

QuCo-RAG: Quantifying Uncertainty from the Pre-training Corpus for Dynamic Retrieval-Augmented Generation

Dehai Min¹, Kailin Zhang², Tongtong Wu³, Lu Cheng¹

¹University of Illinois at Chicago, ²New York University, ³Monash University

dmin10@uic.edu, kz2739@nyu.edu, tongtong.wu@monash.edu, lucheng@uic.edu

Abstract

Dynamic Retrieval-Augmented Generation adaptively determines when to retrieve during generation to mitigate hallucinations in large language models (LLMs). However, existing methods rely on model-internal signals (e.g., logits, entropy), which are fundamentally unreliable because LLMs are typically ill-calibrated and often exhibit high confidence in erroneous outputs. We propose QuCo-RAG, which shifts from **subjective** confidence to **objective** statistics computed from pre-training data. Our method quantifies uncertainty through two stages: (1) before generation, we identify low-frequency entities indicating long-tail knowledge gaps; (2) during generation, we verify entity co-occurrence in the pre-training corpus, where zero co-occurrence often signals hallucination risk. Both stages leverage Infini-gram for millisecond-latency queries over 4 trillion tokens, triggering retrieval when uncertainty is high. Experiments on multi-hop QA benchmarks show QuCo-RAG achieves EM gains of 5–12 points over state-of-the-art baselines with OLMo-2 models, and transfers effectively to models with undisclosed pre-training data (Llama, Qwen, GPT), improving EM by up to 14 points. Domain generalization on biomedical QA further validates the robustness of our paradigm. These results establish corpus-grounded verification as a principled, practically model-agnostic paradigm for dynamic RAG¹.

1 Introduction

Retrieval-Augmented Generation (RAG) (Lewis et al., 2020; Gao et al., 2023b) mitigates LLM hallucinations by grounding generation in external evidence. Early RAG systems employ static strategies with a single retrieval step before generation (Karpukhin et al., 2020; Shi et al., 2024; Min et al., 2025), but fall short for complex multi-step

¹Our code is publicly available at <https://github.com/ZhishanQ/QuCo-RAG>.

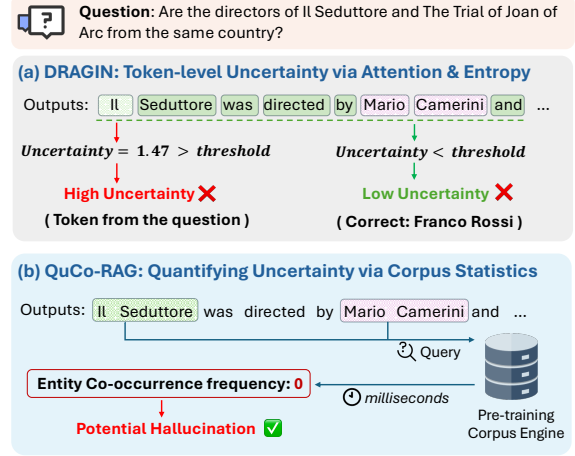


Figure 1: Comparison of retrieval triggering mechanisms. (a) DRAGIN relies on model-internal signals, incorrectly assigning high uncertainty to “Il” (a token from the question) while showing low uncertainty on the hallucinated director name. (b) QuCo-RAG correctly detects the hallucination through zero entity co-occurrence in the pre-training corpus.

tasks where information needs emerge dynamically during generation (Su et al., 2025; Wang et al., 2025, 2023). This has driven the emergence of Dynamic RAG methods that adaptively determine when and what to retrieve based on the generation process (Jiang et al., 2023; Asai et al., 2024).

Current dynamic RAG methods predominantly rely on quantifying uncertainty through model-internal signals such as token probability (Jiang et al., 2023) or entropy (Su et al., 2024; Li et al., 2025a). However, these methods assume internal signals reliably indicate generation correctness—an assumption that is fundamentally flawed (Li et al., 2024b). As illustrated in Figure 1(a), the notable work DRAGIN (Su et al., 2024) exhibits low uncertainty when generating the incorrect director name “Mario Camerini”, yet assigns high uncertainty to “Il”—a token from the question. This failure reflects a well-documented problem: LLMs are poorly calibrated (Guo et al., 2017; Kadavath

et al., 2022; Achiam et al., 2023)—their confidence scores fail to correlate with actual prediction accuracy. This miscalibration leads to “confident hallucinations,” where models produce incorrect content with high confidence (Tian et al., 2023). Furthermore, post-training techniques such as SFT (Dong et al., 2024) and Reinforcement Learning (Ouyang et al., 2022; Guo et al., 2025) often exacerbate this by encouraging decisive answers. More fundamentally, recent theoretical work (Kalai and Vempala, 2024) further shows that for rarely-seen facts, even perfectly calibrated models must hallucinate to maintain statistical consistency.

To bypass the limitations, we propose **QuCo-RAG**, a framework that determines when to retrieve by **Quantifying** uncertainty via pre-training **Corpus** statistics, shifting from subjective internal confidence to objective external evidence. Our key insight is that an LLM’s factual knowledge is fundamentally shaped by its pre-training corpus (Balepur et al., 2025): low-frequency entities correspond to long-tail knowledge that models struggle to memorize reliably, while zero co-occurrence between entity pairs indicates the model has no evidential basis for claims relating them. Based on this insight, QuCo-RAG operates through two-stage detection: **(1) Pre-Generation Knowledge Assessment:** We query entity frequencies in the pre-training corpus, triggering retrieval when entities are low-frequency (long-tail knowledge risks). **(2) Runtime Claim Verification:** We extract knowledge triplets from each generated sentence and verify entity co-occurrence; zero co-occurrence triggers retrieval and regeneration. Both stages leverage Infini-gram (Liu et al., 2024) for millisecond-latency queries over trillion-token corpora.

To validate our approach, we first evaluate QuCo-RAG on multi-hop QA benchmarks using the OLMo-2 model family (7B, 13B, 32B) (OLMo et al., 2024), which provides full access to its 4-trillion token pre-training corpus for precise statistical verification. Results show QuCo-RAG achieves 5–12 point improvements on Exact Match (EM) over state-of-the-art baselines across all model scales, while maintaining competitive efficiency.

Beyond this matched-corpus setting, we demonstrate QuCo-RAG’s broad applicability through two additional dimensions of evaluation. First, for **cross-model transferability**, we show that corpus statistics computed from OLMo-2’s pre-training corpus serve as effective proxies for models with undisclosed training data. Leveraging the

substantial overlap of web-scale pre-training corpora, QuCo-RAG yields up to 14 EM improvements on Llama-3, Qwen2.5, and GPT-4.1/5 series. Second, for **domain generalization**, we evaluate on PubMedQA (Jin et al., 2019), a biomedical QA benchmark requiring specialized knowledge. QuCo-RAG achieves the best accuracy while internal-signal methods either trigger excessive retrievals or fail to improve over no-retrieval baselines, demonstrating that our framework generalizes robustly without domain-specific tuning.

2 Related Work

Dynamic Retrieval-Augmented LLM Dynamic RAG methods have evolved to address the limitations of static retrieval approaches by adaptively determining when and what to retrieve during generation (Xu et al., 2024; Yu et al., 2024; Yang et al., 2025). FLARE (Jiang et al., 2023) pioneered this direction by triggering retrieval when encountering low-probability tokens. Self-RAG (Asai et al., 2024) extended this paradigm by training models to generate special reflection tokens that assess retrieval necessity and response quality, though requiring additional fine-tuning. More recent approaches (Ma et al., 2025) construct more sophisticated uncertainty metrics: DRAGIN (Su et al., 2024) integrates multiple model-internal signals including entropy and attention weights, ETC (Li et al., 2025a) considers first- and second-order entropy differences to capture uncertainty trends, and SeaKR (Yao et al., 2025) extracts self-aware uncertainty from LLMs’ internal FFN states. However, these methods all rely on model-internal signals, which may not reliably indicate correctness.

Reusing LLM Pre-Training Data at Inference Time Recent work explores unlocking additional value from pre-training corpora at inference time. Fang et al. (2025) showed that retrieving from the model’s own pre-training data yields performance gains equivalent to a $5\times$ increase in pre-training compute. Efficient infrastructure has emerged to support trillion-scale corpus access. Infini-gram (Liu et al., 2024) provides millisecond-latency n -gram counting via suffix arrays, while Infini-gram mini (Xu et al., 2025) reduces index size to 44% of the corpus via FM-index (Ferragina and Manzini, 2000). OLMoTrace (Liu et al., 2025) enables real-time tracing of LLM output back to verbatim matches in training documents. Our work leverages this infrastructure for a distinct purpose:

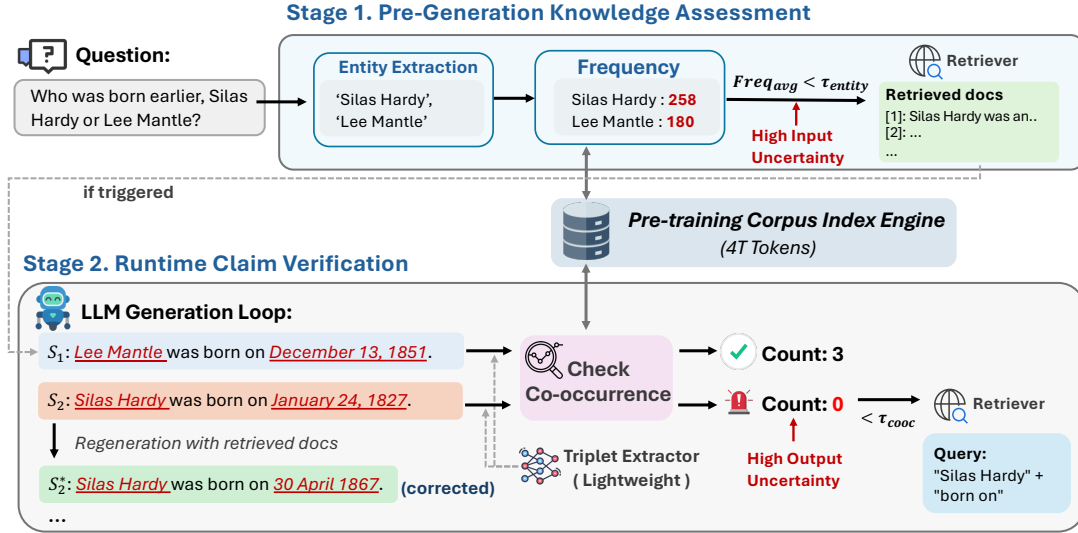


Figure 2: Overview of QuCo-RAG Framework.

using pre-training corpus statistics to *quantify uncertainty and trigger retrieval*, enabling reliable hallucination detection and mitigation.

3 Methodology

3.1 Problem Formulation

We formalize the dynamic RAG problem as follows. Let \mathcal{M} denote an LLM, \mathcal{C} represent an external knowledge base for retrieval (e.g., Wikipedia), and \mathcal{P} denote the pre-training corpus used to train \mathcal{M} . Given an input question Q , the model generates a response $y = (s_1, s_2, \dots, s_N)$, where s_i is the i -th generated sentence. A dynamic RAG system makes two critical decisions during generation:

(1) **When to retrieve.** At each step i , determine whether to trigger retrieval:

$$\delta_i = f_{\text{trigger}}(Q, s_{<i}; \Theta) \in \{0, 1\}, \quad (1)$$

where Θ denotes the source of uncertainty signals. Unlike prior methods that rely on internal model states (i.e., $\Theta = \mathcal{M}$), we ground the decision in pre-training corpus statistics (i.e., $\Theta = \mathcal{P}$).

(2) **What to retrieve.** When $\delta_i = 1$, construct a query $q_i = f_{\text{query}}(Q, s_{<i})$ and retrieve related documents $\mathcal{D}_i = \text{Retrieve}(q_i, \mathcal{C})$, where f_{query} is the query formulation function.

Binary Nature of Retrieval Decisions. Note that the retrieval decision $\delta_i \in \{0, 1\}$ is inherently binary: the system either retrieves or not. This observation motivates our design: rather than estimating *continuous confidence scores* from model-internal signals to infer uncertainty, whose thresholds lack clear semantic grounding, we directly leverage *discrete corpus statistics* to determine whether the

model faces high uncertainty (retrieve) or low uncertainty (proceed without retrieval). Specifically, we consider two high-uncertainty scenarios: (1) **Input uncertainty**: the question contains entities rarely seen during pre-training, indicating insufficient knowledge coverage; (2) **Output uncertainty**: the generated claim relates entities that never co-occur in the corpus, indicating lack of evidential support. Both signals are grounded in corpus statistics, as illustrated in Figure 2.

3.2 Pre-Generation Knowledge Assessment

To quantify input uncertainty, we employ a pre-check mechanism before generation begins. We first use a lightweight entity extractor to identify a set of key entities $\mathcal{E}_Q = \{e_1, e_2, \dots, e_m\}$ from the input question Q . For each entity $e \in \mathcal{E}_Q$, we query its frequency in the pre-training corpus \mathcal{P} , denoted as $\text{freq}(e; \mathcal{P})$. We posit that entities with low frequency in \mathcal{P} represent long-tail knowledge risks, where the model is likely to hallucinate. Retrieval is triggered if the average entity frequency falls below a predefined threshold:

$$\delta_{\text{pre}} = \mathbb{I}(\text{Avg}_{e \in \mathcal{E}_Q} \text{freq}(e; \mathcal{P}) < \tau_{\text{entity}}). \quad (2)$$

We set $\tau_{\text{entity}} = 10^3$ as the default threshold; results remain stable across a wide range (10^3 to 10^7) as shown in Appendix A.2. If $\delta_{\text{pre}} = 1$, we use the original question Q as the search query to retrieve relevant documents \mathcal{D}_0 , which are prepended to the model’s context before generation starts.

3.3 Runtime Claim Verification

To quantify output uncertainty, QuCo-RAG continuously monitors each generated sentence s_i by

verifying whether the claimed facts have evidential support in the pre-training corpus. For a generated sentence s_i , we extract a set of knowledge triplets $\mathcal{T} = \{(h, r, t)\}$, where h, r, t represent the head entity, relation, and tail entity, respectively. We quantify the evidential support for each triplet by computing the co-occurrence frequency of the head and tail entities within a defined window ω (e.g., a document or paragraph) in \mathcal{P} :

$$\text{cooc}(h, t; \mathcal{P}) = |\{\omega \in \mathcal{P} : h \in \omega \wedge t \in \omega\}|. \quad (3)$$

We compute $\text{cooc}(h, t)$ rather than $\text{cooc}(h, r, t)$ because relational predicates exhibit high lexical variability (e.g., “employed by” vs. “worked at”), while named entities are more lexically stable (Galárraga et al., 2014). Retrieval is triggered if the co-occurrence count falls below a threshold τ_{cooc} (default set to 1):

$$\delta_i = \mathbb{I} \left(\min_{(h, r, t) \in \mathcal{T}} \text{cooc}(h, t; \mathcal{P}) < \tau_{\text{cooc}} \right). \quad (4)$$

The rationale for $\tau_{\text{cooc}} = 1$ is intuitive: if two entities never co-occur in the pre-training corpus, the generated claim lacks evidential support and likely constitutes a hallucination (Mallen et al., 2023; Kandpal et al., 2023). Notably, co-occurrence evidence is *asymmetric*: while $\text{cooc}(h, t; \mathcal{P}) > 0$ does not guarantee correctness (entities may co-occur with different relations or in unrelated contexts), $\text{cooc}(h, t) = 0$ strongly indicates hallucination risk (Gao et al., 2023a; Ravichander et al., 2025). When retrieval is triggered ($\delta_i = 1$), we construct a *Semantic-Oriented Query* using the head entity and relation ($q = h \oplus r$) to retrieve supporting documents and regenerate the sentence.

3.4 Implementation Details

Corpus Statistics via Infini-gram. We leverage Infini-gram (Liu et al., 2024), a suffix array-based engine that supports millisecond-latency queries over trillion-token corpora, enabling real-time computation during generation.

Lightweight Triplet Extraction. To minimize overhead while ensuring extraction quality, we distill a specialized 0.5B model from GPT-4o-mini (Hurst et al., 2024). Specifically, we construct 40K annotated examples using in-context learning, then perform full-parameter supervised fine-tuning on Qwen2.5-0.5B-Instruct (Team, 2024). Representative training examples are provided in Appendix A.3.

4 Experimental Setup

4.1 Datasets and Implementation

We evaluate on two widely adopted knowledge-intensive multi-hop QA benchmarks: 2WikiMulti-hopQA (Ho et al., 2020) and HotpotQA (Yang et al., 2018). Following Su et al. (2024), we sample the first 1,000 validation examples from each as our test sets and report Exact Match (EM) and token-level F1 score as evaluation metrics, which are well-suited for these benchmarks as answers are short-form entities that can be reliably extracted and matched. Prior work (Li et al., 2025a) has shown that EM/F1 conclusions align with LLM-as-judge (Li et al., 2025b) evaluations on these datasets. For retrieval, we employ BM25 (Robertson et al., 2009) over the Wikipedia dump from Karpukhin et al. (2020) as our external corpus \mathcal{C} , retrieving top-3 documents per query. We also verify robustness with dense retrievers in Appendix A.4. In our experiments, we query entity frequencies and co-occurrences via the Infini-gram API², which hosts the full OLMo-2 pre-training corpus index. We set the co-occurrence window size to 1,000 tokens, roughly matching passage-level context length. More detailed LLM generation settings and the full prompt template are provided in Appendix A.1. All experiments are conducted on NVIDIA H200 GPUs (141GB HBM3e).

4.2 Baselines

No Retrieval: Wo-RAG generates answers directly without any external retrieval, serving as the lower bound to measure RAG benefits.

Static Retrieval: Single-Round RAG (SR-RAG): performs one-time retrieval using the input question before generation begins. Fixed-Sentence RAG (FS-RAG) (Trivedi et al., 2023) triggers retrieval after every generated sentence, using the last sentence as the query.

Dynamic Retrieval: FLARE (Jiang et al., 2023) triggers retrieval on low-probability tokens. **DRA-GIN** (Su et al., 2024) combines entropy, attention, and semantic signals. **ETC** (Li et al., 2025a) models first- and second-order entropy differences. **SeaKR** (Yao et al., 2025) leverages internal FFN states for uncertainty estimation. All baseline results are reproduced using their released code.

²API Endpoint Documentation: <https://infini-gram.readthedocs.io/en/latest/api.html>. The Infini-gram index supports local deployment for offline environments, requiring primarily CPU and disk storage rather than GPU resources.

Table 1: Performance comparison on multi-hop QA benchmarks across OLMo-2 model scales. **Bold**: best; underline: second-best. *Improv.*: absolute gain over best baseline. 2Wiki: 2WikiMultihopQA.

Method	OLMo-2-7B				OLMo-2-13B				OLMo-2-32B			
	2Wiki		HotpotQA		2Wiki		HotpotQA		2Wiki		HotpotQA	
	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1	EM	F1
Wo-RAG	20.1	26.4	22.6	31.6	28.5	34.5	24.4	33.6	33.3	40.3	22.0	31.3
SR-RAG	23.7	30.7	<u>29.7</u>	<u>40.7</u>	28.9	35.7	<u>29.7</u>	<u>39.5</u>	<u>37.4</u>	<u>46.5</u>	29.5	40.4
FS-RAG	21.1	28.3	14.5	20.7	28.8	35.1	14.6	21.9	34.6	41.0	13.9	19.5
FLARE	22.9	28.9	20.3	28.4	26.2	31.5	15.3	21.9	32.0	39.3	28.3	39.8
DRAGIN	22.8	29.0	17.5	24.7	28.5	33.9	19.5	27.6	33.3	40.2	17.7	24.3
ETC	23.4	29.8	25.1	34.7	<u>29.7</u>	<u>35.9</u>	29.3	<u>39.5</u>	36.0	43.6	30.8	<u>42.2</u>
SeaKR	<u>25.3</u>	<u>32.7</u>	24.8	35.0	29.6	34.6	26.2	37.3	30.2	38.2	28.7	40.4
QuCo-RAG	32.7	41.1	35.3	46.1	41.7	49.1	35.0	46.8	46.8	56.2	41.6	54.2
<i>Improv.</i>	+7.4	+8.4	+5.6	+5.4	+12.0	+13.2	+5.3	+7.3	+9.4	+9.7	+10.8	+12.0

4.3 Models

Primary Models (Matched Corpus). We use the OLMo-2-Instruct family (OLMo et al., 2024) (7B, 13B, and 32B) as our primary evaluation targets. OLMo-2 achieves performance competitive with mainstream models like Qwen2.5 while providing publicly available training data, code, and recipes. The pre-training corpus³ comprises about 4 trillion tokens from diverse sources. This transparency enables precise computation of entity frequencies and co-occurrence statistics, making OLMo-2 ideal for validating our method.

Transferability Models (Proxy Corpus). A key advantage of QuCo-RAG is its applicability to LLMs with undisclosed pre-training data. Given that web-scale pre-training corpora share substantial overlap (Soldaini et al., 2024), statistics derived from a transparent and comprehensive corpus can serve as effective proxies for other models. We demonstrate this by using the OLMo-2 corpus as a proxy for Llama-3-8B-Instruct (Grattafiori et al., 2024), Qwen2.5-32B-Instruct (Team, 2024), and proprietary models (GPT-4.1 (OpenAI, 2025a), GPT-5-chat (OpenAI, 2025b)). For GPT models, we additionally compare against their built-in agentic web search, where the model autonomously invokes web search via the Responses API.

5 Experimental Results

We design experiments to answer three core research questions:

- **RQ1:** How does corpus-based uncertainty compare to model-internal signals? (§5.1)

³<https://huggingface.co/datasets/allenai/olmo-mix-1124>

- **RQ2:** How well does QuCo-RAG transfer to models with undisclosed training data? (§5.2)
- **RQ3:** What is the efficiency-performance trade-off of QuCo-RAG? (§5.3)

5.1 Main Results (RQ1)

Table 1 presents the main results on OLMo-2 models across both benchmarks.

QuCo-RAG Achieves Significant Improvements over Baselines. Across all model scales and datasets, QuCo-RAG consistently outperforms the strongest baselines by significant margins. On OLMo-2-7B, QuCo-RAG achieves 32.7 EM on 2WikiMultihopQA and 35.3 EM on HotpotQA, surpassing the best baseline by +7.4 and +5.6 points respectively. The improvements become even more pronounced with larger models: OLMo-2-13B shows gains of +12.0 EM on 2WikiMultihopQA, while OLMo-2-32B achieves +10.8 EM improvements on HotpotQA. These results demonstrate that grounding retrieval decisions in corpus statistics provides a fundamentally more reliable signal than model-internal uncertainty measures.

Internal-Signal Methods Show Inconsistent Performance. Methods relying on model-internal signals (FLARE, DRAGIN, ETC, SeaKR) show highly variable results across settings. For instance, ETC achieves second-best performance in some configurations, yet underperforms even simple SR-RAG in others. DRAGIN achieves only 17.5–19.5 EM on HotpotQA across all model sizes, substantially underperforming SR-RAG. This inconsistency stems from the fundamental unreliability of internal uncertainty signals. A detailed case study is provided in Appendix A.5.

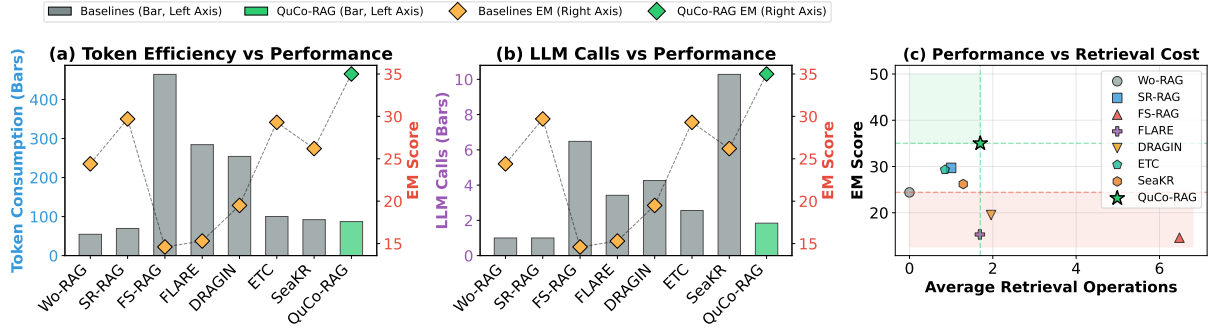


Figure 3: Efficiency-performance trade-off analysis on HotpotQA with OLMo-2-13B-Instruct. (a) EM score versus Token consumption. (b) EM score versus LLM calls. (c) Performance versus Retrieval frequency. QuCo-RAG achieves the highest EM with moderate token usage and LLM calls.

Table 2: Transferability to other model families (EM scores). HPQA: HotpotQA. ‘-’ indicates the method is not applicable due to API limitations. Full results with F1 score are in Appendix A.6.

Method	Qwen2.5-32B		Llama-3-8B	
	2Wiki	HPQA	2Wiki	HPQA
Wo-RAG	26.4	21.6	29.5	20.3
SR-RAG	23.0	31.0	12.9	22.7
FS-RAG	35.9	38.6	28.8	27.0
FLARE	26.4	24.1	26.6	22.2
DRAGIN	28.8	22.2	27.9	20.0
ETC	31.5	21.7	29.9	24.1
SeaKR	22.4	26.7	33.5	33.5
QuCo-RAG	50.0	41.6	38.4	36.2
Improv.	+14.1	+3.0	+4.9	+2.7

Method	GPT-4.1		GPT-5-chat	
	2Wiki	HPQA	2Wiki	HPQA
Wo-RAG	54.7	40.1	50.1	37.7
SR-RAG	60.0	38.8	51.0	42.9
FS-RAG	59.5	25.9	47.3	19.0
FLARE	49.8	38.7	-	-
Web-Tool	42.9	8.9	48.3	19.8
QuCo-RAG	64.6	48.2	59.7	48.4
Improv.	+4.6	+8.1	+8.7	+5.5

5.2 Transferability to Other Models (RQ2)

A critical question for corpus-based methods is whether they generalize to models whose training data is proprietary or undisclosed. We evaluate QuCo-RAG on Qwen2.5, Llama-3, and GPT model families, using the OLMo-2 corpus as a *proxy corpus* for their knowledge distributions (Table 2).

Effectiveness Across Model Families. QuCo-RAG demonstrates remarkable transferability, consistently outperforming all baselines across model families. On open-weight models, it achieves substantial gains; notably, for Qwen2.5-32B on 2Wiki-MultiHopQA, our method obtains a +14.1 EM im-

provement over the strongest baseline. This trend extends to proprietary models: QuCo-RAG improves GPT-5-chat by +8.7 EM on 2WikiMultiHopQA and +5.5 EM on HotpotQA. Conversely, GPT models with agentic web search perform substantially worse than even the no-retrieval baseline, likely due to noisy web results not optimized for complex retrieval demands.

Why Proxy Corpus Works. The effectiveness of cross-model transfer validates our hypothesis that web-scale pre-training corpora share substantial overlap (Soldaini et al., 2024; Li et al., 2024a). Factual knowledge is largely drawn from common sources such as Common Crawl, Wikipedia, and curated web text, making frequency and co-occurrence statistics from one comprehensive corpus a reliable proxy for others. This property renders QuCo-RAG practically *model-agnostic*.

5.3 Efficiency Analysis (RQ3)

Figure 3 illustrates the efficiency-performance trade-off on HotpotQA. QuCo-RAG achieves the highest EM (35.0) while consuming only 87 tokens and 1.84 LLM calls on average, both the lowest among dynamic RAG methods. FS-RAG and DRAGIN consume 2–4× more tokens yet achieve substantially lower performance, while SeaKR incurs excessive LLM calls (10.28) due to repeated hidden-state uncertainty estimation. As shown in Figure 3(c), QuCo-RAG triggers only 1.70 retrievals per question on average, demonstrating precise corpus-grounded detection. Notably, no baseline falls in the green region (higher EM with fewer retrievals than QuCo-RAG), while methods like FLARE and FS-RAG fall in the red region, performing worse than Wo-RAG despite frequent retrieval. Regarding runtime, Figure 4 shows that LLM generation dominates (55–74%),

while corpus-based detection introduces modest overhead, demonstrating favorable scaling for deployment.

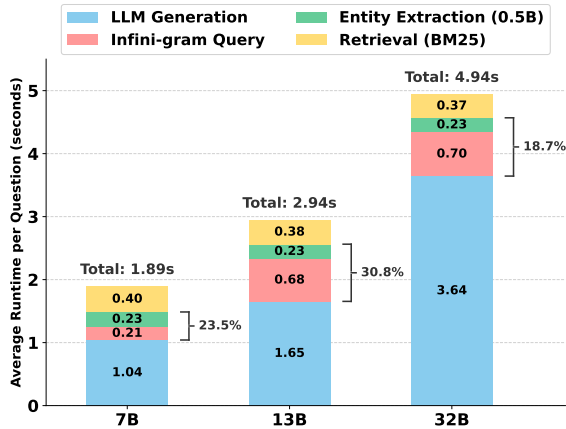


Figure 4: Average runtime breakdown per question for QuCo-RAG components across OLMo-2 model sizes on 2WikiMultihopQA.

6 Analysis and Discussion

We provide additional analyses including ablation studies, domain generalization, and performance breakdown by entity frequency. Threshold sensitivity analysis is provided in Appendix A.2.

6.1 Ablation Studies

Table 3 examines the contribution of each detection stage. Removing Pre-Generation Knowledge Assessment (w/o Initial Check) reduces EM by 2.5 points, confirming that identifying rare entities in the question is valuable for the initial response. Removing Runtime Claim Verification (w/o Runtime Check) causes a larger drop of 5.1 EM points, demonstrating that co-occurrence verification is the more critical component. Interestingly, even w/o Runtime Check (Initial Check only) outperforms SR-RAG by 3.9 EM while triggering fewer retrievals (0.76 vs. 1.00). This suggests that selective retrieval based on entity frequency can be more effective than always-retrieve strategies at the pre-generation stage—not all questions benefit equally from retrieval, and frequency-based detection provides a useful signal for prioritizing retrieval.

6.2 Domain Generalization

To evaluate generalization beyond open-domain QA, we test on PubMedQA (Jin et al., 2019), a biomedical QA benchmark where models answer research questions based on biomedical literature. Following Xiong et al. (2024), we use PubMed

Table 3: Ablation study on two-stage detection (2WikiMultihopQA, OLMo-2-7B). #Ret.: average retrieval count per question.

Configuration	EM	F1	#Ret.
QuCo-RAG (Full)	32.7	41.1	2.61
w/o Initial Check	30.2 _{-2.5}	38.0 _{-3.1}	1.82
w/o Runtime Check	27.6 _{-5.1}	35.6 _{-5.5}	0.76
<i>Baselines</i>			
SR-RAG	23.7	30.7	1.00
Wo-RAG	20.1	26.4	0.00

abstracts and medical textbooks (Jin et al., 2020) as the retrieval corpus \mathcal{C} and report accuracy following the standard benchmark setup (Wu et al., 2025). Notably, we retain the same OLMo-2 pre-training corpus as the statistical signal source \mathcal{P} , without any domain-specific adaptation.

As shown in Table 4, QuCo-RAG achieves the best accuracy (66.4%) while maintaining high efficiency (0.93 retrievals, 54.9 tokens per question). Internal-signal methods exhibit two failure modes in this specialized domain: *over-retrieval* and *under-retrieval*. FLARE suffers from over-retrieval, averaging 2.79 retrievals per question (significantly higher than its typical 1–2 in general-domain QA), achieving decent accuracy but at massive token cost. Conversely, DRAGIN and ETC suffer from under-retrieval, performing no better than Wo-RAG—likely because their internal-signal formulations fail to transfer across domains. QuCo-RAG avoids both pitfalls: large-scale pre-training corpora provide broad coverage of biomedical knowledge, and zero co-occurrence reliably indicates hallucination risks.

Table 4: Domain generalization on PubMedQA (OLMo-2-7B). Δ Acc: improvement over Wo-RAG; #Tok.: average token consumption per question.

Method	Acc	Δ Acc	#Ret.	#Tok.
Wo-RAG	55.2	0.0	0.00	40.3
FS-RAG	61.1	+5.9	5.74	436.1
FLARE	63.4	+8.2	2.79	516.8
DRAGIN	55.2	0.0	1.69	139.0
ETC	55.0	-0.2	0.25	58.8
QuCo-RAG	66.4	+11.2	0.93	54.9

6.3 Performance Across Entity Frequency

To understand how different methods handle knowledge of varying prevalence, we group questions by how often their entities appear in the pre-training corpus. Figure 5 shows EM scores and retrieval

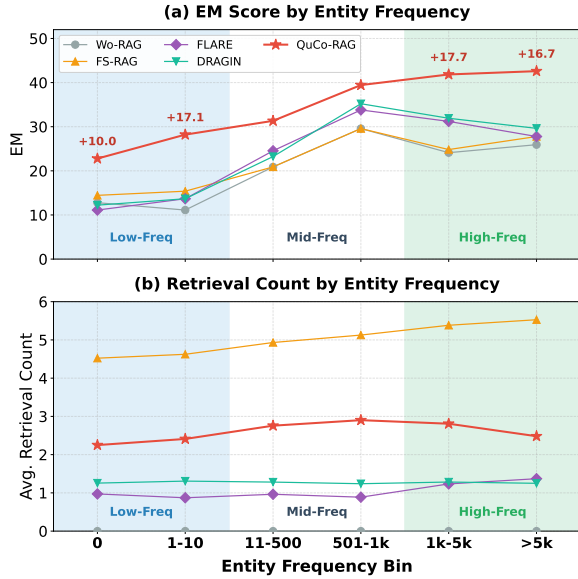


Figure 5: Performance stratified by entity frequency bins on 2WikiMultihopQA (OLMo-2-7B).

counts across frequency bins. Full numerical results are provided in Appendix Table 10. Overall, all methods perform worse in low-frequency bins, confirming that entity frequency correlates with model reliability. In **low-frequency bins (0–10)**, QuCo-RAG demonstrates dominant performance, outperforming Wo-RAG by 10–17 EM points, while DRAGIN and FLARE achieve nearly identical performance to Wo-RAG despite triggering retrievals, suggesting that models lack sufficient signal to recognize uncertainty on rare entities. In **mid-frequency bins (11–1k)**, the gap narrows as internal-signal methods become competitive, likely because mid-frequency entities place models in a “partially learned” state where entropy-based uncertainty is better calibrated. In **high-frequency bins (>1k)**, an interesting divergence emerges: baselines exhibit performance degradation while QuCo-RAG continues to improve. For internal-signal methods, the decline is likely due to overconfidence, failing to trigger retrieval even when generating wrong claims. In contrast, QuCo-RAG benefits from richer knowledge coverage: high-frequency entities have more thoroughly documented relationships in the corpus, making co-occurrence statistics more reliable for uncertainty quantification.

6.4 Broader Impact and Future Directions

Our work establishes corpus statistics as an objective alternative to model-internal uncertainty signals; while this paper focuses on retrieval triggering in RAG systems, the paradigm shift opens several promising avenues in AI safety and robustness.

Enabling Trustworthy AI Applications. Our experiments establish that corpus statistics offer a more reliable uncertainty measure than internal signals. This reliability is critical not only for RAG but also for broader safety-critical tasks, such as *selective answering*, where models can decline to answer when evidential support is absent, and *correctness prediction*, where corpus statistics provide well-grounded confidence scores for generated claims.

From Inference-Time Intervention to Data-Centric AI. Our corpus statistics analysis precisely identifies the model’s knowledge gaps. This signal can inform *training data curation*: rather than only compensating for gaps at inference time via retrieval, developers can proactively collect data for low-frequency entities during continued pre-training or post-training. Similarly, corpus statistics can guide *synthetic data filtering*, where LLM-generated training examples are verified against corpus statistics before inclusion, and *model editing* by distinguishing facts that require targeted injection from those already reliably learned.

Extensions of the Paradigm. Several directions merit exploration: (1) multilingual verification through cross-lingual statistics; (2) temporal dynamics via time-stamped corpora for evolving knowledge; (3) extension beyond entities to events, relations, and numerical claims; and (4) integration into agentic systems as a self-verification tool that agents invoke before acting on generation.

Theoretical Foundations. Our transferability results raise fundamental questions: why do proxy corpora work across model families? Can we formalize information-theoretic bounds on hallucination probability given corpus statistics? These questions connect to broader debates on memorization versus generalization in LLMs.

7 Conclusion

We propose QuCo-RAG, a dynamic RAG framework that quantifies uncertainty from pre-training corpus statistics rather than poorly calibrated model-internal signals. QuCo-RAG achieves state-of-the-art performance on multi-hop QA benchmarks while maintaining superior efficiency, transfers effectively to models with undisclosed training data (Llama, Qwen, GPT), and generalizes robustly to biomedical QA. These results establish corpus-grounded verification as a principled, practically model-agnostic paradigm for dynamic RAG.

Limitations

(1) Lexical Matching Constraints. Our co-occurrence verification relies on exact lexical matching of entity surface forms. This may lead to false positive retrieval triggers when two genuinely related entities co-occur in the corpus under alternative names or aliases (e.g., “NYC” vs. “New York City”), yet show zero co-occurrence for the specific surface forms extracted from the generated text. However, we argue this limitation is acceptable in practice due to the *asymmetric risk* inherent in RAG systems: the cost of an unnecessary retrieval (slightly increased latency) is far lower than that of an undetected hallucination (incorrect output). Our conservative strategy, triggering retrieval when in doubt, thus errs on the side of caution. Moreover, given the massive scale of the pre-training corpus, genuinely related entities typically co-occur in some form, mitigating alias-induced false alarms. Future work could incorporate entity linking (Xin et al., 2025) or canonicalization techniques (Hu et al., 2025) to further reduce unnecessary retrievals.

(2) Temporal Limitations of Static Corpora. Our approach inherits the temporal limitations of static pre-training corpora (Ding, 2025). A corpus indexed at a particular point in time cannot provide meaningful statistics for entities or events that emerge afterward (e.g., a 2024 corpus cannot verify claims about 2025 sports results or newly founded organizations). This limitation can be addressed through periodic corpus updates and index maintenance.

References

- Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, and 1 others. 2023. Gpt-4 technical report. *arXiv preprint arXiv:2303.08774*.
- Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. 2024. *Self-RAG: Learning to retrieve, generate, and critique through self-reflection*. In *The Twelfth International Conference on Learning Representations*.
- Nishant Balepur, Feng Gu, Abhilasha Ravichander, Shi Feng, Jordan Lee Boyd-Graber, and Rachel Rudinger. 2025. Reverse question answering: Can an llm write a question so hard (or bad) that it can’t answer? In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 2: Short Papers)*, pages 44–64.
- Huiyi Chen, Jiawei Peng, Kaihua Tang, Xin Geng, and Xu Yang. 2025. *Enhancing multimodal in-context learning for image classification through coreset optimization*. In *Proceedings of the 33rd ACM International Conference on Multimedia, MM ’25*, page 5130–5139, New York, NY, USA. Association for Computing Machinery.
- Zifeng Ding. 2025. *Inductive representation learning and natural language question answering on temporal knowledge graphs*. Ph.D. thesis, lmu.
- Guanting Dong, Hongyi Yuan, Keming Lu, Chengpeng Li, Mingfeng Xue, Dayiheng Liu, Wei Wang, Zheng Yuan, Chang Zhou, and Jingren Zhou. 2024. How abilities in large language models are affected by supervised fine-tuning data composition. In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 177–198.
- Alex Fang, Thomas Voice, Ruoming Pang, Ludwig Schmidt, and Tom Gunter. 2025. Reusing pre-training data at test time is a compute multiplier. *arXiv preprint arXiv:2511.04234*.
- Paolo Ferragina and Giovanni Manzini. 2000. Opportunistic data structures with applications. In *Proceedings 41st annual symposium on foundations of computer science*, pages 390–398. IEEE.
- Luis Galárraga, Jeremy Heitz, Kevin Murphy, and Fabian M Suchanek. 2014. Canonicalizing open knowledge bases. In *Proceedings of the 23rd acm international conference on conference on information and knowledge management*, pages 1679–1688.
- Tianyu Gao, Howard Yen, Jiatong Yu, and Danqi Chen. 2023a. *Enabling large language models to generate text with citations*. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 6465–6488, Singapore. Association for Computational Linguistics.
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yixin Dai, Jiawei Sun, Haofen Wang, and Haofen Wang. 2023b. Retrieval-augmented generation for large language models: A survey. *arXiv preprint arXiv:2312.10997*, 2(1).
- Aaron Grattafiori, Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Alex Vaughan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. 2017. On calibration of modern neural networks. In *International conference on machine learning*, pages 1321–1330. PMLR.

- Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shitong Ma, Peiyi Wang, Xiao Bi, and 1 others. 2025. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning. *arXiv preprint arXiv:2501.12948*.
- Matthew Ho, Chen Si, Zhaoxiang Feng, Fangxu Yu, Yichi Yang, Zhijian Liu, Zhiting Hu, and Lianhui Qin. 2025. Arcmemo: Abstract reasoning composition with lifelong llm memory. *arXiv preprint arXiv:2509.04439*.
- Xanh Ho, Anh-Khoa Duong Nguyen, Saku Sugawara, and Akiko Aizawa. 2020. [Constructing a multi-hop QA dataset for comprehensive evaluation of reasoning steps](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 6609–6625, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Yujia Hu, Tuan-Phong Nguyen, Shrestha Ghosh, and Simon Razniewski. 2025. Enabling llm knowledge analysis via extensive materialization. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 16189–16202.
- Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.
- Zhengbao Jiang, Frank Xu, Luyu Gao, Zhiqing Sun, Qian Liu, Jane Dwivedi-Yu, Yiming Yang, Jamie Callan, and Graham Neubig. 2023. [Active retrieval augmented generation](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 7969–7992, Singapore. Association for Computational Linguistics.
- Di Jin, Eileen Pan, Nassim Oufattole, Wei-Hung Weng, Hanyi Fang, and Peter Szolovits. 2020. What disease does this patient have? a large-scale open domain question answering dataset from medical exams. *arXiv preprint arXiv:2009.13081*.
- Qiao Jin, Bhuwan Dhingra, Zhengping Liu, William Cohen, and Xinghua Lu. 2019. Pubmedqa: A dataset for biomedical research question answering. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 2567–2577.
- Saurav Kadavath, Tom Conerly, Amanda Askell, Tom Henighan, Dawn Drain, Ethan Perez, Nicholas Schiefer, Zac Hatfield-Dodds, Nova DasSarma, Eli Tran-Johnson, and 1 others. 2022. Language models (mostly) know what they know. *arXiv preprint arXiv:2207.05221*.
- Adam Tauman Kalai and Santosh S Vempala. 2024. Calibrated language models must hallucinate. In *Proceedings of the 56th Annual ACM Symposium on Theory of Computing*, pages 160–171.
- Nikhil Kandpal, Haikang Deng, Adam Roberts, Eric Wallace, and Colin Raffel. 2023. Large language models struggle to learn long-tail knowledge. In *International conference on machine learning*, pages 15696–15707. PMLR.
- Vladimir Karpukhin, Barlas Oguz, Sewon Min, Patrick Lewis, Ledell Wu, Sergey Edunov, Danqi Chen, and Wen-tau Yih. 2020. [Dense passage retrieval for open-domain question answering](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 6769–6781, Online. Association for Computational Linguistics.
- Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen-tau Yih, Tim Rocktäschel, and 1 others. 2020. Retrieval-augmented generation for knowledge-intensive nlp tasks. *Advances in neural information processing systems*, 33:9459–9474.
- Bo Li, Tian Tian, Zhenghua Xu, Hao Cheng, Shikun Zhang, and Wei Ye. 2025a. Modeling uncertainty trends for timely retrieval in dynamic rag. *arXiv preprint arXiv:2511.09980*.
- Dawei Li, Bohan Jiang, Liangjie Huang, Alimohammad Beigi, Chengshuai Zhao, Zhen Tan, Amrita Bhat-tacharjee, Yuxuan Jiang, Canyu Chen, Tianhao Wu, and 1 others. 2025b. From generation to judgment: Opportunities and challenges of llm-as-a-judge. In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 2757–2791.
- Jeffrey Li, Alex Fang, Georgios Smyrnis, Maor Ivgi, Matt Jordan, Samir Yitzhak Gadre, Hritik Bansal, Etash Guha, Sedrick Scott Keh, Kushal Arora, and 1 others. 2024a. Datacomp-1m: In search of the next generation of training sets for language models. *Advances in Neural Information Processing Systems*, 37:14200–14282.
- Siheng Li, Cheng Yang, Taiqiang Wu, Chufan Shi, Yuji Zhang, Xinyu Zhu, Zesen Cheng, Deng Cai, Mo Yu, Lemao Liu, Jie Zhou, Yujiu Yang, Ngai Wong, Xixin Wu, and Wai Lam. 2024b. A survey on the honesty of large language models. *arXiv preprint arXiv:2409.18786*.
- Yu Li, Zhe Yang, Yi Huang, Xin Liu, and Guilin Qi. 2025c. C³TG: Conflict-aware, composite, and collaborative controlled text generation. *arXiv preprint arXiv:2511.09292*.
- Jiacheng Liu, Taylor Blanton, Yanai Elazar, Sewon Min, Yen-Sung Chen, Arnavi Chheda-Kothary, Huy Tran, Byron Bischoff, Eric Marsh, Michael Schmitz, and 1 others. 2025. Olmotrace: Tracing language model outputs back to trillions of training tokens. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 3: System Demonstrations)*, pages 178–188.

- Jiacheng Liu, Sewon Min, Luke Zettlemoyer, Yejin Choi, and Hannaneh Hajishirzi. 2024. [Infini-gram: Scaling unbounded n-gram language models to a trillion tokens](#). In *First Conference on Language Modeling*.
- Huan Ma, Jingdong Chen, Joey Tianyi Zhou, Guangyu Wang, and Changqing Zhang. 2025. Estimating llm uncertainty with evidence. *arXiv preprint arXiv:2502.00290*.
- Alex Mallen, Akari Asai, Victor Zhong, Rajarshi Das, Daniel Khashabi, and Hannaneh Hajishirzi. 2023. [When not to trust language models: Investigating effectiveness of parametric and non-parametric memories](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 9802–9822, Toronto, Canada. Association for Computational Linguistics.
- Dehai Min, Zhiyang Xu, Guilin Qi, Lifu Huang, and Chenyu You. 2025. [UniHGKR: Unified instruction-aware heterogeneous knowledge retrievers](#). In *Proceedings of the 2025 Conference of the Nations of the Americas Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 4577–4594, Albuquerque, New Mexico. Association for Computational Linguistics.
- Niklas Muennighoff. 2022. Sgpt: Gpt sentence embeddings for semantic search. *arXiv preprint arXiv:2202.08904*.
- Team OLMO, Pete Walsh, Luca Soldaini, Dirk Groeneveld, Kyle Lo, Shane Arora, Akshita Bhagia, Yuling Gu, Shengyi Huang, Matt Jordan, and 1 others. 2024. 2 olmo 2 furious. *arXiv preprint arXiv:2501.00656*.
- OpenAI. 2025a. GPT-4.1 Release Information. <https://openai.com/index/gpt-4-1/>.
- OpenAI. 2025b. GPT-5 Release Information. <https://openai.com/index/introducing-gpt-5/>.
- Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, and 1 others. 2022. Training language models to follow instructions with human feedback. *Advances in neural information processing systems*, 35:27730–27744.
- Abhilasha Ravichander, Shruti Ghela, David Wadden, and Yejin Choi. 2025. [HALoGEN: Fantastic LLM hallucinations and where to find them](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1402–1425, Vienna, Austria. Association for Computational Linguistics.
- Stephen Robertson, Hugo Zaragoza, and 1 others. 2009. The probabilistic relevance framework: Bm25 and beyond. *Foundations and Trends® in Information Retrieval*, 3(4):333–389.
- Weijia Shi, Sewon Min, Michihiro Yasunaga, Minjoon Seo, Richard James, Mike Lewis, Luke Zettlemoyer, and Wen-tau Yih. 2024. Replug: Retrieval-augmented black-box language models. In *Proceedings of the 2024 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (Volume 1: Long Papers)*, pages 8371–8384.
- Luca Soldaini, Rodney Kinney, Akshita Bhagia, Dustin Schwenk, David Atkinson, Russell Authur, Ben Bogin, Khyathi Chandu, Jennifer Dumas, Yanai Elazar, and 1 others. 2024. Dolma: An open corpus of three trillion tokens for language model pretraining research. In *Proceedings of the 62nd annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 15725–15788.
- Weihang Su, Qingyao Ai, Yueyue Wu, Anzhe Xie, Changyue Wang, Yixiao Ma, Haitao Li, Zhijing Wu, Yiqun Liu, and Min Zhang. 2025. [Pre-training for legal case retrieval based on inter-case distinctions](#). *ACM Trans. Inf. Syst.*, 43(5).
- Weihang Su, Yichen Tang, Qingyao Ai, Zhijing Wu, and Yiqun Liu. 2024. [DRAGIN: Dynamic retrieval augmented generation based on the real-time information needs of large language models](#). In *Proceedings of the 62nd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 12991–13013, Bangkok, Thailand. Association for Computational Linguistics.
- Qwen Team. 2024. [Qwen2.5: A party of foundation models](#).
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. [Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback](#). In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Harsh Trivedi, Niranjan Balasubramanian, Tushar Khot, and Ashish Sabharwal. 2023. Interleaving retrieval with chain-of-thought reasoning for knowledge-intensive multi-step questions. In *Proceedings of the 61st annual meeting of the association for computational linguistics (volume 1: long papers)*, pages 10014–10037.
- Keheng Wang, Feiyu Duan, Peiguang Li, Sirui Wang, and Xunliang Cai. 2025. LLMs know what they need: Leveraging a missing information guided framework to empower retrieval-augmented generation. In *Proceedings of the 31st International Conference on Computational Linguistics*, pages 2379–2400.
- Yile Wang, Peng Li, Maosong Sun, and Yang Liu. 2023. [Self-knowledge guided retrieval augmentation for large language models](#). In *Findings of the Association for Computational Linguistics: EMNLP*

- 2023, pages 10303–10315, Singapore. Association for Computational Linguistics.
- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, and 1 others. 2022. Chain-of-thought prompting elicits reasoning in large language models. *Advances in neural information processing systems*, 35:24824–24837.
- Junde Wu, Jiayuan Zhu, Yunli Qi, Jingkun Chen, Min Xu, Filippo Menolascina, Yueming Jin, and Vicente Grau. 2025. [Medical graph RAG: Evidence-based medical large language model via graph retrieval-augmented generation](#). In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 28443–28467, Vienna, Austria. Association for Computational Linguistics.
- Amy Xin, Yunjia Qi, Zijun Yao, Fangwei Zhu, Kaisheng Zeng, Bin Xu, Lei Hou, and Juanzi Li. 2025. Llm-ael: Large language models are good context augmenters for entity linking. In *Proceedings of the 34th ACM International Conference on Information and Knowledge Management*, pages 3550–3559.
- Guangzhi Xiong, Qiao Jin, Zhiyong Lu, and Aidong Zhang. 2024. [Benchmarking retrieval-augmented generation for medicine](#). In *Findings of the Association for Computational Linguistics ACL 2024*, pages 6233–6251, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Hao Xu, Jiacheng Liu, Yejin Choi, Noah A. Smith, and Hannaneh Hajishirzi. 2025. [Infini-gram mini: Exact n-gram search at the Internet scale with FM-index](#). In *Proceedings of the 2025 Conference on Empirical Methods in Natural Language Processing*, pages 24955–24980, Suzhou, China. Association for Computational Linguistics.
- Zhipeng Xu, Zhenghao Liu, Yibin Liu, Chenyan Xiong, Yukun Yan, Shuo Wang, Shi Yu, Zhiyuan Liu, and Ge Yu. 2024. Activerag: Revealing the treasures of knowledge via active learning. *CoRR*.
- Diji Yang, Linda Zeng, Jinmeng Rao, and Yi Zhang. 2025. Knowing you don’t know: Learning when to continue search in multi-round rag through self-practicing. In *Proceedings of the 48th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 1305–1315.
- Zhilin Yang, Peng Qi, Saizheng Zhang, Yoshua Bengio, William W. Cohen, Ruslan Salakhutdinov, and Christopher D. Manning. 2018. HotpotQA: A dataset for diverse, explainable multi-hop question answering. In *Conference on Empirical Methods in Natural Language Processing (EMNLP)*.
- Zijun Yao, Weijian Qi, Liangming Pan, Shulin Cao, Linmei Hu, Liu Weichuan, Lei Hou, and Juanzi Li. 2025. Seakr: Self-aware knowledge retrieval for adaptive retrieval augmented generation. In *Proceedings of the 63rd Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 27022–27043.
- Fangxu Yu, Hongyu Zhao, and Tianyi Zhou. 2025. Ts-reasoner: Aligning time series foundation models with llm reasoning. *arXiv preprint arXiv:2510.03519*.
- Tian Yu, Shaolei Zhang, and Yang Feng. 2024. Auto-rag: Autonomous retrieval-augmented generation for large language models. *arXiv preprint arXiv:2411.19443*.
- Yanzhao Zhang, Mingxin Li, Dingkun Long, Xin Zhang, Huan Lin, Baosong Yang, Pengjun Xie, An Yang, Dayiheng Liu, Junyang Lin, and 1 others. 2025. Qwen3 embedding: Advancing text embedding and reranking through foundation models. *arXiv preprint arXiv:2506.05176*.

A Appendix

A.1 Additional Implementation Details

Generation Settings and Prompts. In our experiments, all open-source models use greedy decoding with a 128-token generation limit per step, and GPT models use default parameters via API calls. For generation, we employ 6-to-8-shot Chain-of-Thought prompting (Wei et al., 2022), adopting templates from Trivedi et al. (2023) and Jiang et al. (2023). We use 6 few-shot examples for 2WikiMultihopQA and 8 for HotpotQA, consistent with prior work. The full prompt template is provided in Table 5. We use the Wikipedia dump from Karpukhin et al. (2020) as our external corpus \mathcal{C} , which contains approximately 21 million passages.

Few-shot Examples:

Question: When did the director of film Hypocrite (Film) die?

Answer: The film Hypocrite was directed by Miguel Morayta. Miguel Morayta died on 19 June 2013. So the answer is 19 June 2013.

[... 5–7 more demonstrations ...]

Retrieved Context (if available):

Background information that may be potentially useful in addressing your question:

- [1] <retrieved document 1>
- [2] <retrieved document 2>
- [3] <retrieved document 3>

Instruction:

Please answer the following questions. The format of the answers should be the same as the examples given before. Specifically, you need to think through the answer to this question step by step. Each sentence should only present a fact statement. Avoid using pronouns like He/She/It or possessive pronouns like His/Her/Its, but instead use specific names. At the end of your answer, use “So the answer is” to provide your answer.

Question: <input question>

Table 5: Prompt template used for multi-hop QA experiments. Retrieved context is prepended when retrieval is triggered.

A.2 Threshold Sensitivity

We examine the robustness of QuCo-RAG to its two key hyperparameters: the entity frequency threshold τ_{entity} and the co-occurrence threshold τ_{cooc} . As illustrated in Figure 6(a), EM remains stable (32.2–32.7) across a wide range of τ_{entity} from

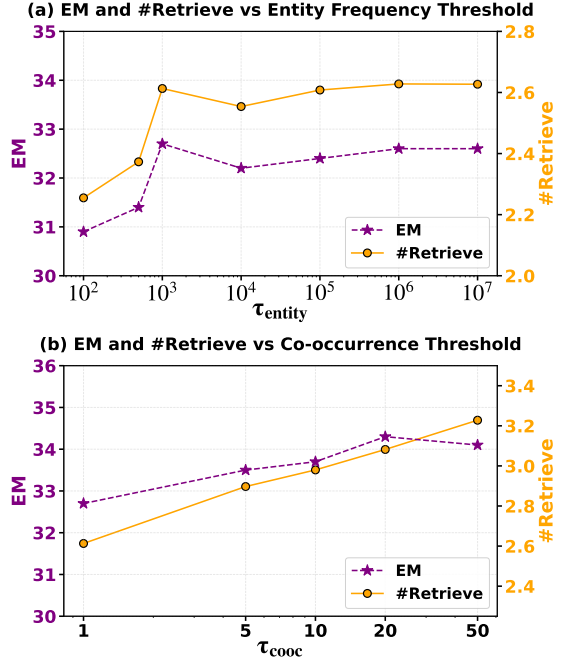


Figure 6: Threshold sensitivity analysis on 2WikiMultihopQA with OLMo-2-7B.

10^3 to 10^7 , with retrieval count also staying consistent (2.5–2.6), demonstrating strong robustness to this hyperparameter. For τ_{cooc} , as shown in Figure 6(b), increasing the threshold imposes a stricter verification standard (requiring more evidential support in the corpus), leading to a monotonic increase in retrieval frequency (from 2.61 to 3.23). While higher thresholds (e.g., $\tau_{\text{cooc}} = 20$) yield marginal EM improvements (reaching 34.3 EM), they incur significantly higher retrieval overhead. We adopt $\tau_{\text{cooc}} = 1$ (i.e., triggering on zero co-occurrence) as our default for its clear interpretability: if two entities never co-occur in the pre-training corpus, the generated claim lacks evidential support and is likely hallucinated.

A.3 Triplet Extractor Training Examples

The quality and diversity of training data are particularly important for robust model training (Li et al., 2025c; Yu et al., 2025). Table 6 shows representative examples from our triplet extractor training data. Each example consists of an input sentence and the extracted output. If the input sentence contains meaningful factual knowledge, the output consists of knowledge triplets in the format (head entity, relation, tail entity); otherwise, the output is empty. We prioritize extracting triplets where the tail entity is a named entity (person, location, organization, date) rather than generic descriptors, as

Table 6: Examples of triplet extractor training data. The model extracts factual triplets from declarative sentences, partial triplets from questions (since the answer is unknown), and returns empty for non-factual statements.

Input Sentence	Extracted Output
<i>Declarative sentences with factual knowledge:</i>	
Kumbasaram was released in 2017.	[["Kumbasaram", "released in", "2017"]]
Beowulf & Grendel was directed by Sturla Gunnarsson.	[["Beowulf & Grendel", "directed by", "Sturla Gunnarsson"]]
Coulson Wallop’s father, Nigel Wallop, studied at Eton College.	[["Coulson Wallop", "father", "Nigel Wallop"], ["Nigel Wallop", "studied at", "Eton College"]]
<i>Questions (answer unknown, extract partial triplets):</i>	
Which film came out first, Kumbasaram or Mystery Of The 13th Guest?	[["Kumbasaram", "came out"], ["Mystery of the 13th Guest", "came out"]]
Where did Diane Meyer Simon’s husband graduate from?	[["Diane Meyer Simon", "husband"]]
<i>Non-factual statements (reasoning conclusions):</i>	
Thus, Kumbasaram came out first.	[]
Therefore, Robert Enrico, the director of The Woman Thou Gavest Me, was born first.	[]

these are more amenable to corpus co-occurrence verification. Non-factual statements such as reasoning conclusions (e.g., sentences starting with "Thus" or "Therefore") return empty outputs since they do not introduce new verifiable facts.

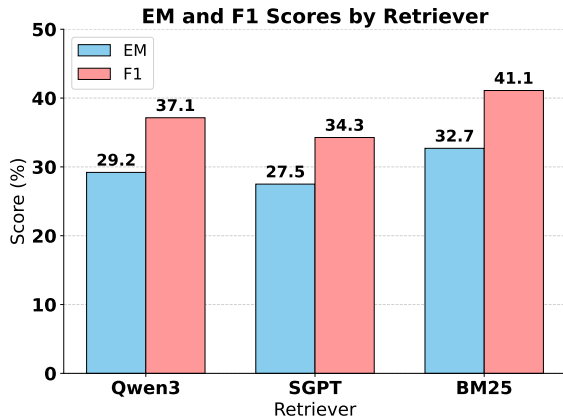


Figure 7: Performance comparison of QuCo-RAG with different retrievers (Qwen3-Embedding, SGPT, and BM25) on 2WikiMultihopQA using OLMo-2-7B.

A.4 Effect of Different Retrievers

To verify that QuCo-RAG is robust to retriever choice, we compare BM25 with dense retrievers SGPT (Muennighoff, 2022) and Qwen3-Embedding-0.6B (Zhang et al., 2025). As shown in Figure 7, QuCo-RAG achieves robust performance across all three retrievers, with EM scores ranging from 27.5 to 32.7 and F1 from 34.3 to 41.1. BM25 achieves the best results (32.7 EM, 41.1 F1), align-

ing with prior findings that sparse retrieval remains highly competitive for RAG tasks (Su et al., 2024). Importantly, even with different retriever backends, QuCo-RAG consistently outperforms baselines (cf. Table 1), confirming that our corpus-based uncertainty quantification mechanism is orthogonal to the choice of retrieval system.

A.5 Case Study

Table 9 presents a detailed case study demonstrating how QuCo-RAG quantifies uncertainty through corpus statistics to detect and correct hallucinations that baseline methods miss. In this multi-hop question, all baselines fail for distinct reasons: Wo-RAG hallucinates without any correction mechanism; SR-RAG retrieves correct director information but cannot perform follow-up retrieval for the mother; FLARE and DRAGIN both detect some uncertainty but their queries contain the hallucinated director name “Igor Maslennikov,” leading to retrieval of irrelevant documents that reinforce the error. Notably, DRAGIN’s internal signals mark this completely fabricated director as low uncertainty, exemplifying the confident hallucination problem. In contrast, QuCo-RAG succeeds through the coordination of two stages: Stage 1 identifies “Polish-Russian War” as a low-frequency entity, triggering initial retrieval that grounds the model to generate the correct director “Xawery Żuławski.” Stage 2 then catches the hallucinated mother “Anna Żuławski” via zero co-occurrence,

Table 7: Comparison of different RAG methods on 2WikiMultihopQA and HotpotQA benchmarks.

Method	2Wiki		HotpotQA	
	EM	F1	EM	F1
Qwen2.5-32B-Instruct				
Wo-RAG	26.4	33.6	21.6	32.4
SR-RAG	23.0	31.8	31.0	41.7
FS-RAG	<u>35.9</u>	<u>45.3</u>	<u>38.6</u>	<u>49.6</u>
FLARE	26.4	33.3	24.1	33.5
DRAGIN	28.8	36.9	22.2	32.4
ETC	31.5	40.2	21.7	32.0
SeaKR	22.4	31.3	26.7	37.5
QuCo-RAG	50.0	58.9	41.6	55.1
Llama-3-8B-Instruct				
Wo-RAG	29.5	37.7	20.3	31.4
SR-RAG	12.9	29.2	22.7	35.4
FS-RAG	28.8	36.8	27.0	38.5
FLARE	26.6	35.1	22.2	31.5
DRAGIN	27.9	36.7	20.0	31.9
ETC	29.9	39.2	24.1	35.1
SeaKR	<u>33.5</u>	<u>40.4</u>	<u>33.5</u>	<u>46.0</u>
QuCo-RAG	38.4	46.6	36.2	48.7
GPT-4.1				
Wo-RAG	54.7	69.9	<u>40.1</u>	<u>56.1</u>
SR-RAG	<u>60.0</u>	72.6	38.8	54.2
FS-RAG	59.5	<u>73.8</u>	25.9	36.5
FLARE	49.8	67.9	38.7	52.1
Web-Tool	42.9	63.2	8.9	16.8
QuCo-RAG	64.6	74.8	48.2	62.2
GPT-5-chat				
Wo-RAG	50.1	67.0	37.7	54.5
SR-RAG	<u>51.0</u>	<u>70.1</u>	<u>42.9</u>	<u>58.6</u>
FS-RAG	47.3	63.3	19.0	31.3
Web-Tool	48.3	69.8	19.8	33.6
QuCo-RAG	59.7	73.3	48.4	62.6

triggering targeted retrieval with a hallucination-free query “Xawery Żułowski mother” that yields the correct answer.

A.6 Full Results for Transferability Experiments

Transferability across different models is crucial for practical deployment (Ho et al., 2025; Chen et al., 2025). Table 7 presents the complete results (EM and F1) for the transferability experiments discussed in Section 5.2. The main paper reports only EM scores for brevity. Across all model families (Qwen2.5-32B, Llama-3-8B, GPT-4.1, and GPT-5-chat), QuCo-RAG consistently achieves the best performance on both metrics. The F1 improvements follow similar patterns to EM, confirming

Table 8: Efficiency comparison of RAG methods across OLMo-2 model sizes. #Tok.: average number of tokens used; #Call: average number of LLM calls; #Ret.: average number of retrieval operations.

Method	2WikiMultihopQA			HotpotQA		
	#Tok.	#Call	#Ret.	#Tok.	#Call	#Ret.
OLMo-2-7B						
Wo-RAG	58.62	1.00	0.00	54.15	1.00	0.00
SR-RAG	49.23	1.00	1.00	69.04	1.00	1.00
FS-RAG	306.09	4.96	4.96	417.77	6.91	6.91
FLARE	132.90	2.33	1.03	436.37	6.89	3.39
DRAGIN	114.09	2.58	1.27	387.54	6.52	3.24
ETC	124.48	3.25	1.25	83.69	2.38	0.79
SeaKR	99.89	11.91	1.39	100.22	10.95	1.29
QuCo-RAG	107.87	2.44	2.61	128.20	3.23	4.47
OLMo-2-13B						
Wo-RAG	53.63	1.00	0.00	54.59	1.00	0.00
SR-RAG	70.65	1.00	1.00	69.57	1.00	1.00
FS-RAG	234.42	4.36	4.36	464.35	6.48	6.48
FLARE	129.67	2.01	0.93	284.34	3.42	1.69
DRAGIN	134.78	2.78	1.27	254.14	4.26	1.96
ETC	126.00	3.23	1.22	100.26	2.56	0.85
SeaKR	78.42	9.42	1.01	92.11	10.28	1.29
QuCo-RAG	105.83	2.50	2.50	87.19	1.84	1.70
OLMo-2-32B						
Wo-RAG	54.72	1.00	0.00	76.19	1.00	0.00
SR-RAG	64.61	1.00	1.00	91.31	1.00	1.00
FS-RAG	266.70	5.02	5.02	593.71	8.59	8.59
FLARE	116.19	2.10	1.01	270.10	3.20	1.59
DRAGIN	103.53	2.69	1.26	554.09	7.49	3.71
ETC	116.85	3.15	1.19	106.24	2.61	0.91
SeaKR	91.08	14.26	2.46	79.43	12.72	1.97
QuCo-RAG	116.29	2.43	2.49	98.09	1.90	1.99

that QuCo-RAG’s gains are robust.

A.7 Detailed Efficiency Metrics

Table 8 presents the complete efficiency comparison across all OLMo-2 model sizes on both datasets. We report three metrics: average token consumption (#Tok.), LLM calls (#Call), and retrieval operations (#Ret.) per question. QuCo-RAG maintains competitive efficiency across all settings. Notably, on HotpotQA with OLMo-2-32B, QuCo-RAG achieves the highest EM (41.6, see Table 1) while using only 98 tokens and 1.90 LLM calls, compared to FS-RAG which consumes 594 tokens and 8.59 calls yet achieves only 13.9 EM. SeaKR consistently incurs the highest number of LLM calls (9–14 per question) due to its iterative hidden-state uncertainty estimation.

A.8 Detailed Performance Breakdown by Entity Frequency Bin

Table 10 presents the full performance breakdown by entity frequency. Entity frequency is defined

Table 9: Case study comparison. **Red** indicates hallucinated/incorrect content; **green** indicates correct content. Only QuCo-RAG produces the correct answer through corpus-grounded uncertainty quantification.

Question: <i>Who is the mother of the director of film Polish-Russian War?</i> Ground Truth: Małgorzata Braunek (Polish-Russian War (film) → Director: Xawery Żuławski → Mother: Małgorzata Braunek)					
Method	Initial Generation	Uncertainty Signal	Retrieval Query	Final Answer	Analysis
Wo-RAG	"...directed by Igor Maslennikov . His mother is Natalia Maslennikova ."	N/A	N/A	Natalia Maslennikova	No retrieval mechanism to correct hallucinated director.
SR-RAG	"...directed by Xawery Żuławski . No information about his mother."	N/A (retrieves once before generation)	Original question	unknown	Single-round retrieval insufficient for multi-hop reasoning.
FLARE	"...directed by Igor Maslennikov . His mother is Svetlana ."	Triggered at sentence-level (probability below threshold)	"Igor Maslennikov..."	unknown	Query included hallucinated director; retrieved irrelevant documents.
DRAGIN	"...directed by Igor Maslennikov . His mother is Natalia Maslennikova ."	Triggered at token "Natalia" (entropy-based); wrong director marked as low uncertainty	"Igor Maslennikov mother"	Natalia Maslennikova	Confident hallucination: internal signals failed to flag the wrong director; query contained error, reinforcing mistake.
QuCo-RAG	S1: "...directed by Xawery Żuławski ." S2: "...mother is Anna Żuławski ."	Stage 1: Low entity freq. → retrieval Stage 2: Co-occurrence = 0 → high uncertainty	Stage 1: Original question Stage 2: "Xawery Żuławski mother"	Małgorzata Braunek	Stage 1 ensured correct director via initial retrieval; Stage 2 caught hallucinated mother via zero co-occurrence.

Table 10: Detailed performance breakdown by entity frequency on 2WikiMultihopQA (OLMo-2-7B). Entity frequency is defined as the average appearance count of all entities in the question within the OLMo-2 pre-training corpus.

Freq. Bin	Count	Wo-RAG		SR-RAG		FS-RAG		FLARE		DRAGIN		QuCo-RAG	
		EM	#Ret.	EM	#Ret.	EM	#Ret.	EM	#Ret.	EM	#Ret.	EM	#Ret.
0	180	12.8	0.00	13.9	1.00	14.4	4.52	11.1	0.97	12.2	1.26	22.8	2.25
1-10	117	11.1	0.00	20.5	1.00	15.4	4.62	13.7	0.87	13.7	1.31	28.2	2.41
11-50	119	13.4	0.00	25.2	1.00	18.5	4.79	17.6	0.84	15.1	1.32	26.9	2.67
51-100	66	27.3	0.00	18.2	1.00	16.7	5.15	25.8	1.17	36.4	1.18	34.8	2.91
101-500	198	23.2	0.00	21.2	1.00	23.7	4.94	28.3	0.97	23.7	1.29	32.8	2.76
501-1k	71	29.6	0.00	40.8	1.00	29.6	5.13	33.8	0.89	35.2	1.24	39.4	2.90
1k-5k	141	24.1	0.00	29.1	1.00	24.8	5.38	31.2	1.23	31.9	1.28	41.8	2.81
>5k	108	25.9	0.00	29.6	1.00	27.8	5.53	27.8	1.37	29.6	1.25	42.6	2.48
Overall	1000	19.9	0.00	23.5	1.00	21.0	4.96	22.8	1.03	22.9	1.27	32.7	2.61

as the average occurrence count of all entities in the question within the OLMo-2 pre-training corpus. QuCo-RAG achieves the best EM in 6 out of 8 frequency bins, with particularly large gains on low-frequency entities (frequency < 50) where internal-signal-based methods (FLARE, DRAGIN) perform similarly to Wo-RAG. This validates our core hypothesis that entity frequency in the pre-training corpus serves as an effective indicator of knowledge gaps.