# Rethinking LLM Parametric Knowledge as Post-retrieval Confidence for Dynamic Retrieval and Reranking

Haoxiang Jin[†]
School of Computer Science and Technology, Xidian University
Shaanxi, China
jinhx@stu.xidian.edu.cn

Ronghan Li[†]
School of Computer Science and Technology, Xidian University
Shaanxi, China
lironghan@stu.xidian.edu.cn

Zixiang Lu[*]
School of Computer Science and Technology, Xidian University
Shaanxi, China
zxlu@xidian.edu.cn

Qiguang Miao[*]
School of Computer Science and Technology, Xidian University
Shaanxi, China
qgmiao@xidian.edu.cn

## Abstract

Large Language Models (LLMs) often generate inaccurate responses (hallucinations) when faced with questions beyond their knowledge scope. Retrieval-Augmented Generation (RAG) addresses this by leveraging external knowledge, but a critical challenge remains: determining whether retrieved contexts effectively enhance the model's ability to answer specific queries. This challenge underscores the importance of knowledge boundary awareness, which current methods-relying on discrete labels or limited signals-fail to address adequately, as they overlook the rich information in LLMs' continuous internal hidden states. To tackle this, we propose a novel post-retrieval knowledge filtering approach. First, we construct a confidence detection model based on LLMs' internal hidden states to quantify how retrieved contexts enhance the model's confidence. Using this model, we build a preference dataset (NQ_Rerank) to fine-tune a reranker, enabling it to prioritize contexts preferred by the downstream LLM during reranking. Additionally, we introduce Confidence-Based Dynamic Retrieval (CBDR), which adaptively triggers retrieval based on the LLM's initial confidence in the original question, reducing knowledge conflicts and improving efficiency. Experimental results demonstrate significant improvements in accuracy for context screening and end-to-end RAG performance, along with a notable reduction in retrieval costs while maintaining competitive accuracy.

## CCS Concepts

• **Do Not Use This Code** → **Generate the Correct Terms for Your Paper**; *Generate the Correct Terms for Your Paper*; Generate

the Correct Terms for Your Paper; Generate the Correct Terms for Your Paper.

## Keywords

Knowledge Boundary, Evaluation, Large Language Models, Retrieval-Augmented Generation, Reranker, Generator

## 1 Introduction

Large Language Models (LLMs) have demonstrated exceptional performance in diverse text generation tasks, such as creative writing and text summarization [1–3]. However, when confronted with questions beyond their knowledge scope, they often generate plausible yet inaccurate responses—a phenomenon termed "hallucination" [4, 5]. To address questions unanswerable by model parameters alone, Retrieval-Augmented Generation (RAG) [6, 7] leverages external knowledge sources to expand the answerable question boundary . Yet this approach introduces a critical challenge : After retrieval, how can we precisely determine whether the acquired knowledge genuinely enhances the model's ability to answer a specific query [8–10]? The core of this challenge lies in effectively coordinating the model's internal parametric knowledge with retrieved contexts to delineate the knowledge utilization scope of RAG system in post-retrieval stages.

This challenge underscores the importance of knowledge boundary awareness. Failures in answering questions typically stem from two causes: 1) Suboptimal prompt design, where inadequate prompting fails to unlock the model's potential [11–13]. This can often be mitigated through prompt optimization strategies such as chain-of-thought [14] or self-verification prompts[15] . 2) Fundamental unawareness of knowledge boundary, where the model cannot recognize its limitations and thus fails to abstain from answering beyond its competence [13, 16–18]. This issue persists in RAG system: poor boundary awareness impedes judgment on the quality of

**(a) Comparison of Reranking Strategies**

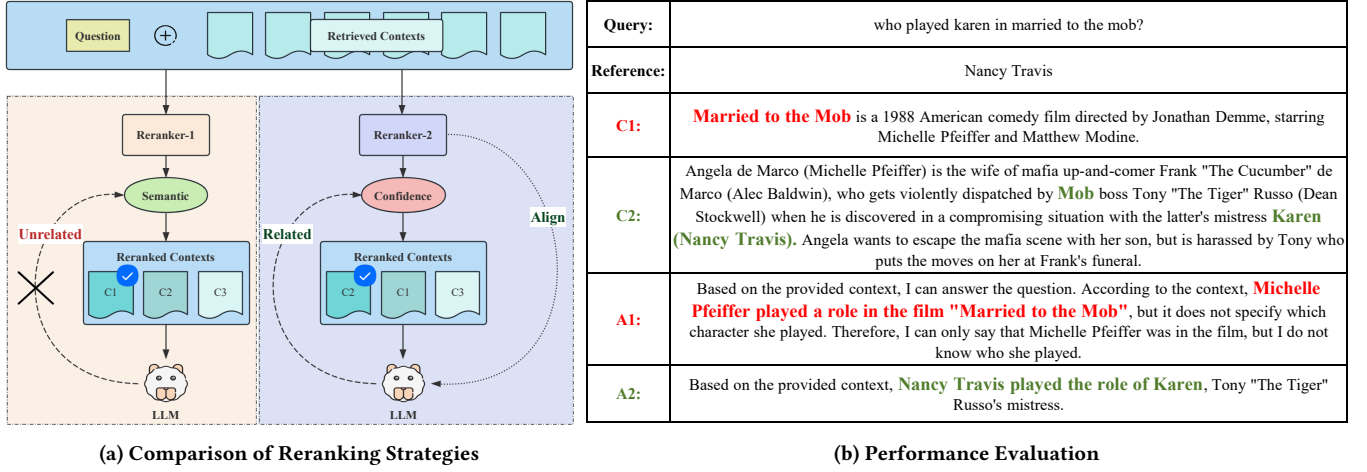| Query: | who played karen in married to the mob? |
|---|---|
| Reference: | Nancy Travis |
| **C1:** | **Married to the Mob** is a 1988 American comedy film directed by Jonathan Demme, starring Michelle Pfeiffer and Matthew Modine. |
| **C2:** | Angela de Marco (Michelle Pfeiffer) is the wife of mafia up-and-comer Frank "The Cucumber" de Marco (Alec Baldwin), who gets violently dispatched by **Mob** boss Tony "The Tiger" Russo (Dean Stockwell) when he is discovered in a compromising situation with the latter's mistress **Karen (Nancy Travis).** Angela wants to escape the mafia scene with her son, but is harassed by Tony who puts the moves on her at Frank's funeral. |
| **A1:** | Based on the provided context, I can answer the question. According to the context, **Michelle Pfeiffer played a role in the film "Married to the Mob"**, but it does not specify which character she played. Therefore, I can only say that Michelle Pfeiffer was in the film, but I do not know who she played. |
| **A2:** | Based on the provided context, **Nancy Travis played the role of Karen**, Tony "The Tiger" Russo's mistress. |

**(b) Performance Evaluation**

**Figure 1: The left part of the diagram (Figure 1a) contrasts two distinct reranking strategies for RAG systems. One employs a conventional similarity-based reranker, which prioritizes contexts solely through textual matching between questions and documents. The other leverages the LLM's intrinsic preference, reranking contexts according to their ability to enhance the model's answer confidence. The right part (Figure 1b) provides a concrete example comparing the effectiveness of contexts reranked by Reranker-1 (similarity-based) and Reranker-2 (LLM-aligned). It demonstrates the responses generated by the same LLM when using each set of contexts, highlighting differences in answer quality, confidence, and relevance.**

retrieved contexts, which in turn makes it difficult to select useful ones that could enhance their reasoning on a given question [19]. Enhancing such awareness is therefore crucial for improving RAG system accuracy.

Current research on knowledge boundary perception follows three primary paths: 1) Prompt-guided confidence estimation: Using engineered prompts to elicit self-assessed confidence scores [20–22]. 2) Multi-sample confidence aggregation: Estimating confidence via correctness rates across multiple responses to the same question [23, 24]. 3) Hidden-state-based confidence: Quantifying uncertainty through the model's internal internal hidden states such as intermediate layer activations [18, 25, 26]. Most existing methods rely on discrete labels such as answerable and unanswerable, lacking analysis of how internal hidden states evolve before and after introducing external knowledge. Crucially, model's internal hidden states, as continuous vector representations, encapsulate richer information than discrete tokens and may more faithfully reflect model confidence [18, 25, 26].

We posit that LLM's internal hidden state reflected confidence serves as a key indicator for evaluating whether retrieved contexts effectively enhances their question-answering capability. Furthermore, the confidence shift observed when LLM processes different retrieved contexts for the same question inherently reveals its preference among retrieved contexts. This intrinsic preference signal can guide reranker model to filter and optimize contexts post-retrieval, significantly boosting RAG system efficacy. As illustrated in Figure 1a, the left side employs an unrelated Reranker + LLM to form a RAG system, which represents the most common current approach. On the right, building upon the left configuration, the Reranker aligns with the LLM's intrinsic preference by ranking contexts that most enhance the LLM's confidence in answering the question during the reranking process. Figure 1b compares the

helpfulness of C2 and C1 in assisting the same LLM in answering the question.

Inspired by work [18], this paper proposes a novel method for post-retrieval knowledge filtering. First, we construct a confidence detection model based on the LLM's internal hidden states to quantify how much the retrieved contexts enhances the LLM's confidence. Leveraging this confidence detection model's analysis, we build a preference dataset, NQ_Rerank, which is then used to fine-tune a Reranker model. This fine-tuning enables the Reranker to prioritize contexts preferred by the downstream LLM during the reranking phase, thereby improving the accuracy of the RAG system in answer generation. Additionally, we introduce Confidence-Based Dynamic Retrieval(CBDR) that adaptively triggers the retrieval process based on the LLM's initial confidence in the original question. This mechanism reduces the risk of knowledge conflicts while enhancing the RAG system's overall efficiency .

Experiments show our approach achieves: 1) 5.19% improvement in post-retrieval contexts screening accuracy. 2) 4.70% higher end-to-end RAG system accuracy. 3) 7.10% reduction in retrieval costs (with dynamic retrieval enabled) while maintaining 5.60% accuracy gains. This work establishes a quantifiable framework for identifying and extending knowledge boundary in RAG system and introduces a confidence-shift-based metric for evaluating retrieval augmentation effectiveness.

## 2 Related Work

### 2.1 Knowledge Boundary of LLM

Since the seminal study [13] introduced the concept of knowledge boundaries to LLM research, it has become a cornerstone for evaluating model's self-awareness. This work categorizes parametric

knowledge into three types: **Prompt-agnostic Knowledge**, Correctly answerable regardless of query phrasing. **Prompt-sensitive Knowledge**, Answerable only under specific prompting strategies. **Unanswerable Knowledge**, Incapable of correct response under any prompt.

Current evaluation paradigms focus on measuring confidence levels regarding answerability:

**Expression-based Confidence:** Leverages LLM's instruction-following capability to elicit self-reported confidence via prompts [15].

**Sampling-based Confidence Estimation:** Uses multi-round sampling such as query paraphrasing or output variations to compute answer entropy and confidence scores [27].

**Internal State-based Confidence Estimation:** Utilizes hidden states—particularly intermediate-layer activations when generating the first token—as continuous confidence indicators [18].

Our work adopts the third approach [18], using Mid_Layer hidden states at the first response token generation as the confidence metric.

## 2.2 Knowledge Boundary in RAG System

The knowledge boundary in RAG system is typically defined as the knowledge space collectively formed by the LLM's internal parametric knowledge and external retrieved knowledge. However, early evaluations of RAG system capabilities predominantly focused on Retriever performance, overemphasizing the system's reliance on external knowledge. As revealed by studies such as [8–10], conflicts between externally retrieved knowledge and internal parametric knowledge may lead the model to produce low-confidence errors.

Consequently, subsequent research has shifted toward coordinating these dual knowledge sources and more precisely delineating RAG's effective knowledge boundary. Approaches like [28] and [29] analyze the internal states of LLM to detect model uncertainty, dynamically determining whether to activate the retrieval process based on this characteristic. Most recently, the DTA [19] framework formally proposes the concept of the Knowledge Boundary of RAG. This study categorizes potential queries into four quadrants based on the LLM's inherent Parametric Knowledge Boundary $KB_p$ and the Retrieved Knowledge Boundary $KB_r$ provided by the Retriever, collectively delineating the holistic effective knowledge boundary of the RAG system.

## 2.3 Preference Alignment in RAG system

To enhance the efficiency of LLM in utilizing externally retrieved knowledge, it is essential to align the preferences between the Retriever and the LLMs within RAG system. Diverse studies have proposed distinct preference signals to guide this alignment: RE-PLUG [30] employs the probability of an LLM generating a correct answer as a preference signal to identify critical contexts; RRR [31] utilizes the overall quality of LLM-generated responses as a preference metric; DPA-RAG [32] introduces a bidirectional alignment strategy to mitigate preference conflicts among different RAG components; RADIO [33] constructs preference indicators based on the correctness of generated rationales, subsequently fine-tuning the reranker model to reconcile preference discrepancies between

Retrievers and LLMs; SEAKR [29] performs multi-round sampling for the same query and leverages the last-layer hidden states of LLMs at the end-of-sequence token (</s>) to compute a Gram matrix. This matrix quantifies model uncertainty, which serves as a preference signal for Reranker optimization.

The core innovation of this paper lies in proposing a novel preference metric: confidence shift, characterized by changes in the internal hidden states of LLM before and after exposure to different external knowledge. This metric is utilized to fine-tune the Reranker, enabling effective filtering of post-retrieval contexts. Compared to methods like SEAKR [29] which similarly leverage internal hidden states but require multi-round sampling, but our confidence shift detection mechanism, based on a single forward pass, significantly reduces computational and temporal overhead. This efficiency is particularly advantageous for real-time application scenarios demanding low latency.

## 3 Methods

In this section, we provide a detailed exposition of the technical methodologies employed to assess the confidence of LLM in their responses to questions by leveraging their internal hidden states. We elaborate on the construction of a preference dataset based on variations in these internal hidden states and demonstrate how this dataset can be utilized to fine-tune a Reranker. Finally, we propose Confidence-Based Dynamic Retrieval(CBDR) that leverages the LLM's self-confidence estimates to optimize information retrieval processes.

## 3.1 Internal State Detection

Human cognitive processes demonstrate that not all thought and reasoning rely on linguistic mediation; the explicit articulation of thought into language often entails a loss of information. This characteristic shares similarities with the operational mechanisms of modern LLM. LLM typically process input information through internal hidden layers and map complex latent representations into sequences of lexical tokens during the final output stage. This transformation, constrained by a fixed vocabulary, inevitably leads to partial information loss.

Recent studies, such as work [34] and work [35], reveal that the internal hidden states of LLMs contain richer information than their final outputs and exhibit stronger latent reasoning capabilities. Meanwhile, [36] demonstrates that internal hidden states in internal hidden layers (specifically at Mid_Layer which is Layer/2) effectively capture the model's self-awareness. Furthermore, [18] indicates that LLMs can perceive their own knowledge boundaries prior to generating responses (Pre-Token).

*3.1.1 Confidence Estimation via Internal Hidden States.* Specifically, the workflow for self-confidence detection based on the internal hidden states of LLM is as follows: For a given target LLM M and a question Q, the model generates internal hidden state representations during inference, denoted as $H_{M,Q}$. Compared to the final token output, this state encapsulates more comprehensive information. Our confidence estimation process is defined as:
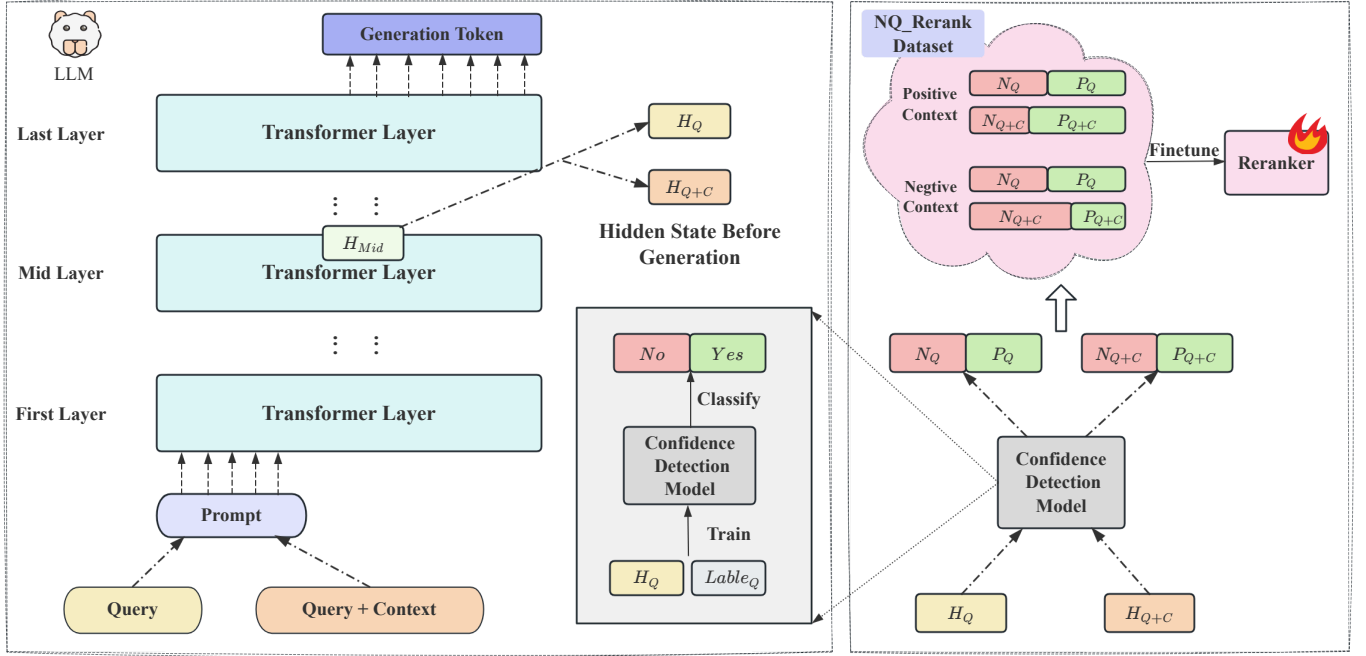
$$C_{M,Q} = E(H_{M,Q}) \tag{1}$$

Figure 2: The complete process of aligning the Reranker with the target LLM involves constructing a preference dataset, NQ_Rerank, by comparing the variations in the LLM's confidence when answering a Question under different contexts. This dataset is then used to fine-tune the Reranker model, aligning it with the target LLM's intrinsic preferences.
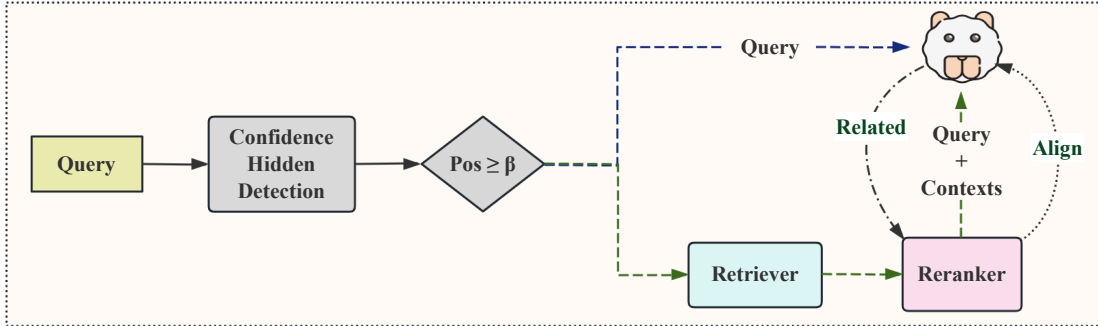


Figure 3: The workflow of Confidence-Based Dynamic Retrieval (CBDR). By varying the confidence threshold $\beta$, we balance the accuracy and retrieval cost of the Retrieval-Augmented Generation (RAG) system.

As illustrated in left side of Figure 2, where E denotes the confidence detection model, and $C_{M,Q}$ is a binary classification label: $C_{M,Q} = 1$ indicates that LLM M is confident in correctly answering question q, whereas $C_{M,Q} = 0$ signifies that the LLM M perceives itself as incapable of responding accurately. Drawing on work [18] and related prior work, we select the internal hidden state vector at Mid_Layer (Layer/2) of LLM M before generating the first answer token (Pre-Token) as $H_{M,Q}$.

The training data for model E is obtained by guiding LLM M to process questions from the NQ dataset [37]. We collect the internal hidden state $H_{M,Q}$ during inference and determine the correctness of the LLM M's response based on the ground-truth answer to question Q, thereby constructing binary training samples ($H_{M,Q}$,

$Label_Q$). Here, $Label_Q = 1$ indicates that the model answers question Q correctly, while $Label_Q = 0$ denotes an incorrect response. The training methodology for model E follows the approach described in work [18]. We performed data cleaning on the training dataset NQ [37] and analyzed the impact of different prompt designs on model reasoning.

## 3.2 Preference Dataset

*3.2.1 Preference Definition.* This study focuses on the post-retrieval processing stage within RAG system, with the aim of exploring how to rerank the retrieved contexts to maximize RAG system's utility in enhancing the answer reasoning capabilities of downstream LLM.

Conventional Reranker are typically trained on datasets constructed based on semantic similarity between a question and contexts, and compute relevance scores by capturing complex semantic interactions through interactive encoding. While such general-purpose methods ensure model transferability and compatibility with diverse LLMs, they often fail to adequately incorporate the preferences of specific downstream LLM, thereby limiting the full potential of RAG system.

$$
\begin{aligned}
Conf(H_{M,Q}) &= P(Label = 1 \mid E(H_{M,Q})) \\
&= Softmax(Z_1) \\
&= \frac{e^{z_1}}{e^{z_0} + e^{z_1}}
\end{aligned}
\tag{2}
$$

As illustrated in Figure 2, this paper defines the following preference criterion: a context $C$ is considered to exhibit a positive preference for the target LLM $M$ in answering question $Q$ if and only if it provides effective informational enhancement, satisfying the condition $Conf(H_{M,Q+C}) > Conf(H_{M,Q})$. Conversely, if it leads to a decrease in LLM's confidence $Conf(H_{M,Q+C}) < Conf(H_{M,Q})$, the context $C$ is regarded as having a negative preference. As shown in Equation 2, the output of the Conf(·) function is defined as the probability of the $Label = 1$ assigned by model E. A softmax layer is appended to the final layer of model E to produce this probabilistic output.

*3.2.2 Dataset Construction.* We preprocess the NQ dataset [37] to obtain a series of $(Query, Contexts)$ tuple samples. For each sample, we record the internal hidden state at Mid_Layer when the target LLM M generates its first token under the following two scenarios: 1) The state $H_{M,Q}$ when only the query $Q$ is provided; 2) The state $H_{M,Q+C_i}$ when both the query $Q$ and a context $C_i$ are provided (Where i iterates over the Contexts).

This yields a sequence of internal hidden states:

$$[H_{M,Q}, H_{M,Q+C_1}, H_{M,Q+C_2}, ..., H_{M,Q+C_i}]$$

This sequence of states is then fed into the confidence hidden detection model E to obtain the probability value for the $Label = 1$ output by the softmax layer, resulting in a probability sequence:

$$[Conf(H_{M,Q}), Conf(H_{M,Q+C_1}), ..., Conf(H_{M,Q+C_i})]$$

The enhancement effect of each context $C_i$ on LLM M's response to question $Q$ is determined by comparing the change in model confidence after incorporating the context $C_i$:

$$
Inc(Q, C_i) = Conf(H_{M,Q+C_i}) - Conf(H_{M,Q})
\tag{3}
$$

If $Inc(Q, C_i) > 0$, the sample is labeled as a positive preference sample. If $Inc(Q, C_i) < 0$, it is labeled as a negative preference sample.

For each $(Query, Contexts)$ sample, all context $C_i$ are ranked according to $Inc(Q, C_i)$. The Top-K(K = 5) contexts with the highest increase are selected as positive examples, and the Top-K with the largest decrease are taken as negative examples. As illustrated in right side of Figure 2, this process constructs the final preference dataset, denoted as NQ_Rerank.

## 3.3 Reranker Fine-tuning

To enhance the ability of the Reranker to identify the utility of contexts for the target LLM, we performed supervised fine-tuning on a base Reranker using the constructed preference dataset NQ_Rerank. This fine-tuning process aims to achieve effective alignment between the Reranker and the target LLM's preferences, enabling it to more accurately select contexts that significantly enhance the target LLM's confidence in answering a given question $Q$.

During fine-tuning, the InfoNCE (Noise Contrastive Estimation) loss function was employed as the optimization objective:

$$
f(Q, C) = exp(\phi(Q, C)/\tau)
\tag{4}
$$

$$
L = -log\frac{f(Q, C^+)}{f(Q, C^+) + \sum_{i=1}^{N} f(Q, C_i^-)}
\tag{5}
$$

Where: $f(Q, C)$ denotes the relevance score between question $Q$ and context $C$ computed by the Reranker; $C^+$ represents the positive context; $C^-$ denotes the negative context; $\tau$ is the temperature parameter.

This loss function forces the model to increase the score margin between the positive context $C^+$ and a set of negative contexts $\{C^-\}$, thereby learning a ranking criterion consistent with the target LLM's preferences.

## 3.4 Confidence-Based Dynamic Retrieval

Although the fine-tuned Reranker has achieved considerable alignment with the target LLM's preferences and effectively improved the ranking priority of beneficial contexts, it still exhibits the following limitations:

**Risk of Irrelevant Context Interference:** The context set returned by the Retriever may contain few or no beneficial contexts. Under such circumstances, even after reranking, the Top-K results output by the Reranker may still include irrelevant or misleading context, which can interfere with the reasoning process of the downstream LLM and even lead to hallucinations.

**Redundant Computational Overhead:** Typical Reranker often contain hundreds of millions to billions of parameters. Although single inference is relatively efficient, the cumulative computational cost of executing a full retrieval-reranking process for every query can be substantial, especially when the retrieved contexts have low relevance, resulting in significantly diminished cost-effectiveness.

To mitigate these issues and enhance the efficiency and reliability of the RAG system, we propose CBDR. The workflow of this strategy is illustrated in Figure 3: 1) If the target LLM exhibits high confidence in responding to the current query $Q$ that $Conf(H_{M,Q}) > \beta$, where $\beta$ is a predefined threshold, the retrieval and reranking steps are skipped, and the LLM generates the answer directly. 2) If the confidence score falls below the threshold $Conf(H_{M,Q}) < \beta$, the full retrieval process is initiated: the Retriever fetches a set of contexts, which are reranked by the fine-tuned Reranker, and the Top-K contexts are fed into the LLM along with the query for reasoning.

This strategy aims to maintain answer quality whenever possible while significantly reducing redundant computation for high-confidence questions and avoiding potential interference from low-quality retrieval results for known questions. The effectiveness of this strategy will be thoroughly validated in Section 4.2.3 of the experiments.

## 4 Experiments

This section aims to systematically evaluate the effectiveness of the proposed method, focusing on the following three aspects: 1) We assess the impact of preference-aligned fine-tuning on the performance of the Reranker by examining the improvement in ranking tasks after fine-tuning with the constructed dataset NQ_Rerank. 2) We evaluate the influence of the preference-aligned Reranker on the RAG system, analyzing whether the internal hidden states of the LLM can reliably reflect its preference for retrieved contexts. 3) We examine the effect of the CBDR on the efficiency and accuracy of the RAG system, verifying the performance and advantages of the confidence-based retrieval scheduling strategy in practical applications.

It should be noted that, since this study employs the Reranker to align with the preferences of downstream LLM and designs experiments accordingly to evaluate the improvements, the Retriever module has been intentionally excluded from the experimental setup. Instead, all Rerankers are provided with the same set of retrieved contexts to ensure a fair evaluation of their reranking performance under consistent contextual inputs.

## 4.1 Experimental Setup

*4.1.1 Datasets.* We take two representative open-domain QA benchmark datasets, including the Natural Questions (NQ) dataset [37] and the HotpotQA dataset [38]. The NQ dataset is constructed from real Google search questions and contains retrieved relevant contexts along with human-annotated short and long answers. HotpotQA is a question-answering dataset comprising examples that require multi-step reasoning.

All data used for training in this study were derived from the NQ dataset. The specific data partitions are described below:

**Dataset for Hidden Detection Model E:** The training set (Train) consists of 1,000 positive samples (questions answered correctly by the target LLM M) and 1,000 negative samples (questions answered incorrectly by the target LLM M), selected from the NQ dataset; the development set (Dev) contains 300 positive and 300 negative samples and was used for hyperparameter tuning and early stopping; the test set (Test) includes 500 positive and 500 negative samples for the final performance evaluation of the mode.

**The preference dataset NQ_Rerank:** Constructed based on the NQ_Retrieval[*] dataset—a structured version of the NQ dataset adapted for retrieval tasks. The final version of the dataset excluded data items lacking valid positive-context or negative-context examples. The resulting dataset comprises 7,622 training samples and 1,216 evaluation samples.

---

*4.1.2 LLMs.* We conducted experiments on two representative open-source models, including Llama3-8B-Instruct [39] and Qwen2.5-7B-Instruct [40]. The experiments uniformly employed Llama3-8B-Instruct as the base LLM for downstream task reasoning. During inference, the temperature of the model was set to 1.0, and tokens were selected using a greedy decoding strategy.

*4.1.3 Rerankers.* For a comprehensive comparison, we selected four representative Rerankers as backbone models and conducted preference alignment fine-tuning experiments specifically on bge-reranker-v2-m3:

**gte_passage-ranking_multilingual-base:** A multilingual Reranker proposed by Alibaba DAMO Academy. Public evaluations indicate that it outperforms other models of similar scale in multilingual retrieval tasks.

**Qwen3-Reranker-4B:** A 4-billion-parameter version based on the Qwen foundation model, specifically optimized for text embedding and reranking tasks.

**Qwen3-Reranker-8B:** A 8-billion-parameter version based on the Qwen foundation model, Qwen3-Reranker-8B represents one of the state-of-the-art publicly available rerank models.

**bge-reranker-v2-m3:** A lightweight reranker known for its strong multilingual support and efficient inference speed.

**bge-reranker-v2-m3-ft(Ours):** To validate the effectiveness of aligning with the target LLM's preferences via confidence estimation, we conducted supervised fine-tuning on the bge-reranker-v2-m3 model using the constructed NQ_Rerank preference dataset.

*4.1.4 Evaluation Metrics.*

$$MRR@K = \frac{1}{|Q|} \sum_{i=1}^{|Q|} \frac{1}{rank_i} \tag{6}$$

Following mainstream evaluation practices for Rerankers , we adopt the following metrics to assess performance: 1) Precision@K measures the proportion of effectively positive contexts among the Top-K ranked results. 2) Recall@K evaluates the ratio of successfully ranked positive contexts within the Top-K results relative to all relevant contexts. 3) Mean Reciprocal Rank (MRR@K) represents the average reciprocal rank of the first relevant context across all questions. Where $rank_i$ denotes the position of the first relevant context for the i-th query (positions beyond K are excluded from calculation).

Given that the number of positive contexts or negative contexts per query item in the NQ_Rerank dataset ranges from [1, 5], we select Top-k values of k=1, 3, 5 for evaluation to align with the data characteristics.

*4.1.5 Implementation Details.* During the training of the internal hidden detection model E, the initial learning rate was set to $5e^{-5}$, the dropout rate was configured to 0.5, and the training was conducted over 30 epochs.

For the fine-tuning of the bge-reranker-v2-m3 model, the initial learning rate was set to $6e^{-5}$, weight decay was configured to 0.01, the maximum query length (query_max_len) was set to 128, the maximum passage length (passage_max_len) was set to 512, and the training was performed for 1 epoch.

**Table 1: Performance Comparison of Different Rerankers on the NQ_Rerank Test Set.**

| Reranker | Params | Top-1 | | | Top-3 | | | Top-5 | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | Precision | Recall | MRR | Precision | Recall | MRR | Precision | Recall | MRR |
| gte_passage-ranking_multilingual-base | 304M | 85.52 | 29.47 | 85.52 | 71.45 | 62.66 | 90.37 | 60.98 | 82.53 | 90.99 |
| Qwen3-Reranker | 4B | 81.74 | 27.62 | 81.74 | 70.92 | 62.33 | 88.15 | 61.71 | 83.53 | 88.93 |
| Qwen3-Reranker | 8B | <u>87.25</u> | <u>30.47</u> | <u>87.25</u> | <u>74.35</u> | <u>65.15</u> | <u>91.65</u> | <u>64.22</u> | <u>86.42</u> | <u>92.19</u> |
| bge-reranker-v2-m3 | 568M | 86.01 | 29.45 | 86.01 | 72.62 | 63.61 | 90.47 | 62.40 | 84.01 | 91.07 |
| bge-reranker-v2-m3-ft (Ours) | 568M | **91.20** | **32.01** | **91.20** | **76.98** | **67.14** | **94.40** | **65.64** | **87.97** | **94.72** |

## 4.2 Results

*4.2.1 Reranker Performance.* To evaluate whether the fine-tuned Reranker bge-reranker-v2-m3-ft can more effectively select retrieved contexts suitable for the downstream LLM Llama3-8B-Instruct, comparative experiments were conducted on the NQ_Rerank test set. The experimental setup was as follows: each Reranker was provided with a query and its corresponding set of contextual documents (including both positive contexts and negative contexts). Each Reranker then output a reranked list of Top-K relevant documents (where K ∈ 1, 3, 5). The performance of each Reranker was comprehensively evaluated using the Precision@K, Recall@K, and MRR@K metrics.

The experimental results, summarized in Table 1, are as follows:

1. The fine-tuned model bge-reranker-v2-m3-ft, which underwent preference alignment training using the NQ_Rerank dataset, achieved optimal performance across all evaluation metrics (K ∈ 1, 3, 5).
2. Notable improvements were observed particularly in Top-1 performance. Compared to the second-best model Qwen3 _Rerank_8B, Precision@1 and MRR@1 increased by +3.95 percentage points (pp) while Recall@1 improved by +1.54 pp; when compared to the baseline model bge-reranker-v2-m3 before fine-tuning, Precision@1 and MRR@1 increased by +5.19 pp and Recall@1 improved by +2.56 pp.

*4.2.2 RAG System Accuracy.* To further investigate whether a fine-tuned Reranker can enhance the final performance of a Retrieval-Augmented Generation (RAG) system, we constructed multiple "Reranker + LLM" combined systems for comparative experiments. The experimental setup was as follows: each Reranker performed re-ranking on the query and its corresponding set of context documents, selected the Top-K (K ∈ 1, 3) documents, and fed them to the downstream LLM for answer generation. The overall system performance was ultimately evaluated by assessing the precision of the generated answers.

The experimental results, as shown in Table 2, yielded the following main findings:

1. When the downstream LLM was Llama3-8B-Instruct, the RAG system utilizing the bge-reranker-v2-m3-ft as the Reranker

consistently achieved higher accuracy than the system using the original bge-reranker-v2-m3, with a maximum improvement of up to +4.7 pp. Furthermore, this combination achieved optimal or near-optimal performance on both the NQ and HotpotQA datasets.
2. When the downstream LLM was Qwen2.5-7B-Instruct, the RAG systems employing either the fine-tuned or the original Reranker demonstrated comparable performance in terms of accuracy, with no significant differences observed. This to some extent demonstrates the robustness of bge-reranker-v2-m3-ft.

*4.2.3 Dynamic Retrieval Efficiency.* We also evaluated the impact of a dynamic retrieval strategy based on the confidence level of the downstream LLM on the performance of the RAG system. The experiment measured the system's accuracy (Precision) and the proportion of overhead saved due to skipped retrieval under different confidence thresholds, as summarized in Table 3. The main findings are as follows:

1. Enabling the dynamic retrieval mechanism allowed the RAG system to significantly reduce retrieval overhead by assessing the LLM's confidence in addressing the query. Furthermore, under the setting with a reranking scope of Top-3 contexts, the system's accuracy improved by 0.9 percentage points.
2. In the Top-1 scenario, the strategy with dynamic retrieval enabled resulted in a slightly lower accuracy—by 0.2 percentage points—compared to the strategy without dynamic retrieval.

## 5 Discussion

Based on the experimental results presented, we draw the following conclusions:

## 5.1 Model Scale Positively Correlates with Reranker Performance

A positive correlation is observed between model scale and performance among the four un-tuned base Rerankers. Generally, performance improves as model size and capability increase. Taking the Qwen3_Reranker series as an example: although its absolute performance remains relatively low, Qwen3_Rerank_8B(8

**Table 2: Accuracy of RAG Systems with Different Reranker and LLM Combinations.**

| LLM | Reranker | Params | HotpotQA | | NQ | |
|---|---|---|---|---|---|---|
| | | | Top-1 | Top-3 | Top-1 | Top-3 |
| | gte_passage-ranking_multilingual-base | 304M | 47.20 | 51.80 | <u>63.80</u> | 67.60 |
| | Qwen3-Reranker | 4B | 42.30 | 50.10 | 50.70 | 64.00 |
| Qwen2.5-7B-Instruct | Qwen3-Reranker | 8B | <u>47.50</u> | <u>51.90</u> | 56.30 | 68.80 |
| | bge-reranker-v2-m3 | 568M | 47.20 | **53.30** | **64.20** | <u>69.70</u> |
| | bge-reranker-v2-m3-ft (Ours) | 568M | **48.70** | **53.30** | 63.40 | **69.90** |
| | gte_passage-ranking_multilingual-base | 304M | **48.80** | 50.20 | 60.10 | 60.70 |
| | Qwen3-Reranker | 4B | 40.70 | 48.00 | 49.70 | 62.30 |
| Llama3-8B-Instruct | Qwen3-Reranker | 8B | <u>48.40</u> | 50.10 | 55.20 | **68.80** |
| | bge-reranker-v2-m3 | 568M | 46.60 | <u>51.40</u> | <u>61.50</u> | 62.20 |
| | bge-reranker-v2-m3-ft (Ours) | 568M | 48.00 ↑ | **52.20** ↑ | **62.60** ↑ | <u>66.90</u> ↑ |

**Table 3: Efficiency-Accuracy Trade-off of the Confidence-Based Dynamic Retrieval Strategy.**

| Reranker | NQ | | | |
|---|---|---|---|---|
| | Top-1 | Top-3 | $\beta$ | RR↓ |
| bge-reranker-v2-m3 | 61.50 | 62.20 | 0 | 100 |
| bge-reranker-v2-m3-ft (Ours) | **62.60** | <u>66.90</u> | 0 | 100 |
| bge-reranker-v2-m3-ft (Ours) | <u>62.40</u> | 66.10 | 0.95 | **83.30** |
| bge-reranker-v2-m3-ft (Ours) | 61.70 | **67.80** | 0.98 | <u>92.90</u> |

to select contexts that increase the LLM's confidence in answering a given query. Significant improvements in Top-1 accuracy (Precision@1) and Mean Reciprocal Rank (MRR@1) further highlight the fine-tuned Reranker's enhanced capability in high-precision retrieval of critical documents.

### 5.3 Effectiveness of Preference Alignment Depends on the Downstream LLM

The effectiveness of preference alignment varies depending on the downstream LLM. Results from Table 2 show that when using the same fine-tuned Reranker bge-reranker-v2-m3-ft, accuracy improvements reached up to 4.7 pp with Llama3-8B-Instruct, whereas negligible gains were observed with Qwen2.5-7B-Instruct. This suggests that the success of preference alignment is partially dependent on the architecture and training methodology of the target LLM, and that the preferences embedded during fine-tuning must be compatible with the specific LLM being used.

billion parameters) significantly outperforms its lighter counterpart Qwen3_Rerank_4B(4 billion parameters) and achieves the best results among the base models. This trend demonstrates that 1) the preference dataset NQ_Rerank, constructed based on LLM confidence, possesses intrinsic validity in distinguishing model capabilities; and 2) the proposed confidence-change-based preference definition (Section 3.2.1) is well-founded.

### 5.2 Preference-Aligned Fine-Tuning Yields Substantial Gains

Preference-aligned fine-tuning leads to substantial performance gains. The fine-tuned Reranker bge-reranker-v2-m3-ft outperforms all baseline Rerankers across all evaluation metrics. This strongly indicates that supervised fine-tuning using the NQ_Rerank dataset effectively aligns the Reranker model with the preferences of the target downstream LLM (Llama3-8B-Instruct), making it more inclined

### 5.4 Dynamic Retrieval Balances Performance and Efficiency

The dynamic retrieval strategy effectively balances performance and efficiency. By incorporating CBDR, the system significantly reduces retrieval overhead while maintaining competitive accuracy—improving by 0.9 pp under the Top-3 setting. Although a slight decrease in accuracy (0.2 pp) was observed in the Top-1 scenario, this approach demonstrates practical potential for balancing efficiency and effectiveness in real-world applications.

# 6 Conclusion

This paper establishes internal hidden state confidence dynamics as a principled signal for optimizing RAG systems. By quantifying confidence shifts induced by retrieved contexts, we enable precise Reranker alignment and adaptive retrieval activation. Our framework CBDR has brought more efficient performance to the RAG system, which has practical application value. Our work bridges parametric and external knowledge while providing a generalizable paradigm for evaluating knowledge boundary interactions. Future work will extend this approach to multimodal RAG and multi-documents knowledge scenarios.

# References

[1] Wayne Xin Zhao, Kun Zhou, Junyi Li, Tianyi Tang, Xiaolei Wang, Yupeng Hou, Yingqian Min, Beichen Zhang, Junjie Zhang, Zican Dong, et al. A survey of large language models. *arXiv preprint arXiv:2303.18223*, 1(2), 2023.

[2] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.

[3] Binyuan Hui, Jian Yang, Zeyu Cui, Jiaxi Yang, Dayiheng Liu, Lei Zhang, Tianyu Liu, Jiajun Zhang, Bowen Yu, Keming Lu, et al. Qwen2. 5-coder technical report. *arXiv preprint arXiv:2409.12186*, 2024.

[4] Ziwei Ji, Tiezheng Yu, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. Towards mitigating llm hallucination via self reflection. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 1827–1843, 2023.

[5] Ariana Martino, Michael Iannelli, and Coleen Truong. Knowledge injection to counter large language model (llm) hallucination. In *European Semantic Web Conference*, pages 182–185. Springer, 2023.

[6] Gautier Izacard and Edouard Grave. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.

[7] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, et al. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474, 2020.

[8] Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan Berant. Making retrieval-augmented language models robust to irrelevant context. *arXiv preprint arXiv:2310.01558*, 2023.

[9] Feiteng Fang, Yuelin Bai, Shiwen Ni, Min Yang, Xiaojun Chen, and Ruifeng Xu. Enhancing noise robustness of retrieval-augmented language models with adaptive adversarial training. *arXiv preprint arXiv:2405.20978*, 2024.

[10] Florin Cuconasu, Giovanni Trappolini, Federico Siciliano, Simone Filice, Cesare Campagnano, Yoelle Maarek, Nicola Tonellotto, and Fabrizio Silvestri. The power of noise: Redefining retrieval for rag systems. In *Proceedings of the 47th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 719–729, 2024.

[11] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM computing surveys*, 55(12):1–38, 2023.

[12] Qingxiu Dong, Jingjing Xu, Lingpeng Kong, Zhifang Sui, and Lei Li. Statistical knowledge assessment for large language models. *Advances in Neural Information Processing Systems*, 36:29812–29830, 2023.

[13] Xunjian Yin, Xu Zhang, Jie Ruan, and Xiaojun Wan. Benchmarking knowledge boundary for large language models: A different perspective on model evaluation. *arXiv preprint arXiv:2402.11493*, 2024.

[14] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837, 2022.

[15] Ruiyang Ren, Yuhao Wang, Yingqi Qu, Wayne Xin Zhao, Jing Liu, Hao Tian, Hua Wu, Ji-Rong Wen, and Haifeng Wang. Investigating the factual knowledge boundary of large language models with retrieval augmentation. *arXiv preprint arXiv:2307.11019*, 2023.

[16] Moxin Li, Yong Zhao, Wenxuan Zhang, Shuaiyi Li, Wenya Xie, See-Kiong Ng, Tat-Seng Chua, and Yang Deng. Knowledge boundary of large language models: A survey. *arXiv preprint arXiv:2412.12472*, 2024.

[17] Hang Zheng, Hongshen Xu, Yuncong Liu, Lu Chen, Pascale Fung, and Kai Yu. Enhancing llm reliability via explicit knowledge boundary modeling. *arXiv preprint arXiv:2503.02233*, 2025.

[18] Shiyu Ni, Keping Bi, Jiafeng Guo, Lulu Yu, Baolong Bi, and Xueqi Cheng. Towards fully exploiting llm internal states to enhance knowledge boundary perception.

[19] Xin Sun, Jianan Xie, Zhongqi Chen, Qiang Liu, Shu Wu, Yuehe Chen, Bowen Song, Weiqiang Wang, Zilei Wang, and Liang Wang. Divide-then-align: Honest alignment based on the knowledge boundary of rag. *arXiv preprint arXiv:2505.20871*, 2025.

[20] Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu, Xipeng Qiu, and Xuanjing Huang. Do large language models know what they don't know? *arXiv preprint arXiv:2305.18153*, 2023.

[21] Shiyu Ni, Keping Bi, Jiafeng Guo, and Xueqi Cheng. When do llms need retrieval augmentation? mitigating llms' overconfidence helps retrieval augmentation. *arXiv preprint arXiv:2402.11457*, 2024.

[22] Shiyu Ni, Keping Bi, Lulu Yu, and Jiafeng Guo. Are large language models more honest in their probabilistic or verbalized confidence? In *China Conference on Information Retrieval*, pages 124–135. Springer, 2024.

[23] Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald Clark, Quoc V Le, Christopher Ré, and Azalia Mirhoseini. Large language monkeys: Scaling inference compute with repeated sampling. *arXiv preprint arXiv:2407.21787*, 2024.

[24] Rachel Longjohn, Giri Gopalan, and Emily Casleton. Statistical uncertainty quantification for aggregate performance metrics in machine learning benchmarks. *arXiv preprint arXiv:2501.04234*, 2025.

[25] Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu, Zhijing Wu, Yujia Zhou, and Yiqun Liu. Unsupervised real-time hallucination detection based on the internal states of large language models. *arXiv preprint arXiv:2403.06448*, 2024.

[26] Chao Chen, Kai Liu, Ze Chen, Yi Gu, Yue Wu, Mingyuan Tao, Zhihang Fu, and Jieping Ye. Inside: Llms' internal states retain the power of hallucination detection. *arXiv preprint arXiv:2402.03744*, 2024.

[27] Zihao Feng, Xiaoxue Wang, Ziwei Bai, Donghang Su, Bowen Wu, Qun Yu, and Baoxun Wang. Improving generalization in intent detection: Grpo with reward-based curriculum sampling. *arXiv preprint arXiv:2504.13592*, 2025.

[28] Maria Marina, Nikolay Ivanov, Sergey Pletenev, Mikhail Salnikov, Daria Galimzianova, Nikita Krayko, Vasily Konovalov, Alexander Panchenko, and Viktor Moskvoretskii. Llm-independent adaptive rag: Let the question speak for itself. *arXiv preprint arXiv:2505.04253*, 2025.

[29] Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao, Linmei Hu, Weichuan Liu, Lei Hou, and Juanzi Li. Seakr: Self-aware knowledge retrieval for adaptive retrieval augmented generation. *arXiv preprint arXiv:2406.19215*, 2024.

[30] Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Rich James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. Replug: Retrieval-augmented black-box language models. *arXiv preprint arXiv:2301.12652*, 2023.

[31] Youan Cong, Cheng Wang, Pritom Saha Akash, and Kevin Chen-Chuan Chang. Query optimization for parametric knowledge refinement in retrieval-augmented large language models. *arXiv preprint arXiv:2411.07820*, 2024.

[32] Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen Wang, Ji-Rong Wen, and Zhicheng Dou. Understand what llm needs: Dual preference alignment for retrieval-augmented generation. In *Proceedings of the ACM on Web Conference 2025*, pages 4206–4225, 2025.

[33] Pengyue Jia, Derong Xu, Xiaopeng Li, Zhaocheng Du, Xiangyang Li, Yichao Wang, Yuhao Wang, Qidong Liu, Maolin Wang, Huifeng Guo, et al. Bridging relevance and reasoning: Rationale distillation in retrieval-augmented generation. *arXiv preprint arXiv:2412.08519*, 2024.

[34] Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Patel, Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-Ziv. Layer by layer: Uncovering hidden representations in language models. *arXiv preprint arXiv:2502.02013*, 2025.

[35] Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Aurojit Panda, Jinyang Li, and He He. Reasoning models know when they're right: Probing hidden states for self-verification. *arXiv preprint arXiv:2504.05419*, 2025.

[36] Amos Azaria and Tom Mitchell. The internal state of an llm knows when it's lying. *arXiv preprint arXiv:2304.13734*, 2023.

[37] Tom Kwiatkowski, Jennimaria Palomaki, Olivia Redfield, Michael Collins, Ankur Parikh, Chris Alberti, Danielle Epstein, Illia Polosukhin, Matthew Kelcey, Jacob Devlin, Kenton Lee, Kristina N. Toutanova, Llion Jones, Ming-Wei Chang, Andrew Dai, Jakob Uszkoreit, Quoc Le, and Slav Petrov. Natural questions: a benchmark for question answering research. *Transactions of the Association of Computational Linguistics*, 2019.

[38] Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W Cohen, Ruslan Salakhutdinov, and Christopher D Manning. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.

[39] Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. *arXiv e-prints*, pages arXiv–2407, 2024.

[40] Qwen Team. Qwen2.5: A party of foundation models, September 2024.