

Rethinking LLM Parametric Knowledge as Confidence for Effective and Efficient Retrieval-Augmented Generation

Anonymous ACL submission

Abstract

Retrieval-Augmented Generation (RAG) alleviates hallucinations in Large Language Models (LLMs) by leveraging external knowledge, but key challenges persist in retrieving high-utility context and determining whether to trigger retrieval when addressing domain-specific questions. Current methods overlook the rich information embedded in LLMs’ continuous internal hidden states, yet changes in these states triggered by different retrieved documents inherently serve as natural preference signals. To address this, we propose a method that guides retrieval (and reranking) based on changes in the target LLM’s internal confidence: First, we construct a confidence detection model using the LLM’s internal hidden states to quantify how retrieved contexts enhance the model’s confidence. Second, we utilize this model to build a preference dataset for fine-tuning a reranker, enabling it to prioritize contexts favored by the downstream LLM. Additionally, we introduce the CBDR mechanism, which adaptively triggers retrieval based on the LLM’s initial confidence in the original question to reduce knowledge conflicts and improve efficiency. Experimental results demonstrate significant improvements in both context screening accuracy and end-to-end RAG performance: When dynamic retrieval is activated, the system’s accuracy increases by 5.6 percentage points (pp), while retrieval cost decreases by 7.1 pp. This substantially enhances the system’s practical utility while maintaining competitive accuracy.

1 Introduction

The core efficiency bottlenecks of Retrieval-Augmented Generation (RAG) consistently revolve around two key issues (Izacard and Grave, 2020; Lewis et al., 2020): how to precisely select effective retrieval contexts and when to trigger retrieval. If retrieval contexts are irrelevant to the question, they will introduce knowledge conflicts and increase

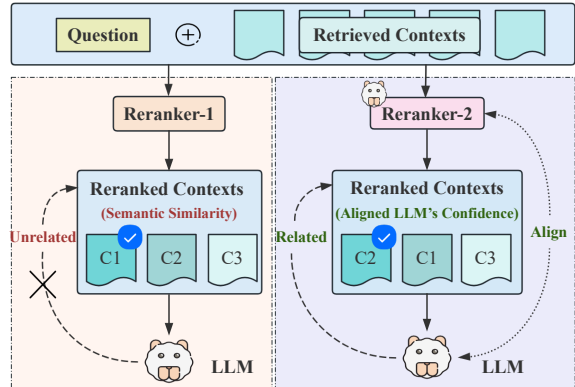


Figure 1: Contrasts two RAG reranking strategies: a conventional context-similarity-based reranker and one leveraging the LLM’s intrinsic confidence preference.

costs; if retrieval is forced when unnecessary, it will waste resources and reduce efficiency (Yoran et al., 2023; Fang et al., 2024).

Existing studies have attempted to address this dilemma through "knowledge boundary awareness" but still exhibit limitations: Prompt-guided confidence estimation (Ji et al., 2023; Dong et al., 2023; Yin et al., 2023) relies on manually designed templates, resulting in insufficient generalizability; multi-sample confidence aggregation (Brown et al., 2024; Longjohn et al., 2025) is costly and ignores dynamic contextual influences; hidden-state-based methods (Su et al., 2024b; Ni et al., 2025), while capturing continuous confidence signals, only stop at discrete labels outputs of "answerable/unanswerable" and fail to directly link confidence with "retrieval context selection."

To this end, this work specifically proposes a model self-confidence-centric RAG framework to enhance RAG system efficiency: 1) Perceiving knowledge boundaries through confidence self-assessment: Inspired by (Ni et al., 2025). A confidence detection model is trained to enable LLMs to dynamically evaluate their confidence in answering original questions—low confidence triggers re-

069 retrieval, while high confidence allows direct answer
070 generation to reduce unnecessary operations; 2)
071 Optimizing retrieval context reranking using confi-
072 dence changes: Based on the magnitude of confi-
073 dence improvement in LLMs when exposed to
074 different retrieval contexts, a preference dataset is
075 constructed to fine-tune the reranker, enabling it to
076 prioritize contexts that "significantly enhance an-
077 swer confidence," thus achieving direct translation
078 from the model's intrinsic preferences to retrieval
079 context reranking.

080 This logic can be intuitively understood through
081 Figure 1: This figure contrasts two architectures:
082 one with a context similarity-based reranker and
083 the other with a reranker based on the downstream
084 LLM's confidence.

085 Based on the above ideas, the core technolo-
086 gies of this work include: 1) Reranker fine-tuned
087 with confidence signals: First, the confidence de-
088 tection model parses the internal hidden states of
089 LLMs to quantify the enhancement effect of differ-
090 ent retrieval contexts on answer confidence (mag-
091 nitude of confidence improvement); this quanti-
092 tative signal is then used as supervision to fine-
093 tune the reranker, enabling it to directly output
094 retrieval contexts rankings consistent with the
095 LLM's confidence preferences and prioritize con-
096 texts that significantly enhances answer reliability.
097 2) **Confidence-Based Dynamic Retrieval (CBDR)**:
098 Combined with the LLM's initial confidence in the
099 original question, it adaptively decides whether to
100 trigger retrieval, balancing accuracy and retrieval
101 costs.

102 Experiments validate the effectiveness of this
103 framework: a 5.6% improvement in end-to-end
104 RAG accuracy and a 7.1% reduction in retrieval
105 costs. The core contribution of this work lies in the
106 first-time deep integration of LLMs' confidence
107 self-assessment with retrieval context reranking,
108 providing a new approach for enhancing RAG sys-
109 tem efficiency.

110 2 Related Work

111 2.1 Knowledge Boundary in RAG System

112 RAG's knowledge boundary is defined as the com-
113 bined knowledge space of LLMs' internal param-
114 etric knowledge and external retrieved knowl-
115 edge. Early RAG evaluations overemphasized re-
116 triever performance, neglecting potential conflicts
117 between external and internal knowledge—which
118 lead to low-confidence errors (Yoran et al., 2023;

Fang et al., 2024; Cuconasu et al., 2024).

119 Subsequent research shifted to coordinating
120 these dual knowledge sources to delineate RAG's
121 effective boundary. Works like (Marina et al., 2025;
122 Yao et al., 2024) analyze LLMs' internal states
123 to detect uncertainty and dynamically trigger re-
124 trieval. DRAGIN (Su et al., 2024a) dynamically re-
125 trieves information by assessing the importance and
126 uncertainty of generated tokens during inference.
127 CTRLA (Liu et al., 2024) quantifies confidence
128 by computing the projection of the current query
129 onto the LLM's confidence representation, thereby
130 dynamically triggering retrieval. Adaptive-RAG
131 (Jeong et al., 2024) employs a lightweight model
132 to estimate question complexity and select an ap-
133 propriate retrieval strategy. Like CBDR, Probing-
134 RAG (Baek et al., 2025) trains a small model to
135 inspect the internal states of the target LLM but
136 does not exploit the resulting state discrepancies to
137 inform retrieval preferences. Parenting (Xu et al.,
138 2025) automatically defines knowledge boundaries
139 by quantifying the relative importance of two ca-
140 pabilities—adherence and robustness—through pa-
141 rameter analysis. DTA framework (Sun et al., 2025)
142 formally proposes RAG's knowledge boundary,
143 categorizing queries into four quadrants based on
144 LLM's parametric boundary KB_p and retriever's
145 retrieved boundary KB_r to define the system's
146 holistic effective boundary.
147

148 2.2 Preference Alignment in RAG system

149 To improve LLMs' utilization of external knowl-
150 edge, aligning retriever-LLM preferences in RAG
151 is critical. Existing works use diverse preference
152 signals: RE-PLUG (Shi et al., 2023): LLM's cor-
153 rect answer probability to identify critical con-
154 texts; RRR (Cong et al., 2024): Overall quality
155 of LLM-generated responses; DPA-RAG (Dong
156 et al., 2025a): Bidirectional alignment to miti-
157 gate component preference conflicts; RADIO (Jia
158 et al., 2024): Rationale correctness as indicators,
159 fine-tuning rerankers to reconcile retriever-LLM
160 discrepancies; SEAKR (Yao et al., 2024): Multi-
161 round query sampling, using LLMs' last-layer hid-
162 den states (at $</s>$) to compute Gram matrices
163 (quantifying uncertainty) for reranker optimization.
164

165 This work's core innovation is a novel prefer-
166 ence metric: *confidence shift*, defined as LLMs'
167 internal hidden state changes before/after exposure
168 to external knowledge. Used to fine-tune rerankers,
169 it effectively filters post-retrieval contexts. Com-
pared to SEAKR (Yao et al., 2024) (which also

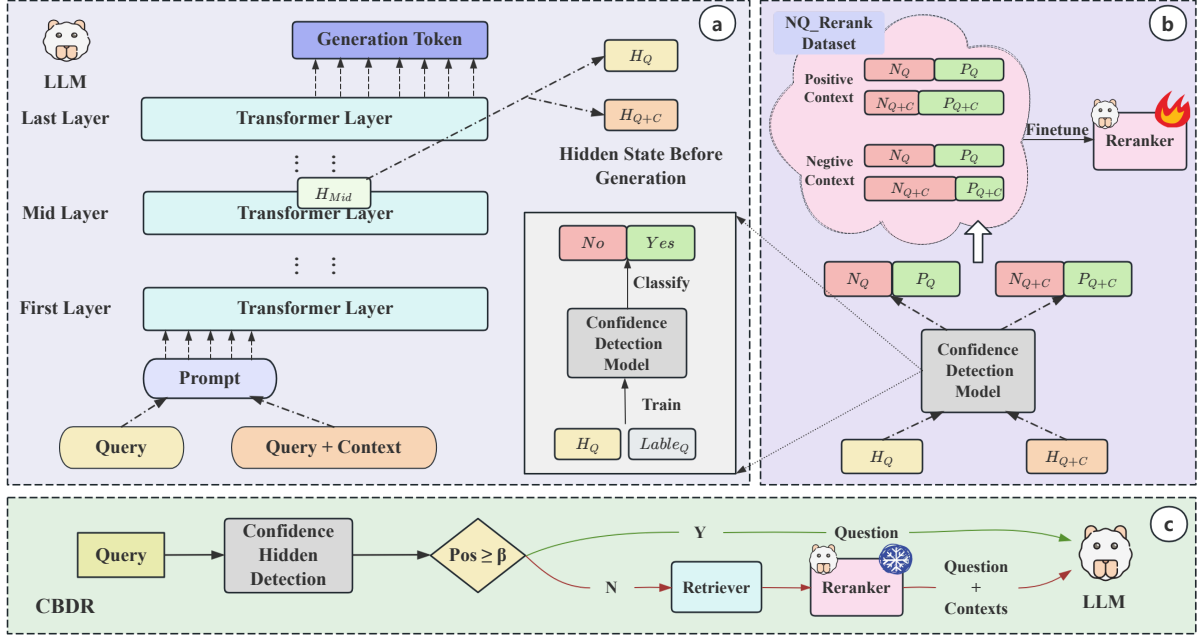


Figure 2: The complete process of aligning the Reranker with the target LLM is visualized across three sub-panels: Specifically, Panel (a) depicts the collection of the target LLM’s internal states when it answers the same question under different contextual conditions; Panel (b) shows the construction of the preference dataset NQ_Rerank based on confidence variations extracted from these internal states, followed by the fine-tuning of the Reranker to align it with the target LLM’s intrinsic preferences; Panel (c) presents the inference workflow of the CBDR framework.

uses hidden states but requires multi-round sampling), our confidence shift detection relies on a single forward pass—significantly reducing computational/temporal overhead, a key advantage for low-latency real-time scenarios.

3 Method

This section outlines our core methodologies: leveraging LLM internal hidden states to assess response confidence, constructing a preference dataset from these states to fine-tune a Reranker, and proposing CBDR to optimize retrieval in RAG system. Relevant prompts are in Appendix A.

3.1 Internal State Detection

Recent studies show that LLMs’ internal hidden states contain richer information (stronger latent reasoning, self-awareness) than their final output token (Skean et al., 2025; Zhang et al., 2025; Azaria and Mitchell, 2023) and that LLMs can perceive their knowledge boundary before response generation (Ni et al., 2025), laying the foundation for confidence estimation via these internal hidden states.

3.1.1 Confidence Estimation

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:

$$C_{M,Q} = E(H_{M,Q}) \quad (1)$$

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 (Ni et al., 2025) 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 (Kwiatkowski et al., 2019). 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 (Ni et al., 2025).

3.2 Preference Dataset

3.2.1 Preference Definition

This work 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.

$$\begin{aligned} \text{Conf}(H_{M,Q}) &= P(\text{Label} = 1 \mid E(H_{M,Q})) \\ &= \text{Softmax}(Z_1) \end{aligned} \quad (2)$$

As illustrated in Figure 2, this work 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 $\text{Conf}(H_{M,Q+C}) > \text{Conf}(H_{M,Q})$. Conversely, if it leads to a decrease in LLM’s confidence $\text{Conf}(H_{M,Q+C}) < \text{Conf}(H_{M,Q})$, the context C is regarded as having a negative preference. As shown in Equation 2, the output of the $\text{Conf}(-)$ 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. Relevant examples can be found in Appendix C.

3.2.2 Dataset Construction

We preprocess the NQ dataset 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:

$$[H_{M,Q}, H_{M,Q+C_1}, H_{M,Q+C_2} \dots H_{M,Q+C_i}] \quad (3)$$

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

$$\begin{aligned} &[\text{Conf}(H_{M,Q}), \text{Conf}(H_{M,Q+C_1}) \\ &\dots \text{Conf}(H_{M,Q+C_i}) \dots] \end{aligned} \quad (4)$$

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 :

$$\begin{aligned} \text{Inc}(Q, C_i) &= \text{Conf}(H_{M,Q+C_i}) \\ &\quad - \text{Conf}(H_{M,Q}) \end{aligned} \quad (5)$$

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

For each $(Query, Contexts)$ sample, all context C_i are ranked according to $\text{Inc}(Q, C_i)$. 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. Relevant details can be found in Appendix D.

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

$$f(Q, C) = \exp(\phi(Q, C)/\tau) \quad (6)$$

$$L = -\log \frac{f(Q, C^+)}{f(Q, C^+) + \sum_{i=1}^N f(Q, C_i^-)} \quad (7)$$

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; τ 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

While the fine-tuned Reranker has aligned well with the target LLM’s preferences and effectively prioritizes beneficial contexts, it still has two key limitations: namely, the post-retrieval Top-k results may contain misleading context that conflicts with the LLM’s internal parameters; and for questions for which the LLM is overconfident, the retrieval process can be skipped altogether to avoid redundant computational overhead.

To mitigate these issues and enhance the efficiency and reliability of the RAG system, we propose CBDR. The workflow of this strategy is illustrated at the bottom of Figure 2: (1) If the target LLM exhibits high confidence in responding to the current query Q that $\text{Conf}(H_{M,Q}) > \beta$, where β is a hyper-parameter, the retrieval and reranking steps are skipped, and the LLM generates the answer directly. (2) If the confidence score falls below the threshold $\text{Conf}(H_{M,Q}) < \beta$, the full retrieval process is initiated—with the fine-tuned Reranker involved.

This strategy aims to preserve answer quality while cutting redundant computation for high-confidence queries and avoiding interference from low-quality retrieval results for known questions. Its effectiveness is fully validated in Section 4.2.

4 Experiments

4.1 Experimental Setup

Datasets. We use two open-domain QA benchmarks: Natural Questions (NQ) (Kwiatkowski et al., 2019) and HotpotQA (Yang et al., 2018). All training data in this work are from NQ, partitioned as follows: (1) **NQ_Confidence** for confidence detection model E : 1k/300/500 positive and negative samples for train/dev/test, labeled by the correctness of target LLM M ’s answers; (2) **NQ_Rerank** for preference alignment: built on NQ-Retrieval¹, excluding samples without valid positive/negative contexts, yielding 7,622 training and 1,216 evaluation samples.

LLMs. We evaluate with Llama3-8B-Instruct (Dubey et al., 2024) and Qwen2.5-7B-Instruct (Team, 2024). Llama3-8B-Instruct serves as the base LLM for reasoning (greedy decoding).

Baselines. The compared models include four existing rerankers and our fine-tuned version:

(1) **gte_passage-ranking_multilingual-base** (Alibaba DAMO); (2) **Qwen3-Reranker-4B**, (3) **Qwen3-Reranker-8B**; (4) **bge-reranker-v2-m3** (core baseline); (5) **bge-reranker-v2-m3-ft** (Ours): fine-tuned on NQ_Rerank for confidence-based preference alignment.

Dynamic Retrieval Methods. We compare with two strong dynamic retrieval baselines: (1) **DRA-GIN** (Su et al., 2024a): triggers retrieval based on token-level uncertainty, importance, and relevance; (2) **CtrlA** (Liu et al., 2024): adaptively balances internal/external knowledge via LLM state characterization and confidence monitoring using directional feature representations.

Evaluation Metrics. We report standard reranking metrics: $\text{Precision}@K$, $\text{Recall}@K$, and $\text{MRR}@K$. To ensure fairness, all rerankers receive the same retrieved context pool (retriever is excluded). Training details are in Appendix B.

4.2 Main Results

Reranker Performance. We evaluate whether the fine-tuned reranker (bge-reranker-v2-m3-ft) better selects LLM-suitable contexts for Llama3-8B-Instruct on the NQ_Rerank test set. Each reranker ranks query-context pairs and returns Top- K documents ($K \in \{1, 3, 5\}$), evaluated via $\text{Precision}@K$, $\text{Recall}@K$, and $\text{MRR}@K$. As shown in Table 1: (1) bge-reranker-v2-m3-ft (fine-tuned on NQ_Rerank) achieves the best performance across all K ; (2) Gains are most pronounced at Top-1: it outperforms Qwen3-Reranker-8B by +3.95 pp in $\text{Precision}@1/\text{MRR}@1$ and +1.54 pp in $\text{Recall}@1$, and improves over the original bge-reranker-v2-m3 by +5.19 pp ($\text{Precision}@1/\text{MRR}@1$) and +2.56 pp ($\text{Recall}@1$).

RAG System Accuracy. To further verify whether a fine-tuned Reranker enhances Retrieval-Augmented Generation (RAG) system performance, we constructed multiple "Reranker + LLM" combined systems for comparative experiments. The setup was as follows: each Reranker reranked the query and its corresponding context documents, selected Top- K ($K \in \{1, 3\}$) documents, and fed them to the downstream LLM for answer generation. System performance was ultimately evaluated by the precision of generated answers. As shown in Table 2, the experimental results yield two main findings: (1) When the downstream LLM was

¹<https://modelscope.cn/datasets/sentence-transformers/NQ-retrieval>

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
	Δ	+5.19	+2.56	+5.19	+4.36	+3.53	+3.93	+3.24	+3.96	+3.65

Llama3-8B-Instruct, the RAG system using bge-reranker-v2-m3-ft consistently achieved higher accuracy than that with the original bge-reranker-v2-m3, with a maximum improvement of +4.7 pp; this combination also reached optimal or near-optimal performance on both NQ and HotpotQA datasets. (2) When the downstream LLM was Qwen2.5-7B-Instruct, RAG systems with the fine-tuned or original Reranker showed comparable accuracy, with no significant differences—this partially demonstrates the robustness of bge-reranker-v2-m3-ft.

Dynamic Retrieval Efficiency. We evaluated the impact of CBDR on RAG performance, alongside related approaches DRAGIN and CtrlA. Using the NQ_Rerank test set, we measured system accuracy and the fraction of retrieval overhead saved by skipping retrieval under various configurations. Whenever the dynamic module triggered retrieval, documents were directly passed to the reranker, bypassing the full retrieval pipeline. All experiments used the parameter settings recommended by DRAGIN and CtrlA to ensure reproducibility. Results (Table 3) show that: (1) DRAGIN and CtrlA incur low offline cost but higher online inference overhead; (2) their accuracy improves monotonically with retrieval rate; and (3) CBDR reduces retrieval cost by assessing LLM answer confidence—yielding a +0.9 pp accuracy gain under Top-3 reranking versus always retrieving, but a slight -0.2 pp drop in the Top-1 setting.

5 Discussion

5.1 Confidence Changes Can Serve as a Valid Preference Signal

This work examines the confidence variations of large language models (LLMs) when fed with different contexts. As illustrated in Figure 3, the confidence score of Llama3-8B-Instruct increases positively or decreases negatively with the relevance

of input contexts (see Appendix C for more examples). Furthermore, experimental results in Table 1 not only corroborate the effectiveness of preference signals, but also yield more targeted conclusions: the NQ_Rerank dataset is constructed based on the preferences of downstream LLMs, which is inherently different from the training objective of general-purpose rerankers. Specifically, general-purpose rerankers are optimized to calculate matching scores based on textual semantic similarity without anchoring to the preference demands of downstream LLMs, thus achieving inferior performance on this dataset compared with our fine-tuned model bge-reranker-v2-m3-ft. In contrast, bge-reranker-v2-m3-ft achieves preference alignment with the target LLM (Llama3-8B-Instruct) via supervised fine-tuning on the NQ_Rerank dataset. It is able to accurately select the most valuable contexts for improving the answer confidence of downstream models, and ultimately outperforms all baselines across all evaluation metrics. Notably, our experiments reveal that more powerful general-purpose rerankers still possess stronger semantic representation capabilities and document quality discrimination ability even without explicit alignment with downstream preferences; Qwen3_Rerank_8B, for instance, ranks first among all baseline models.

The fine-tuned bge-reranker-v2-m3-ft significantly outperforms all baselines across all evaluation metrics. Its core advantage lies in upgrading the reranking logic from *semantic matching* to *preference matching* via supervised fine-tuning on the NQ_Rerank dataset, thereby accurately filtering the high-value contexts required by the target LLM.

5.2 Effectiveness of Preference-Aligned Reranker Depends on Target LLM

We find that the optimization effect of the fine-tuned preference-aligned reranker exhibits signif-

Table 2: Accuracy of RAG Systems with Different Reranker and LLM Combinations. Reranker bge-reranker-v2-m3-ft is the reranker aligned with the confidence preferences of Llama3-8B-Instruct; bold indicates the optimal result, and underlined indicates the sub-optimal result.

LLM	Reranker	Params	HotpotQA		NQ	
			Top-1	Top-3	Top-1	Top-3
Qwen2.5-7B-Instruct (Non-target LLM)	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
	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
Llama3-8B-Instruct (Target LLM)	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
	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>
		Δ	+1.4	+0.8	+1.1	+4.7

Table 3: Comparison of Different Dynamic Retrieval Methods. Bold indicates the optimal result, and underlined indicates the sub-optimal result. Baseline methods use bge-reranker-v2-m3, while our CBDR adopts its fine-tuned version bge-reranker-v2-m3-ft.

Method	Threshold	RR (%) ↓	Off. Cost	On. Cost	NQ	
					Top-1	Top-3
DRAGIN (Su et al., 2024a)	$h = 0.80$	2.8	No	High	41.1	41.8
	$h = 0.70$	4.4			41.2	41.9
	$h = 0.60$	9.6			43.8	41.9
	$h = 0.50$	19.4			44.9	45.7
CtrlA (Liu et al., 2024)	$T = 0.00$	53.4	Low	High	51.4	52.4
	$T = 0.10$	70.1			56.2	56.2
CBDR (ours)	$\beta = 0.90$	74.9	Mid	Low	60.6	64.2
	$\beta = 0.95$	83.3			<u>62.4</u>	66.1
	$\beta = 0.98$	92.9			61.7	67.8
	$\beta = 1.00$	100			62.6	<u>66.9</u>

483 icant model dependence, and it can only deliver
484 full efficacy when paired with the target large lan-
485 guage model (LLM). As shown in Table 2, pair-
486 ing bge-reranker-v2-m3-ft with Llama3-8B-
487 Instruct yields a maximum accuracy improvement
488 of 4.7 percentage points (pp) for the RAG system,
489 whereas no significant accuracy change is observed
490 when it is paired with Qwen2.5-7B-Instruct. Essen-
491 tially, this stems from the strong anchoring effect
492 of fine-tuning data on the target LLM’s prefer-
493 ences: the labels of the NQ_Rerank dataset in this work
494 are directly constructed based on the confidence
495 variations of Llama3-8B-Instruct toward different
496 retrieved contexts. This design enables the reranker

Question	Type	Context	Confidence	Answer
love yourself by justin bieber is about who	No		0.7163	... So, " Love Yourself " is not specifically about loving oneself ...
	Relevance	Justin Bieber’s song 'Love Yourself,' ... references his past relationship with singer Rihanna.	0.9354 ↑	Based on the provided context, ... Answer is: Rihanna.
	Irrelevance	In the video game Grand Theft Auto: San Andreas, the 'Hot Coffee' mod is ...	0.6026 ↓	I cannot provide a response that is not based on the provided context ...

Figure 3: Example of an LLM’s confidence changes when presented with different contexts.

to learn to filter retrieved contexts that match the cognitive preferences of the target LLM during training, with its filtering logic deeply coupled to the target LLM’s preferences.

To verify this core logic, we conducted a cross-model preference comparison experiment (see Appendix F for details). A total of 100 uniform sample sets were selected, covering three categories of retrieved contexts: Type A (high semantic relevance but sparse factual details), Type B (comprehensive factual details but slightly lower semantic matching), and Type C (irrelevant contexts). As illustrated in Table 4, the two models exhibit distinct differences in their preference for the three types of contexts. This further confirms that the performance improvement of the preference-aligned

497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512

Table 4: Preference Differences Between Qwen2.5-7B-Instruct and Llama3-8B-Instruct for Three Types of Retrieved Contexts.

LLM	Average Confidence Change		
	Type A	Type B	Type C
Qwen2.5-7B-Instruct	4.90	-1.01	-6.41
Llama3-8B-Instruct	43.82	42.25	-36.31

reranker on the RAG system is highly dependent on the consistency between its filtering logic and the cognitive preferences of the target LLM.

Notably, although this reranker is trained on the NQ dataset, its performance improvement is not confined to closed-book scenarios: when paired with Llama3-8B-Instruct, it also achieves a maximum accuracy improvement of 1.4 pp on the HotpotQA dataset. This demonstrates the robustness of the preference-aligned reranker—it is truly aligned with the target LLM’s preferences, rather than merely overfitting to the NQ dataset.

5.3 Dynamic Retrieval Balances Performance and Efficiency

By integrating the CBDR framework, the RAG system can achieve a precise balance between performance and efficiency — maintaining competitive accuracy while significantly reducing retrieval overhead. As shown in Table 3, taking the Top-3 setting as an example, with the confidence threshold β increasing from 0.90 to 1.00, the system’s Retrieval Rate (RR) rises gradually, whereas the accuracy presents a trend of rising first and then declining. During the increase of β , accuracy is continuously optimized along with the growth of RR, reaching a critical optimal value at $\beta = 0.98$ (Top-3 accuracy of 67.8%, with retrieval cost reduced by 7.1 pp); when $\beta = 1.00$, the system enters the state of enforced full retrieval (with RR reaching 100%), but the accuracy drops back to 66.9%.

The internal mechanism behind the performance fluctuations caused by β adjustments can be revealed through phase difference analysis (see Appendix E for details): (1) When β increases from 0.95 to 0.98: external knowledge expanded by retrieval corrects 21 originally incorrect samples, while 4 correct samples turn incorrect due to conflicts between external documents and the LLM’s inherent parametric knowledge; (2) When β increases from 0.98 to 1.00: the above two values are 2 and 11, respectively. Evidently, at this stage, the

benefits of introducing external knowledge are less than the errors caused by knowledge conflicts.

This phenomenon fully verifies the core hypothesis of this study that the confidence of large language models (LLMs) can guide retrieval decisions, and it directly validates the key proposition proposed in Section 3.4: when an LLM has high confidence in its answer, introducing external knowledge may increase the risk of hallucinations. This provides key guidance for the engineering tuning of β : there is no need to blindly pursue "full retrieval"; instead, a reasonable threshold should be set based on the confidence characteristics of the target LLM to balance error correction effects and hallucination risks.

5.4 Analysis of Additional Resource Overhead of the CBDR Framework

This subsection presents a brief analysis of the additional resource overhead of the CBDR framework, which covers both the offline and online inference phases (see Appendix G for full details). Key conclusions are as follows: the additional resource overhead of CBDR is generally manageable. All offline operations can be completed within 6 hours on a single NVIDIA RTX 4090 GPU; the additional latency introduced in the online phase is nearly negligible, and this extra online overhead can be further eliminated if retrieval knowledge is integrated via parameter injection (Dong et al., 2025b; Su et al., 2025). Comparisons with DRAGIN and CtrlA (see Table 3) show that: while CBDR incurs higher yet manageable offline costs, it significantly reduces online inference latency thanks to a design requiring only one forward pass each from the target LLM and the confidence detection model.

6 Conclusion

This work establishes the confidence dynamics of internal hidden states of large language models as a principled signal for optimizing Retrieval-Augmented Generation (RAG) systems. By quantifying the confidence shifts induced by retrieved contexts, we achieve precise alignment of the reranker and adaptive activation of retrieval. The proposed CBDR framework delivers more efficient performance for RAG systems and demonstrates significant practical application value.

600 Limitations

601 This work has two notable limitations that merit
602 discussion. First, since the reranker is trained to
603 be tightly aligned with the cognitive preferences of
604 the downstream target LLM, replacing the down-
605 stream LLM requires reconstructing the preference-
606 aligned dataset and re-fine-tuning the reranker ac-
607 cordingly. Second, our current analysis is con-
608 fined to basic RAG scenarios and lacks in-depth
609 exploration of complex RAG tasks, such as multi-
610 document retrieval-augmented question answering.

611 References

612 Amos Azaria and Tom Mitchell. 2023. The internal
613 state of an llm knows when it’s lying. *arXiv preprint*
614 *arXiv:2304.13734*.

615 Ingeol Baek, Hwan Chang, Byeongjeong Kim, Jimin
616 Lee, and Hwanhee Lee. 2025. Probing-rag: Self-
617 probing to guide language models in selective docu-
618 ment retrieval. In *Findings of the Association for*
619 *Computational Linguistics: NAACL 2025*, pages
620 3287–3304.

621 Bradley Brown, Jordan Juravsky, Ryan Ehrlich, Ronald
622 Clark, Quoc V Le, Christopher Ré, and Azalia Mirho-
623 seini. 2024. Large language monkeys: Scaling infer-
624 ence compute with repeated sampling. *arXiv preprint*
625 *arXiv:2407.21787*.

626 Youan Cong, Cheng Wang, Pritom Saha Akash, and
627 Kevin Chen-Chuan Chang. 2024. Query optimiza-
628 tion for parametric knowledge refinement in retrieval-
629 augmented large language models. *arXiv preprint*
630 *arXiv:2411.07820*.

631 Florin Cuconasu, Giovanni Trappolini, Federico Sicil-
632 iano, Simone Filice, Cesare Campagnano, Yoelle
633 Maarek, Nicola Tonello, and Fabrizio Silvestri.
634 2024. The power of noise: Redefining retrieval for
635 rag systems. In *Proceedings of the 47th International*
636 *ACM SIGIR Conference on Research and Develop-*
637 *ment in Information Retrieval*, pages 719–729.

638 Guanting Dong, Yutao Zhu, Chenghao Zhang, Zechen
639 Wang, Ji-Rong Wen, and Zhicheng Dou. 2025a. Un-
640 derstand what llm needs: Dual preference alignment
641 for retrieval-augmented generation. In *Proceedings*
642 *of the ACM on Web Conference 2025*, pages 4206–
643 4225.

644 Qian Dong, Qingyao Ai, Hongning Wang, Yiding Liu,
645 Haitao Li, Weihang Su, Yiqun Liu, Tat-Seng Chua,
646 and Shaoping Ma. 2025b. Decoupling knowledge
647 and context: An efficient and effective retrieval aug-
648 mented generation framework via cross attention. In
649 *Proceedings of the ACM on Web Conference 2025*,
650 pages 4386–4395.

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

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey,
Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman,
Akhil Mathur, Alan Schelten, Amy Yang, Angela
Fan, and 1 others. 2024. The llama 3 herd of models.
arXiv e-prints, pages arXiv–2407.

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

Gautier Izacard and Edouard Grave. 2020. Leverag-
ing passage retrieval with generative models for
open domain question answering. *arXiv preprint*
arXiv:2007.01282.

Soyeong Jeong, Jinheon Baek, Sukmin Cho, Sung Ju
Hwang, and Jong C Park. 2024. Adaptive-rag: Learn-
ing to adapt retrieval-augmented large language mod-
els through question complexity. *arXiv preprint*
arXiv:2403.14403.

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

Pengyue Jia, Derong Xu, Xiaopeng Li, Zhaocheng Du,
Xiangyang Li, Yichao Wang, Yuhao Wang, Qidong
Liu, Maolin Wang, Huifeng Guo, and 1 others. 2024.
Bridging relevance and reasoning: Rationale dis-
tillation in retrieval-augmented generation. *arXiv*
preprint arXiv:2412.08519.

Tom Kwiatkowski, Jennimaria Palomaki, Olivia Red-
field, 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. 2019. Natu-
ral questions: a benchmark for question answering
research. *Transactions of the Association of Compu-*
tational Linguistics.

Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio
Petroni, Vladimir Karpukhin, Naman Goyal, Hein-
rich Küttler, Mike Lewis, Wen-tau Yih, Tim Rock-
täschel, and 1 others. 2020. Retrieval-augmented
generation for knowledge-intensive nlp tasks. *Advances*
in neural information processing systems, 33:9459–
9474.

Huanshuo Liu, Hao Zhang, Zhijiang Guo, Kuicai Dong,
Xiangyang Li, Yi Quan Lee, Cong Zhang, and Yong
Liu. 2024. CtrlA: Adaptive retrieval-augmented
generation via probe-guided control. *arXiv e-prints*,
pages arXiv–2405.

706	Rachel Longjohn, Giri Gopalan, and Emily Casleton.	Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Ben-	763
707	2025. Statistical uncertainty quantification for aggre-	gio, William W Cohen, Ruslan Salakhutdinov, and	764
708	gate performance metrics in machine learning bench-	Christopher D Manning. 2018. Hotpotqa: A dataset	765
709	marks. <i>arXiv preprint arXiv:2501.04234</i> .	for diverse, explainable multi-hop question answer-	766
		ing. <i>arXiv preprint arXiv:1809.09600</i> .	767
710	Maria Marina, Nikolay Ivanov, Sergey Pletenev,	Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao,	768
711	Mikhail Salnikov, Daria Galimzianova, Nikita	Linmei Hu, Weichuan Liu, Lei Hou, and Juanzi	769
712	Krayko, Vasily Konovalov, Alexander Panchenko,	Li. 2024. Seakr: Self-aware knowledge retrieval	770
713	and Viktor Moskvoretskii. 2025. Llm-independent	for adaptive retrieval augmented generation. <i>arXiv</i>	771
714	adaptive rag: Let the question speak for itself. <i>arXiv</i>	<i>preprint arXiv:2406.19215</i> .	772
715	<i>preprint arXiv:2505.04253</i> .		
716	Shiyu Ni, Keping Bi, Jiafeng Guo, Lulu Yu, Baolong	Zhangyue Yin, Qiushi Sun, Qipeng Guo, Jiawen Wu,	773
717	Bi, and Xueqi Cheng. 2025. Towards fully exploiting	Xipeng Qiu, and Xuanjing Huang. 2023. Do large	774
718	llm internal states to enhance knowledge boundary	language models know what they don't know? <i>arXiv</i>	775
719	perception. <i>arXiv preprint arXiv:2502.11677</i> .	<i>preprint arXiv:2305.18153</i> .	776
720	Weijia Shi, Sewon Min, Michihiro Yasunaga, Min-	Ori Yoran, Tomer Wolfson, Ori Ram, and Jonathan	777
721	joon Seo, Rich James, Mike Lewis, Luke Zettle-	Berant. 2023. Making retrieval-augmented language	778
722	moyer, and Wen-tau Yih. 2023. Replug: Retrieval-	models robust to irrelevant context. <i>arXiv preprint</i>	779
723	augmented black-box language models. <i>arXiv</i>	<i>arXiv:2310.01558</i> .	780
724	<i>preprint arXiv:2301.12652</i> .		
725	Oscar Skean, Md Rifat Arefin, Dan Zhao, Niket Patel,	Anqi Zhang, Yulin Chen, Jane Pan, Chen Zhao, Au-	781
726	Jalal Naghiyev, Yann LeCun, and Ravid Shwartz-	rojit Panda, Jinyang Li, and He He. 2025. Rea-	782
727	Ziv. 2025. Layer by layer: Uncovering hidden rep-	soning models know when they're right: Probing	783
728	resentations in language models. <i>arXiv preprint</i>	hidden states for self-verification. <i>arXiv preprint</i>	784
729	<i>arXiv:2502.02013</i> .	<i>arXiv:2504.05419</i> .	785
730	Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu,	A Prompt	786
731	and Yiqun Liu. 2024a. Dragin: dynamic retrieval	In this work, distinct prompts were utilized to	787
732	augmented generation based on the information	guide LLMs in reasoning for different tasks. In	788
733	needs of large language models. <i>arXiv preprint</i>	the question-answering scenario, we classify the	789
734	<i>arXiv:2403.10081</i> .	prompts into two types:	790
735	Weihang Su, Yichen Tang, Qingyao Ai, Junxi Yan,	(1) As shown on the left side of Figure 4, this	791
736	Changyue Wang, Hongning Wang, Ziyi Ye, Yujia	prompt guides the model to directly answer ques-	792
737	Zhou, and Yiqun Liu. 2025. Parametric retrieval	tion using its parametric knowledge. When con-	793
738	augmented generation. In <i>Proceedings of the 48th</i>	structing the dataset for the Confidence Detection	794
739	<i>International ACM SIGIR Conference on Research</i>	Model E, this prompt was consistently used to	795
740	<i>and Development in Information Retrieval</i> , pages	guide the Llama3-8B-Instruct model in generat-	796
741	1240–1250.	ing answers. Additionally, this prompt is applied	797
742	Weihang Su, Changyue Wang, Qingyao Ai, Yiran Hu,	in scenarios within the CBDR system where the	798
743	Zhijing Wu, Yujia Zhou, and Yiqun Liu. 2024b. Un-	LLM has high confidence in answering this ques-	799
744	supervised real-time hallucination detection based on	tion and thus skips the retrieval step to provide a	800
745	the internal states of large language models. <i>arXiv</i>	direct response.	801
746	<i>preprint arXiv:2403.06448</i> .	(2) As shown on the right side of Figure 4,	802
747	Xin Sun, Jianan Xie, Zhongqi Chen, Qiang Liu, Shu Wu,	this prompt guides the LLM to answer question	803
748	Yuehe Chen, Bowen Song, Weiqiang Wang, Zilei	by combining external knowledge with its para-	804
749	Wang, and Liang Wang. 2025. Divide-then-align:	metric knowledge. When constructing the prefer-	805
750	Honest alignment based on the knowledge boundary	ence dataset NQ_Rerank, this prompt was used to	806
751	of rag. <i>arXiv preprint arXiv:2505.20871</i> .	guide Llama3-8B-Instruct to generate the first To-	807
752	Qwen Team. 2024. Qwen2.5: A party of foundation	ken based on the provided context; meanwhile, this	808
753	models.	prompt is also utilized in the CBDR system when	809
754	Yongxin Xu, Ruizhe Zhang, Xinke Jiang, Yujie Feng,	the LLM needs to answer question with reference	810
755	Yuzhen Xiao, Xinyu Ma, Runchuan Zhu, Xu Chu,	to retrieved documents.	811
756	Junfeng Zhao, and Yasha Wang. 2025. Parenting: Op-	B Implementation Details.	812
757	timizing knowledge selection of retrieval-augmented	During the training of the Confidence detection	813
758	language models with parameter decoupling and tai-	model E, the initial learning rate was set to $5e^{-5}$,	814
759	lored tuning. In <i>Proceedings of the 63rd Annual</i>		
760	<i>Meeting of the Association for Computational Lin-</i>		
761	<i>guistics (Volume 1: Long Papers)</i> , pages 11643–		
762	11662.		

QA Prompt	RAG Prompt
<p>You need to read the question carefully and answer it based on your own knowledge.</p> <p>Question: {question}</p>	<p>You are a rigorous language model. Please answer the question based on the provided context. If the context does not support reasoning about the answer, please answer the question based on your own knowledge.</p> <p>Contexts: {contexts}</p> <p>Question: {question}</p>

Figure 4: The prompt on the left side of the figure guides the LLM to answer question using its parametric knowledge; the prompt on the right side requires the LLM to answer question by combining external knowledge with its parametric knowledge.

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.

C LLM Confidence

The training of the Confidence Detection Model E fully follows the approach described in (Ni et al., 2025), and the downstream LLM employed is Llama3-8B-Instruct. After obtaining Model E, we conducted an initial verification to examine how the LLM’s confidence in answering questions varies when provided with relevant versus irrelevant documents. As illustrated in Figure 5, the observed changes in confidence confirm our hypothesis: changes in the LLM’s internal hidden states can guide the selection of external knowledge. It should be noted that in this work, external knowledge refers specifically to contexts obtained through retrieval.

D Preference Dataset

In this work, to align the preferences of the Reranker with the target LLM, we constructed the preference dataset NQ_Rerank using changes in the LLM’s confidence. The NQ_Rerank dataset is divided into a training set with 7,622 items and an evaluation set with 1,216 items. As shown in Figure 6, each data item contains four fields: query, pos, neg, and prompt. Among these, pos and neg are lists of positive and negative contexts, respectively. Specifically, each context in pos enhances the LLM’s confidence, whereas each context in neg

reduces the LLM’s confidence. The prompt field refers to the default prompt used by the Reranker model for the reranking task.

E Dynamic Retrieval

When CBDR is adopted for dynamic retrieval in the RAG system (The algorithm can be found in Algorithm 1), the system achieves the optimal performance at $\beta = 0.98$. As illustrated in Figure 7, this figure comprehensively depicts the impact of the threshold β (ranging from 0.90 to 1.00) on both RAG system performance (quantified by the system score on the left y-axis) and retrieval efficiency (measured by the skip retrieval ratio on the right y-axis). To further elucidate the performance characteristics around this optimal threshold, we analyzed the differential responses corresponding to two threshold intervals: β ranging from 0.95 to 0.98 and from 0.98 to 1.00. Specifically, these differential responses refer to the questions where the integration of external knowledge (enabled by dynamic retrieval) alters the correctness of the LLM’s final answers.

Our observations of these data reveal the following: When the introduction of external knowledge enables the LLM to change from answering incorrectly to correctly, this is always because the introduced external knowledge expands the knowledge boundary of the RAG system, thereby leading the LLM to generate correct answers. However, when the introduction of external knowledge instead causes the LLM to shift from answering correctly to incorrectly, there are multiple reasons for this: (1) As shown in Figure 8, the most common reason is the introduction of incorrect external knowledge, which causes conflicts between the model’s internal parametric knowledge and the ex-

Question	Context_type	Context	Confidence
when was the last time anyone was on the moon	No		0.7454
	Relevance	The Apollo program by NASA included the last human Moon landing during Apollo 17 in December 1972. Astronauts Eugene Cernan and Harrison Schmitt landed on the lunar surface on 14 December 1972 UTC, conducting three days of exploration. No subsequent human missions have reached the Moon since then, making this the final	0.8562 ↑
	Irrelevance	In the video game Grand Theft Auto: San Andreas, the 'Hot Coffee' mod is an unauthorized user modification that accesses a hidden mini-game featuring sexual interactions. This content was originally inaccessible in the official release but was discovered in the game's code, leading to controversy and a re-rating of the game by the ESRB.	0.6293 ↓
when did the eagles win last super bowl	No		0.9288
	Relevance	The Philadelphia Eagles of the NFL last won the Super Bowl in February 2018, which corresponded to the 2017 league season. They defeated the New England Patriots 41-33 in Super Bowl LII, securing their first championship since 1960. As of the current time in 2025, this remains their most recent Super Bowl victory.	0.9462 ↑
	Irrelevance	In the video game Grand Theft Auto: San Andreas, the 'Hot Coffee' mod is an unauthorized user modification that accesses a hidden mini-game featuring sexual interactions. This content was originally inaccessible in the official release but was discovered in the game's code, leading to controversy and a re-rating of the game by the ESRB.	0.4606 ↓
how many seasons of the bastard executioner are there	No		0.5645
	Relevance	FX's historical drama series 'The Bastard Executioner,' created by Kurt Sutter, premiered in September 2015 but received low ratings and mixed reviews. Consequently, the network canceled it after the first season's conclusion in November 2015, with no renewal for additional seasons.	0.9410 ↑
	Irrelevance	Justin Bieber's song 'Love Yourself,' released in 2015, features lyrics co-written with Ed Sheeran that are interpreted as addressing an ex-partner. Although unconfirmed directly by Bieber, widespread media reports and fan speculation suggest it references his past relationship with singer Rihanna, contributing to the song's narrative.	0.5078 ↓

Figure 5: This figure presents three examples of confidence changes. For each example, the confidence levels indicated by the LLM’s internal hidden states are provided under three scenarios, namely: without external documents provided, with relevant documents provided, and with irrelevant documents provided.

Query	who does sam neil play in peter rabbit
Pos	["Peter Rabbit is a 2018 live-action/computer-animated comedy film directed by Will Gluck and written by Rob Lieber and Gluck, based on the stories of Peter Rabbit created by Beatrix Potter. ...", "... The accusations focused on a scene where Thomas McGregor \u2014 whose character has a known severe allergy to blackberries \u2014 is pelted with the berries until one enters his mouth, causing him to enter anaphylactic shock and grab for his EpiPen.[35][36][37] ...", "... A local toy shop on Compston Road, Ambleside, was adapted to be Mr McGregor's.[citation needed]"]
Neg	["The film was first revealed in April 2015 through email leaks as a result of the Sony Pictures hack.[8] The official announcement of the film came that December.[9]", "The Singing Sparrows were voiced by Jessica Freedman, Shana Halligan, Katharine Hoye, Chris Mann, Chad Reisser, and Fletcher Sheridan", "Peter feels bad for what he has done, and upon learning that Bea intends to leave the neighborhood, he and Benjamin head to London to find Thomas at Harrods...", "Thomas and Peter start a war with each other by setting up traps and other offensive nuisance..."]
Prompt	Given a question, retrieve Wikipedia passages that answer the question.

Figure 6: This figure presents an example of the preference dataset NQ_Rerank; each data item contains four fields: Query, Pos, Neg, and Prompt.

886 ternal knowledge, thereby triggering hallucinations.
887 (2) As shown in Figure 9, the introduced external
888 knowledge contains correct documents, but the
889 model exhibits attention bias and fails to focus on

these correct documents. (3) As shown in Figure 890
891 10, the questions are time-sensitive. This issue
892 arises due to the inherent temporal limitations of
893 the NQ dataset; thus, the NQ dataset provides out-

dated external documents and reference answers. As a result, the LLM would originally answer correctly, but ends up answering incorrectly due to the erroneous external knowledge.

F Cross-Model Preference Comparison Experiment

This experiment aims to quantify the preference differences between Llama3-8B-Instruct (the target model) and Qwen2.5-7B-Instruct (the comparison model) toward different types of retrieved contexts. It verifies the strong binding characteristic between the filtering logic of the preference-aligned reranker and the target LLM at the cognitive mechanism level, providing empirical support for the core conclusion in the main text that the reranker’s effectiveness is model-dependent.

F.1 Target LLMs

1) **Llama3-8B-Instruct**: Fine-tuning target of the preference-aligned reranker and baseline model in the main context experiments; 2) **Qwen2.5-7B-Instruct**: An open-source LLM of the same scale, used to verify the generality of preference differences.

F.2 Retrieved Context Corpus

The corpus is randomly selected and constructed from candidate retrieved contexts in the NQ dataset, consisting of 100 sample sets. Each set corresponds to one unified question and includes 3 types of retrieved contexts (1 context per type, generated by rewriting gold contexts), ensuring a single controllable experimental variable: 1) **Type A**: Highly matches the core semantics of the question but only mentions core concepts without specific data, logical chains, or supplementary details; 2) **Type B**: Contains complete factual information (data, logic, etc.) required to answer the question but has slightly lower semantic matching with the question’s expression (consistent core concepts, but expressed through synonymous substitution and sentence structure reconstruction); 3) **Type C**: Extremely low semantic relevance to the question and no question-related factual information, serving only as a control benchmark for preference judgment.

F.3 Controlled Variables

To avoid interference from irrelevant factors, the following variables are strictly controlled for all contexts: 1) **Context Length**: Uniformly limited

to within 200 characters; 2) **Format Standardization**: Pure context format without special symbols, lists, formulas, etc.; 3) **Language Style**: Formal written language, avoiding colloquial or emotional expressions; 4) **Question Consistency**: The 3 context types in each sample set correspond to the same question, ensuring that preference differences arise only from context types rather than the question itself.

F.4 Experiment

This experiment adopts a baseline-test control design to quantify model preferences through confidence changes. The specific steps are as follows:

Sample Preprocessing: Standardize the 100 sets of questions and their corresponding Type A, B, and C contexts: Use Prompts to guide LLMs in generating contexts based on the questions and Gold context; the Type A prompts are shown in Figure 11.

Confidence Collection: Adopt a dual-scenario comparison of "context-free baseline" and "context-augmented test" to collect the generation confidence of the two models: Context-Free Baseline (Base_Conf): Input the 100 standardized questions individually into Llama3-8B-Instruct and Qwen2.5-7B-Instruct. The models generate answers relying solely on internal parameter knowledge, and the confidence during generation is recorded (calculated as the maximum value of the softmax outputs of model logits, with a value range of [0,1]). Context-Augmented Test (Test_Conf): For each sample set, construct three input formats—"question + Type A context", "question + Type B context", and "question + Type C context"—and input them into the two LLMs sequentially. After generating answers, the corresponding confidence is recorded (using the same calculation method as the baseline).

Data Calculation and Statistics: Conduct quantitative analysis on the collected confidence data, with the core calculation logic as follows: Per-Sample Confidence Improvement Value: For each model, sample set, and text type, calculate Improvement Value = Test_Conf - Base_Conf. Positive value indicates the context enhances model confidence, i.e., "preference"; Negative value indicates the context reduces model confidence, i.e., "rejection". Average Confidence Improvement Value: For each model and context type, compute the total sum of improvement values across the 100 sample sets.

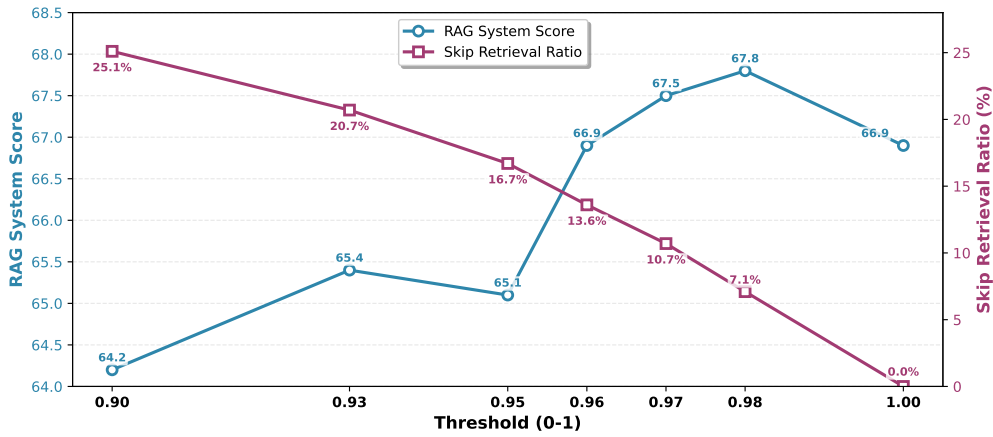


Figure 7: The Impact of Threshold (0–1) on RAG System Performance and Retrieval Efficiency under the Top-3 Setting. X-axis: Threshold (0.90–1.00); Left Y-axis: RAG System Score; Right Y-axis: Skip Retrieval Ratio (%).

Question	Context	Retrieval	Reference	Answer
when did brazil last host the fifa world cup	["The 1951 World Cup hosting duty was decided on 26 July 1946 , the same day that Brazil was selected for the 1949 World Cup, in Luxembourg City. On 27 July ...",	No	2014	According to my knowledge, Brazil last hosted the FIFA World Cup in 2014.
	"Brazil made the first unopposed bid since the initial selection of the 1986 FIFA World Cup (when Colombia was selected as host, but later withdrew for financial problems) ... "Only Mexico, Italy, France, Germany (West Germany until shortly after the 1990 World Cup) and Brazil have hosted the event on two occasions ..."]	Top-3		Based on the provided context, it can be inferred that Brazil last hosted the FIFA World Cup in 1950 .

Figure 8: Since all the introduced external knowledge is incorrect, it causes the LLM’s answer to shift from correct to incorrect.

Question	Context	Retrieval	Reference	Answer
when was the art institute of chicago founded	["The Art Institute of Chicago, founded in 1879 and located in Chicago's Grant Park, is one of the oldest and largest art museums in the United States ...",	No	1879	The Art Institute of Chicago was founded in 1879 ...
	"In 1866 , a group of 35 artists founded the Chicago Academy of Design in a studio on Dearborn Street, with the intent to run a free school with its own art gallery ... The Academy's charter was granted in March 1867 ." "In 1882 , the Chicago Academy of Fine Arts changed its name to the current Art Institute of Chicago and elected as its first president the banker and philanthropist Charles L ..."]	Top-3		Based on the provided context, the Art Institute of Chicago was founded in 1866 as the Chicago Academy of Design ...

Figure 9: Although the external knowledge contains correct context, the LLM focuses on incorrect contexts, resulting in its answer shifting from correct to incorrect.

G Detailed Analysis of Resource Consumption for the CBR Framework

The additional resource consumption required to construct the CBR framework is divided into the offline preparation phase and online inference phase, with overall computational costs well-controlled. Details are as follows:

G.1 Offline Preparation Phase

The offline phase involves four core tasks, with resource consumption concentrated on target LLM forward computations and lightweight model training (no generation operations), ensuring manageable computational costs: 1) **Constructing the NQ_Confidence dataset**: Provides training data for the confidence detection model E, including

993
994
995
996
997
998
999
1000

1001
1002
1003
1004
1005
1006
1007
1008

Question	Context	Retrieval	Reference	Answer
who was the most recent president to be impeached	["... Andrew Johnson and Bill Clinton are the only two presidents to have been successfully impeached ...",	No	Bill Clinton	The most recent president to be impeached was Donald Trump , the 45th President of the United States ...
	"... President Richard Nixon resigned from office after the House Judiciary Committee had already reported articles of impeachment to the floor.", "... especially of presidents, are common,[22][23] going back to the administration of George Washington in the mid-1790s ..."]	Top-3		... the most recent president to be impeached was Bill Clinton , who was successfully impeached by the House of Representatives in 1998 ...

Figure 10: Since the questions are time-sensitive, the contexts and reference provided by the NQ dataset are all incorrect; the retrieved external knowledge thus causes the LLM to answer incorrectly.

Type A Prompt
<p>Based on the following Query and corresponding key factual text (Gold content), generate an English auxiliary response text. Requirements: 1) Fully retain high semantic relevance to the query without deviating from the core topic; 2) Only keep core concepts and conclusions, and delete all specific data, logical reasoning processes, and case details; 3) Control the length to approximately 1/3 of the Gold content with concise language; 4) Return the result in JSON format.</p> <p>Example: Query: "What is the approximate length of the Amazon River in kilometers?" Content: "The Amazon River is the second longest river in the world, with a total length of its main stream about 6,400 kilometers, a drainage basin area of 7.05 million square kilometers, and flowing through 9 countries including Brazil and Peru." Output: ... json { "result": "The Amazon River is one of the world's major long rivers, with its main stream totaling approximately several thousand kilometers and flowing through multiple South American countries." } ... Input: Query: {question} Gold content: {content} Output:</p>

Figure 11: Type A prompt is used to generate content that is semantically relevant to the gold context but with concise factual details.

1009 input contexts and corresponding target LLM confidence labels. Sample scale: 2,000 training samples,
1010 1,000 development (Dev) samples, and 500 test samples, totaling 3,500 samples. Computational load: Each sample requires one forward pass
1011 of the target LLM (e.g., Llama3-8B-Instruct) to extract confidence features, resulting in a total of 3,500 target LLM forward passes. 2) **Training**
1012 **the confidence detection model E:** Model architecture: 5-layer MLP with 2M parameters. Training configuration: Trained on the NQ_Confidence
1013 training set for 100 epochs with a batch size of 32. Resource characteristics: Lightweight model with extremely low training overhead; completable
1014 within 10 minutes on a single consumer-grade GPU (e.g., NVIDIA RTX 4090). 3) **Constructing the NQ_Rerank fine-tuning dataset:** Pro-

vides fine-tuning data for the preference-aligned reranker, with labels generated based on confidence changes of the target LLM toward different
1026 retrieved contexts. Sample scale: 7,622 training samples and 1,216 test samples, totaling 8,838 samples. Computational load: Each sample requires
1027 8 forward passes of the target LLM (matching 8 types of retrieved contexts) to generate preference labels, resulting in a total of 70,704 target LLM
1028 forward passes. 4) **Fine-tuning the reranker:** Base model: bge-reranker-v2-m3 (568M parameters). Fine-tuning configuration: Fine-tuned on
1029 the NQ_Rerank training set for 5 epochs (weights after 1 epoch are used in practice) with a batch size of 32. Resource characteristics: Moderate-scale
1030 model with few fine-tuning epochs; the entire process can be completed within 2 hours on a single
1031 GPU. Sample scale: 7,622 training samples and 1,216 test samples, totaling 8,838 samples. Computational load: Each sample requires
1032 8 forward passes of the target LLM (matching 8 types of retrieved contexts) to generate preference labels, resulting in a total of 70,704 target LLM
1033 forward passes. 4) **Fine-tuning the reranker:** Base model: bge-reranker-v2-m3 (568M parameters). Fine-tuning configuration: Fine-tuned on
1034 the NQ_Rerank training set for 5 epochs (weights after 1 epoch are used in practice) with a batch size of 32. Resource characteristics: Moderate-scale
1035 model with few fine-tuning epochs; the entire process can be completed within 2 hours on a single
1036 GPU. Sample scale: 7,622 training samples and 1,216 test samples, totaling 8,838 samples. Computational load: Each sample requires
1037 8 forward passes of the target LLM (matching 8 types of retrieved contexts) to generate preference labels, resulting in a total of 70,704 target LLM
1038 forward passes. 4) **Fine-tuning the reranker:** Base model: bge-reranker-v2-m3 (568M parameters). Fine-tuning configuration: Fine-tuned on
1039 the NQ_Rerank training set for 5 epochs (weights after 1 epoch are used in practice) with a batch size of 32. Resource characteristics: Moderate-scale
1040 model with few fine-tuning epochs; the entire process can be completed within 2 hours on a single
1041 GPU. Sample scale: 7,622 training samples and 1,216 test samples, totaling 8,838 samples. Computational load: Each sample requires
1042 8 forward passes of the target LLM (matching 8 types of retrieved contexts) to generate preference labels, resulting in a total of 70,704 target LLM
1043 forward passes. 4) **Fine-tuning the reranker:** Base model: bge-reranker-v2-m3 (568M parameters). Fine-tuning configuration: Fine-tuned on
1044 the NQ_Rerank training set for 5 epochs (weights after 1 epoch are used in practice) with a batch size of 32. Resource characteristics: Moderate-scale
1045 model with few fine-tuning epochs; the entire process can be completed within 2 hours on a single GPU.

Algorithm 1 CBDR Inference

Input: Query q , Target LLM M , confidence detection model E , retriever \mathcal{R} , document corpus \mathcal{D} , Fine-tuned reranker FR , confidence threshold β , Top- K parameter K , QA prompt \mathcal{P}_{qa} , RAG prompt \mathcal{P}_{rag}

Output: Final answer y

- 1: Construct pure QA prompt: $\text{prompt}_{qa} \leftarrow \mathcal{P}_{qa}(q)$
 - 2: Tokenize input: $\mathbf{x} \leftarrow \text{Tokenize}(\text{prompt}_{qa})$
 - 3: Forward through M with hidden states output: $\text{outputs} \leftarrow M(\mathbf{x}, \text{output_hidden_states} = \text{True})$
 - 4: Extract last context hidden state: $\mathbf{V} \leftarrow \text{outputs.hidden_states}[-1][0, -1, :]$
 - 5: Compute confidence score: $S \leftarrow E(\mathbf{V})$
 - 6: **if** $S \geq \beta$ **then**
 - 7: Generate answer directly: $y \leftarrow M.\text{generate}(\mathbf{x}, \text{max_new_tokens} = L_{\text{max}})$
 - 8: **else**
 - 9: Retrieve candidate documents: $\mathcal{D}_{\text{raw}} \leftarrow \mathcal{R}.\text{retrieve}(q, \text{top_k} = K \cdot r)$
 - 10: Rerank using M -aligned reranker: $\mathcal{D}_{\text{reranked}} \leftarrow FR.\text{rerank}(q, \mathcal{D}_{\text{raw}}, K)$
 - 11: Construct context: $\text{context} \leftarrow \bigoplus_{i=1}^K \mathcal{D}_{\text{reranked}}[i].\text{text}$
 - 12: Build RAG prompt: $\text{prompt}_{rag} \leftarrow \mathcal{P}_{rag}(\text{context}, q)$
 - 13: Tokenize RAG input: $\mathbf{x}_{rag} \leftarrow \text{Tokenize}(\text{prompt}_{rag})$
 - 14: Generate answer with evidence: $y \leftarrow M.\text{generate}(\mathbf{x}_{rag}, \text{max_new_tokens} = L_{\text{max}})$
 - 15: **end if**
 - 16: Return y
-

1043 NVIDIA RTX 4090.

1044 G.2 Online Inference Phase

1045 The additional resource consumption in the online
1046 phase only involves two lightweight forward passes,
1047 which do not affect inference efficiency: 1) **Work-**
1048 **flow:** For input query, the target LLM first performs
1049 one forward pass to obtain hidden states, while the
1050 confidence detection model E (2M parameters) ex-
1051 ecutes one forward pass to determine the necessity
1052 of retrieval. If retrieval is deemed unnecessary, the
1053 generation process can be directly connected using
1054 the results of this target LLM forward pass,
1055 eliminating the need for additional LLM calls and
1056 redundant computations. 2) **Optimization:** Re-
1057 trieval knowledge is directly injected into the target
1058 LLM’s Attention module (Dong et al., 2025b) or
1059 FNN module (Su et al., 2025) via parameter injec-
1060 tion. Thus, even if retrieval is required after the
1061 LLM’s forward pass, the generation process can
1062 be seamlessly connected following parameterized
1063 knowledge injection, with no redundant compu-
1064 tational overhead. 3) **Latency:** The total latency
1065 of the two forward passes is less than 100ms (on
1066 a single NVIDIA RTX 4090), which negligibly
1067 increases inference latency and meets real-time re-
1068 quirements.