accurate and free of irrelevant content. Consequently, all three passages in the second round of utility judgments have utility in answering the question. For the second question in Figure 9, the first pseudo-answer is correct, but two passages without utility are judged as “utility”. The second pseudo-answer, with slight rewording, results in all selected passages being correct.

## D Datasets and Evaluation

Detailed statistics of the experimental datasets are shown in Table 15. We use the *trec\_eval* official tool for evaluation of NDCG. The annotation instruction for different points from Craswell et al. (2020):

- Perfectly relevant: The passage is dedicated to the query and contains the exact answer.
- Highly relevant: The passage has some answer for the query, but the answer may be a bit unclear, or hidden amongst extraneous information.

Table 14: Performance comparison of different methods on TREC DL and WebAP datasets using listwise approach (T in temperature sampling, k in top-k sampling, p in nucleus sampling).

Method	TREC DL (F1)	WebAP (F1)
UJ-ExpA (T=0)	52.83	26.99
5-sampling (T=0)	52.68	28.97
Self-consistency (T=0.7, k=40)	54.17	26.11
Self-consistency (T=0.5, k=40)	52.71	27.08
Self-consistency (T=0.3, k=40)	53.76	25.79
Self-consistency (T=0.7, k=20)	54.39	26.24
Self-consistency (T=0.7, k=40, p=0.95)	53.31	27.08
Self-consistency (T=0.7, k=40, p=0.90)	54.32	26.67
ITEM-A <sub>s</sub> (T=0)	<b>56.03</b>	<b>29.32</b>

<p><b>Question:</b> what type of tissue are bronchioles.</p>	<p><b>First answer:</b> The bronchioles are tubes made of epithelial tissue and smooth muscle tissue. They are...</p>	<p><b>Second answer:</b> The bronchioles are lined with epithelial tissue. Specifically, they are lined with ciliated cuboidal epithelium.</p>	<p><b>Third answer:</b> Simple cuboidal epithelium.</p>
	<p><b>First utility judgment:</b></p> <p>[2] What type of muscles control the size of the bronchioles in the lungs? The type of muscles...</p> <p>[3] Simple cuboidal epithelium is a type of epithelial tissue found in the body. It is a single layer thick and made of...</p> <p>[15] A pulmonary lobule is the portion of the lung ventilated by one bronchiole...</p>	<p><b>Second utility judgment:</b></p> <p>[3] Simple cuboidal epithelium is a type of epithelial tissue found in the body. It is a single layer thick and made of...</p> <p>[12] Simple cuboidal epithelium is a type of tissue that is found lining parts of organs...</p> <p>[13] Cuboidal epithelium is a type of tissue in animal anatomy. It is a very important...</p>	<p><b>Third utility judgment:</b></p> <p>[3] Simple cuboidal epithelium is a type of epithelial tissue found in the body. It is a single layer thick and made of...</p> <p>[12] Simple cuboidal epithelium is a type of tissue that is found lining parts of organs...</p> <p>[13] Cuboidal epithelium is a type of tissue in animal anatomy. It is a very important...</p>
	<p><b>First answer:</b> Family Feud has been on air since 1976.</p>	<p><b>Second answer:</b> The original Family Feud debuted in 1976.</p>	<p><b>Third answer:</b> The Family Feud debuted in 1976.</p>
<p><b>Question:</b> when did family feud come out?</p>	<p><b>First utility judgment:</b></p> <p>[1] What time is Celebrity Family Feud on tonight? Tonight's episode of Celebrity Family Feud ...</p> <p>[3] When will Celebrity Family Feud Season 3 start? Looking for the premiere date for Celebrity Family Feud...</p> <p>[13] ..Family Feud...The original version began in 1976...</p>	<p><b>Second utility judgment:</b></p> <p>[13] ..Family Feud...The original version began in 1976...</p>	<p><b>Third utility judgment:</b></p> <p>[13] ..Family Feud...The original version began in 1976...</p>

Figure 9: An example of our ITEM-A<sub>s</sub> using Mistral on the TREC dataset. Green means the passage has utility, and orange means the passage does not have utility.

- Related: The passage seems related to the query but does not answer it.
- Irrelevant: The passage has nothing to do with the query.

## E Answer Passage Retrieval

Non-factoid questions usually expected longer answers, such as sentence-level or passages-level (Keikha et al., 2014a; Yang et al., 2016; Keikha et al., 2014b). Yang et al. (2016) developed an annotated dataset for answer passage retrieval called WebAP, which has an average of 76.4 qrels per query. Cohen et al. (2018) and Hashemi et al. (2020) introduced the WikiPassageQA dataset and ANTIQUE dataset for answer passage retrieval, respectively. Compared to the WebAP dataset, WikiPassageQA and ANTIQUE have incomplete annotations, with an average of 1.7 qrels and 32.9 qrels per query (Hashemi et al., 2019,

Dataset	#Psg	#PsgLen	#Q	#Rels/Q
TREC	8.8M	93	82	212.8
WebAP	379k	45	73	76.4
NQ	21M	100	1868	1.0
GTI-NQ	10	100	1863	1.0

Table 15: Statistics of experimental datasets.

2020). Bi et al. (2019) created the PsgRobust dataset for answer passage retrieval, which is built on the TREC Robust collection (Voorhees et al., 2003) following a similar approach to WebAP but without manual annotation.

## F $k$ -sampling

The output of  $k$ -sampling each time contains explicit answers and utility judgments. If the question length is  $l_q$ , the total length of the input passages is  $l_p$ , and the average length of a single passage is  $l_{avg}$ , then the  $k$ -sampling input cost is  $(k + 1) \times (l_q + l_p)$ . If the average length of the output explicit answer is  $l_a$ , and the average length of the output utility judgments is  $l_u$ , then the  $k$ -sampling output cost is  $(k + 1) \times (l_a + l_u)$ . Taking ITEM-As as an example, with a maximum of three iterations, the maximum input cost for utility judgments is  $3 \times (l_q + l_p)$ . For answer generation, the longest input is  $l_q + l_p$  and the shortest is  $l_q + l_{avg}$ . Therefore, the maximum input cost for ITEM-As is  $6 \times (l_q + l_p)$  and the minimum is  $4 \times (l_q + l_p) + 2 \times (l_q + l_{avg})$ . The maximum output cost is  $3 \times (l_a + l_u)$ . Therefore, in order to ensure fairness in the calculation of large language model parameters, we choose  $k=5$ .

## G Retrievers

We use two representative retrievers to gather candidate passages in  $\mathcal{D}$  for utility judgments. Following with previous works (Zhang et al., 2024; Sun et al., 2023), we employ RocketQAv2 (Ren et al., 2021) and BM25 (Robertson et al., 2009) for the NQ dataset and retrieval datasets(i.e., TREC DL and WebAP datasets), respectively. Based on the retrieval results to build the  $\mathcal{D}$  we have two settings: (i) For TREC DL and WebAP datasets, we select the top-20 BM25 retrieval results. If these do not include passages with utility, we replaced the last one with a utility-annotated passage. (ii) For the NQ dataset, we use the top-10 dense retrieval results to form the candidate list  $\mathcal{D}$ , following the GTU setting of Zhang et al. (2024).

## H Instruction Details

### H.1 Instruction of Listwise and Pointwise Approaches

For the prompts of the NQ dataset using ChatGPT, we follow the setting of Zhang et al. (2024), otherwise, we use the following prompts. Following Sun et al. (2023), we input N passages using the form of multiple rounds of dialogue in the listwise approach. The prompts we used in our experiments are shown in Figure 10 and Figure 11.

### H.2 Instruction of the Ranking Approach

For RankGPT, we directly use the instruction of Sun et al. (2023) for relevance ranking, as shown in Figure 13. For the relevance ranking in our ITEM, the instructions are shown in Figure 12 and Figure ??.

### H.3 Instruction of Answer Generation

Li et al. (2023) utilize LLM to generate the missing information in the provided documents for the current question and then re-retrieve it as relevant feedback. Therefore, we have also designed two kinds of pseudo answers for utility judgments, i.e., (i) the explicit answer, which produces an answer based on the given information, and (ii) the implicit answer, which does not answer the question directly but gives the information necessary to answer the question. The two instructions are shown in Figure 14 and Figure 15.

**user:**  
You are the utility judge, an intelligent assistant that can select the passages that have utility in answering the question.

**assistant:**  
Yes, i am the utility judge.

**user:**  
I will provide you with {num} passages, each indicated by number identifier [].  
I will also provide you with a reference answer to the question.  
Select the passages that have utility in generating the reference answer to the following question from the {num} passages: {query}.

**assistant:**  
Okay, please provide the passages and the reference answer.

**user:**  
[1] {{passage\_1}}

**assistant:**  
Received passage [1]

**user:**  
[2] {{passage\_2}}

**assistant:**  
Received passage [2]  
(more passages) ...

**user:**  
Question: {query}.  
Reference answer: {answer}.

The requirements for judging whether a passage has utility in answering the question are: The passage has utility in answering the question, meaning that the passage not only be relevant to the question, but also be useful in generating a correct, reasonable and perfect answer to the question. Directly output the passages you selected that have utility in generating the reference answer to the question. The format of the output is: 'My selection:[i],[j],...'. Only response the selection results, do not say any word or explain.

Figure 10: Instruction in the listwise approach.

**user:**  
You are the utility judge, an intelligent assistant that can judge whether a passage has utility in answering the question or not.

**assistant:**  
Yes, i am the utility judge.

**user:**  
I will provide you with a passage and the reference answer to the question. \n Judge whether the passage has utility in generating the reference answer to the following question or not: {query}.

**assistant :**  
Okay, please provide the passage and the reference answer to the question.

**user:**  
Question: {query}.  
Reference answer: {answer}.  
Passage: {passage}

The requirements for judging whether a passage has utility in answering the question are: The passage has utility in answering the question, meaning that the passage not only be relevant to the question, but also be useful in generating a correct, reasonable and perfect answer to the question.

The reference answer may not be the correct answer, but it provides a pattern of the correct answer. Directly output whether the passage has utility in generating the reference answer to the question or not. If the passage has utility in generating the reference answer, output 'My judgment: Yes, the passage has utility in answering the question.'; otherwise, output 'My judgment: No, the passage has no utility in answering the question.'.

Figure 11: Instruction in the pointwise approach.

**user:**  
You are RankGPT, an intelligent assistant that can rank passages based on their relevance to the query.

**assistant:**  
Yes, i am RankGPT.

**user:**  
I will provide you with {num} passages, each indicated by number identifier []. I will also give you a reference answer to the query. Rank the passages based on their relevance to query: {query}.

**assistant :**  
Okay, please provide the passages and the reference answer.

**user:**  
[1] {{passage\_1}}

**assistant :**  
Received passage [1]

**user:**  
[1] {{passage\_2}}

**assistant :**  
Received passage [2]

(more passages) ...

**user:**  
Query: {query}.  
Reference answer: {answer}

Rank the {num} passages above based on their relevance to the query. The passages should be listed in descending order using identifiers. The most relevant passages should be listed first. The output format should be [] > [] > [] > ..., e.g., [i] > [j] > [k] > ... Only response the ranking results, do not say any word or explain.

Figure 12: Instruction of the relevance ranking approach in our ITEM.

**user:**  
You are RankGPT, an intelligent assistant that can rank passages based on their relevance to the query.

**assistant:**  
Yes, i am RankGPT.

**user:**  
I will provide you with {num} passages, each indicated by number identifier []. Rank the passages based on their relevance to query: {query}.

**assistant :**  
Okay, please provide the passages.

**user:**  
[1] {{passage\_1}}

**assistant :**  
Received passage [1]

**user:**  
[1] {{passage\_2}}

**assistant :**  
Received passage [2]

(more passages) ...

**user:**  
Query: {query}.

Rank the {num} passages above based on their relevance to the query. The passages should be listed in descending order using identifiers. The most relevant passages should be listed first. The output format should be [] > [] > [] > ..., e.g., [i] > [j] > [k] > ... Only response the ranking results, do not say any word or explain.

Figure 13: Instruction of the ranking approach in Sun et al. (2023).

**user:**  
You are a faithful question and answer assistant. Answer the question based on the given information with one or few words/sentences without the source.

**assistant:**  
Yes, i am the faithful question and answer assistant.

**user:**  
Given the information:  
{passage}  
Answer the following question based on the given information with one or few words/sentences without the source.  
Question: {question}  
Answer:

Figure 14: Instruction of the explicit answer generation.

**user:**  
You are a faithful question and answer assistant. Given a question and references. To answer the question, output which information is necessary to answer the question based on the references.

**assistant:**  
Yes, i am the faithful question and answer assistant.

**user:**  
References: {pas}  
Question: {question}  
To answer the question, output which information is necessary to answer the question based on the references. Do not mention references when printing out necessary information. The format of the output is: 'Necessary information: [xxx]'.

Figure 15: Instruction of the implicit answer generation.